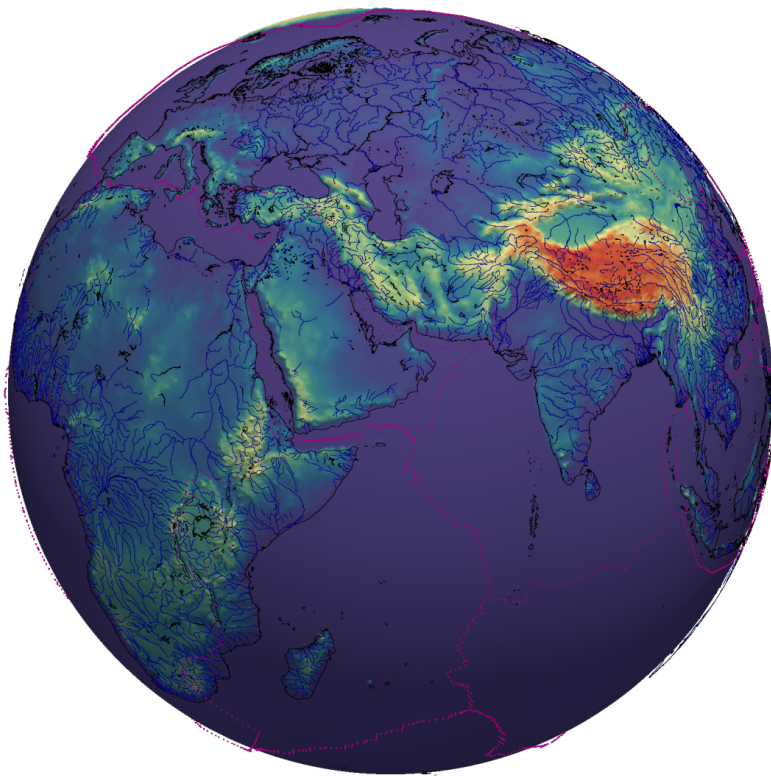


FieldStone



C. Thieulot

With (direct or indirect, small or large) contributions from (in alphabetical order): Wolfgang Bangerth, Daniëlle Bintanja, Marjolein Blasweiler, Taco Broerse, Daniel Douglas, Rens Elbertsen, Zoltán Erdős, Frederic Gueydan, Riad Hassani, Sverre Hassing, Agnes Hendrickx, Jort Jansen, Wouter Klessens, Gilles Mercier, Job Mos, Bob Myhill, Taka Shinohara, Alessandro Regorda, Neil Ribe, Ashim Rijal, Bart Root, Thomas Sanders, Wim Spakman, Thomas Theunissen, Marcel Thielmann, Arie van den Berg, Michiel van den Berg, Erik van der Wiel, Lukas van de Wiel, Eric van den Hoogen, Iris van Zelst, Jan Veenhof, Alraune Zech.

If you find anything in this document useful for your research please cite it as follows (*please* include the doi!):

```
@article{fieldstone,  
author = "Cedric Thieulot",  
title = "{Fieldstone: a computational geodynamics (self-)teaching tool}",  
year = "2023",  
doi = "10.5194/egusphere-egu23-14212"}
```

Why do I have to promise where I am going while I am not there yet?

You can't google something you don't know exists.

You can be correct or you can get stuff done.

Contents

| | | |
|----------|-------------------------------------------------------------------------------------------------------------|-----------|
| 1 | Introduction | 17 |
| 1.1 | Philosophy | 17 |
| 1.2 | ambition & motivation | 18 |
| 1.3 | Acknowledgements | 18 |
| 1.4 | About the author | 18 |
| 1.5 | Essential/relevant literature | 19 |
| 1.6 | Installing packages | 19 |
| 1.7 | What is a (real) fieldstone? | 20 |
| 1.8 | Why the Finite Element method? | 21 |
| 1.9 | Notations | 21 |
| 1.10 | Colour maps for visualisation | 21 |
| 1.11 | How my bibliography works | 22 |
| 1.12 | Youtube resources | 23 |
| 1.13 | How to download a single stone | 24 |
| 1.14 | Oldies but goodies | 25 |
| 2 | Physics and a bit of mathematics | 29 |
| 2.1 | Some maths | 29 |
| 2.1.1 | About vectors | 29 |
| 2.1.2 | dot products, cross products and dyadic products | 30 |
| 2.1.3 | Rotation matrix | 30 |
| 2.2 | Units | 31 |
| 2.3 | Coordinate systems | 33 |
| 2.3.1 | Cartesian coordinates | 33 |
| 2.3.2 | Polar coordinates | 34 |
| 2.3.3 | Cylindrical coordinates | 35 |
| 2.3.4 | Spherical coordinates | 36 |
| 2.3.5 | Converting tensors between Cartesian and Cylindrical bases | 38 |
| 2.3.6 | Converting tensors between Cartesian and Spherical bases | 39 |
| 2.4 | A continuum mechanics primer | 39 |
| 2.4.1 | Forces | 39 |
| 2.4.2 | Stress tensor and tractions | 40 |
| 2.4.3 | Strain rate and spin tensor | 41 |
| 2.5 | Viscous Newtonian rheology | 42 |
| 2.6 | The heat transport equation - energy conservation equation | 43 |
| 2.7 | The momentum conservation equations | 43 |
| 2.8 | The mass conservation equations | 43 |
| 2.9 | The equations in ASPECT manual | 44 |
| 2.10 | Equations for thermal convection in an anelastic, compressible, self-gravitating spherical mantle | 46 |

| | | |
|---------|------------------------------------------------------------------|-----|
| 2.11 | Non-dimensionalisation of the Navier-Stokes equations | 48 |
| 2.11.1 | Approach # 1 - isothermal flow | 49 |
| 2.11.2 | Approach # 2 - Temperature dependent | 50 |
| 2.12 | The Navier-Stokes equations in cylindrical coordinates | 54 |
| 2.13 | The Stokes equations in spherical coordinates | 55 |
| 2.14 | The equations for axisymmetric geometries | 57 |
| 2.15 | The Boussinesq approximation | 60 |
| 2.16 | The Extended Boussinesq approximation | 60 |
| 2.17 | Stokes equation for elastic medium | 61 |
| 2.18 | The strain rate tensor in all coordinate systems | 63 |
| 2.18.1 | Cartesian coordinates | 63 |
| 2.18.2 | Polar coordinates | 64 |
| 2.18.3 | Cylindrical coordinates | 64 |
| 2.18.4 | Spherical coordinates | 64 |
| 2.18.5 | Relationship between Cartesian and polar coordinates expressions | 64 |
| 2.19 | Boundary conditions | 66 |
| 2.19.1 | The Stokes equations | 66 |
| 2.19.2 | The heat transport equation | 66 |
| 2.20 | Meaningful physical quantities | 67 |
| 2.21 | Principal stress and principal invariants | 71 |
| 2.21.1 | In two dimensions | 71 |
| 2.21.2 | In three dimensions | 73 |
| 2.21.3 | About the 2nd principal invariant of the deviatoric stress | 77 |
| 2.22 | Tensor (moment) invariants | 78 |
| 2.23 | Stress & strain rate invariants | 79 |
| 2.24 | Two-dimensional plane strain calculations | 81 |
| 2.25 | Alternative principal stresses notations | 85 |
| 2.26 | Recap of notations and definitions of stress invariants | 88 |
| 2.27 | Rheology in geodynamics | 90 |
| 2.27.1 | Linear viscous aka Newtonian | 90 |
| 2.27.2 | Power-law model | 91 |
| 2.27.3 | Carreau model | 91 |
| 2.27.4 | Bingham model | 92 |
| 2.27.5 | Herschel-Bulkley visco-plastic model | 93 |
| 2.27.6 | The Casson model | 94 |
| 2.27.7 | The Ellis model | 94 |
| 2.27.8 | One model to rule them all? | 94 |
| 2.27.9 | Dislocation and Diffusion creep | 95 |
| 2.27.10 | The von Mises failure criterion | 99 |
| 2.27.11 | The Tresca failure criterion | 102 |
| 2.27.12 | The Mohr-Coulomb failure criterion | 104 |
| 2.27.13 | The Drucker-Prager failure criterion | 107 |
| 2.27.14 | The Griffith-Murrell failure criterion | 114 |
| 2.27.15 | The Cam-clay failure criterion | 114 |
| 2.27.16 | The failure envelope, or yield surface | 114 |
| 2.27.17 | Peierls creep | 117 |
| 2.27.18 | Stress limiting rheology | 117 |
| 2.27.19 | Arrhenius law | 118 |
| 2.27.20 | Simple parametrisation of the mantle | 118 |
| 2.27.21 | Glen's law for ice | 119 |

| | | |
|----------|---------------------------------------------------------------------------------------------------|------------|
| 2.27.22 | Strain rate partitioning across deformation mechanisms | 120 |
| 2.27.23 | Anisotropic viscosity | 131 |
| 2.27.24 | Rheology of the lithosphere | 133 |
| 2.28 | The Perzyna model | 135 |
| 2.28.1 | von Mises plasticity following Zienkiewicz (1975) | 138 |
| 2.28.2 | Dissecting Choi & Petersen (2015) | 141 |
| 2.28.3 | my take on this in 3D for Drucker-Prager | 143 |
| 2.28.4 | my take on this in 3D for MC | 144 |
| 2.28.5 | Revisiting Lemiale et al (2008) and Spiegelman et al (2016) | 145 |
| 2.29 | Moment of inertia | 147 |
| 2.30 | The need for numerical modelling | 148 |
| 2.31 | Important mathematical concepts and equations | 149 |
| 2.31.1 | Taylor expansion | 149 |
| 2.31.2 | Divergence theorem | 149 |
| 3 | The Finite Difference Method | 150 |
| 3.1 | Back to basics: what is a derivative? | 150 |
| 3.2 | Welcome to the discrete world | 151 |
| 3.3 | FDM basics in 1D | 151 |
| 3.4 | Solving the 1D diffusion equation | 157 |
| 3.5 | Solving the 1D advection equation | 165 |
| 3.6 | FDM basics in 2D | 168 |
| 3.7 | Solving the 2D diffusion equation | 170 |
| 3.7.1 | Explicit scheme | 171 |
| 3.7.2 | Implicit scheme | 172 |
| 3.7.3 | The 9-point stencil for the Laplace operator | 177 |
| 3.8 | Solving the 2D advection-diffusion equation | 177 |
| 3.9 | FEM vs FDM? | 180 |
| 4 | Numerical integration | 183 |
| 4.1 | In 1 dimension | 184 |
| 4.1.1 | Midpoint and Trepezoidal rules | 184 |
| 4.1.2 | in 1D - Gauss-Legendre quadrature | 186 |
| 4.1.3 | A probably naive way of finding the quadrature points coordinates and weights | 188 |
| 4.1.4 | Examples | 189 |
| 4.2 | In 2 & 3 dimensions | 191 |
| 4.2.1 | On the reference square | 191 |
| 4.2.2 | On a generic quadrilateral | 191 |
| 4.2.3 | Exercises | 192 |
| 4.2.4 | Quadrature on triangles | 194 |
| 4.2.5 | A mathematical recreation: computing the volume of a tetrahedron | 198 |
| 4.2.6 | Quadrature on tetrahedra | 199 |
| 4.2.7 | The Gauss-Lobatto approach | 200 |
| 4.2.8 | Computing the 'real' coordinates of the quadrature points and other consider- ations | 201 |
| 5 | The building blocks of the Finite Element Method | 203 |
| 5.1 | A bit of FE terminology | 203 |
| 5.2 | Elements and basis functions in 1D | 205 |
| 5.2.1 | Linear basis functions (Q_1) | 205 |
| 5.2.2 | Quadratic basis functions (Q_2) | 206 |

| | | |
|----------|------------------------------------------------------------------------|------------|
| 5.2.3 | Cubic basis functions (Q_3) | 207 |
| 5.2.4 | Quartic basis functions (Q_4) | 210 |
| 5.2.5 | Fifth-order basis functions (Q_5) | 212 |
| 5.2.6 | Sixth-order basis functions (Q_6) | 213 |
| 5.2.7 | A generic approach to 1D basis functions | 217 |
| 5.3 | Elements and basis functions in 2D | 219 |
| 5.3.1 | Bilinear basis functions in 2D (Q_1) | 220 |
| 5.3.2 | Biquadratic basis functions in 2D (Q_2) | 222 |
| 5.3.3 | Bicubic basis functions in 2D (Q_3) | 224 |
| 5.3.4 | Eight node serendipity basis functions in 2D ($Q_2^{(8)}$) | 226 |
| 5.3.5 | Eight node serendipity basis functions in 2D ($QH8 - C1$) | 228 |
| 5.3.6 | Biquartic basis functions in 2D (Q_4) | 229 |
| 5.3.7 | Linear basis functions for triangles in 2D (P_1) | 229 |
| 5.3.8 | Linear basis functions for quadrilaterals in 2D (P_1) | 232 |
| 5.3.9 | Enriched linear basis functions in triangles (P_1^+) | 233 |
| 5.3.10 | Quadratic basis functions for triangles in 2D (P_2) | 235 |
| 5.3.11 | Enriched quadratic basis functions in triangles (P_2^+) | 237 |
| 5.3.12 | Cubic basis functions for triangles (P_3) | 239 |
| 5.3.13 | Quartic basis functions for triangles (P_4) | 243 |
| 5.3.14 | Enriched linear basis functions in quadrilaterals (Q_1^+) -WIP | 248 |
| 5.3.15 | The rotated Q_1 (Rannacher-Turek element) | 254 |
| 5.3.16 | The 2D enriched $Q_1^+ \times P_0$ of Fortin | 257 |
| 5.3.17 | The P_1^{NC} space | 258 |
| 5.4 | Elements and basis functions in 3D | 260 |
| 5.4.1 | Linear basis functions in tetrahedra (P_1) | 260 |
| 5.4.2 | Enriched linear in tetrahedra (P_1^+) | 261 |
| 5.4.3 | Triquadratic basis functions in 3D (Q_2) | 262 |
| 5.4.4 | Enriched quadratic basis functions in tetrahedra (P_2^+) | 263 |
| 5.4.5 | Linear basis functions for hexahedra (P_1) | 264 |
| 5.4.6 | 20-node serendipity basis functions in 3D ($Q_2^{(20)}$) | 266 |
| 5.4.7 | The rotated Q_1 | 266 |
| 5.4.8 | The 3D enriched $Q_1^+ \times P_0$ of Fortin | 269 |
| 5.4.9 | The $Q_1^{++} \times Q_1$ of Karabelas et al (2020) | 276 |
| 5.4.10 | The DSSY element | 283 |
| 5.5 | Low order elements recap | 286 |
| 5.6 | On the meaning of basis functions | 289 |
| 6 | Solving the heat transport equation with linear Finite Elements | 292 |
| 6.1 | The diffusion equation in 1D | 292 |
| 6.2 | The advection-diffusion equation in 1D | 301 |
| 6.3 | The advection-diffusion equation in 2D | 305 |
| 6.4 | Another approach to solving the advection diffusion | 311 |
| 6.5 | The advection-diffusion eq in axisymmetric cylindrical coordinates | 313 |
| 6.6 | The SUPG formulation for the energy equation | 314 |
| 7 | Solving the Stokes equations with the FEM | 322 |
| 7.1 | A quick tour of similar literature | 322 |
| 7.2 | Strong and weak forms | 323 |
| 7.3 | Which velocity-pressure pair for Stokes? | 324 |
| 7.3.1 | The compatibility condition (or LBB condition, or inf-sup condition) | 324 |

| | | |
|--------|--------------------------------------------------------------------------------------------------|-----|
| 7.3.2 | Families | 325 |
| 7.3.3 | The bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$) | 325 |
| 7.3.4 | The bi/tri-quadratic velocity - bi/tri-linear pressure element ($Q_2 \times Q_1$) | 326 |
| 7.3.5 | The bi/tri-quadratic velocity - discontinuous linear pressure element ($Q_2 \times P_{-1}$) | 327 |
| 7.3.6 | The biquadratic velocity - discontinuous bilinear pressure element ($Q_2 \times Q_{-1}$) | 327 |
| 7.3.7 | The stabilised bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$ -stab) | 328 |
| 7.3.8 | The stabilised bi/tri-linear velocity - bi/tri-linear pressure element ($Q_1 \times Q_1$ -stab) | 339 |
| 7.3.9 | The Rannacher-Turek element - rotated $Q_1 \times P_0$ | 340 |
| 7.3.10 | The $P_1 \times P_0$ pair | 341 |
| 7.3.11 | The $P_2 \times P_0$ pair | 342 |
| 7.3.12 | The $Q_2 \times Q_0$ pair | 342 |
| 7.3.13 | The $P_1 \times P_1$ -stabilised pair | 342 |
| 7.3.14 | The $P_1^+ \times P_1$ (MINI) pair in 2D & 3D | 342 |
| 7.3.15 | The $P_2 \times P_1$ pair | 343 |
| 7.3.16 | The $P_2^+ \times P_{-1}$ pair (Crouzeix-Raviart) | 343 |
| 7.3.17 | The $P_2^+ \times P_1$ pair | 344 |
| 7.3.18 | The $P_2 \times (P_1 + P_0)$ pair | 344 |
| 7.3.19 | The $Q_2 \times (Q_1 + Q_0)$ pair | 346 |
| 7.3.20 | The $P_3 \times P_2$ pair | 346 |
| 7.3.21 | The Raviart-Thomas family | 346 |
| 7.3.22 | The Bernaud-Raugel pair | 346 |
| 7.3.23 | The Scott-Vogelius pair | 347 |
| 7.3.24 | The BDM (Brezzi-Douglas-Marini) pair | 347 |
| 7.3.25 | The DSSY pair | 347 |
| 7.3.26 | The Han pair | 350 |
| 7.3.27 | The Divergence-free nonconforming $P_1^{NC} \times P_0$ pair | 352 |
| 7.3.28 | The Chen nonconforming $Q_1 \times Q_0$ pair (?) | 353 |
| 7.3.29 | Other FE element pairs | 354 |
| 7.3.30 | A note about incompressibility and standard mixed methods | 354 |
| 7.4 | The penalty approach for viscous flow | 355 |
| 7.5 | The mixed FEM for viscous flow | 360 |
| 7.5.1 | In three dimensions | 360 |
| 7.5.2 | Revisiting the penalty method | 368 |
| 7.5.3 | A much more compact derivation of the Stokes matrix blocks | 368 |
| 7.5.4 | Pressure scaling | 370 |
| 7.5.5 | Going from 3D to 2D | 371 |
| 7.5.6 | The cylindrical axisymmetric case | 373 |
| 7.6 | Mappings & Jacobians | 377 |
| 7.6.1 | General case | 377 |
| 7.6.2 | Linear mapping on a triangle | 378 |
| 7.6.3 | Bilinear mapping (Q_1) on a quadrilateral | 379 |
| 7.6.4 | Biquadratic mapping of a straight-edge face Q_2 element | 383 |
| 7.6.5 | Biquadratic mapping of a not-so straight-line face Q_2 element | 385 |
| 7.6.6 | Bilinear, biquadratic and bicubic mapping in an annulus | 385 |
| 7.6.7 | Biquadratic mapping - the middle node conundrum | 389 |
| 7.6.8 | The Double Jacobian approach | 393 |
| 7.7 | Solving the elastic equations | 400 |
| 7.8 | The case against the $Q_1 \times P_0$ element | 405 |
| 7.9 | Isoviscous Stokes for incompressible flow | 407 |
| 7.10 | $Q_1 \times P_0$ macro-elements | 409 |

| | | |
|----------|-----------------------------------------------------------------------|------------|
| 7.11 | Solving the Stokes system | 411 |
| 7.11.1 | When using the penalty formulation | 412 |
| 7.11.2 | Uzawa algorithms and the Schur complement approach | 412 |
| 7.11.3 | Conjugate gradient and the Schur complement approach | 415 |
| 7.11.4 | Generalized Conjugate Residual approach (Geenen <i>et al.</i> (2009)) | 419 |
| 7.11.5 | Using MINRES a la Burstedde <i>et al.</i> (2008) | 420 |
| 7.11.6 | The Augmented Lagrangian approach | 421 |
| 7.11.7 | The SIMPLE method | 423 |
| 7.11.8 | The GMRES approach - NOT FINISHED | 427 |
| 7.12 | Boundary conditions | 428 |
| 7.12.1 | Imposing Dirichlet boundary conditions | 428 |
| 7.12.2 | In-out flux boundary conditions for lithospheric models | 430 |
| 7.12.3 | Periodic boundary conditions | 431 |
| 7.12.4 | Free-slip boundary conditions on annulus | 432 |
| 7.13 | Open boundary conditions | 434 |
| 7.13.1 | Two-dimensional case - $Q_1 \times P_0$ elements | 435 |
| 7.13.2 | Three-dimensional case - $Q_1 \times P_0$ elements | 436 |
| 7.13.3 | Two-dimensional case - $Q_2 \times Q_1$ elements | 437 |
| 7.13.4 | Two-dimensional case - Linear triangle elements | 441 |
| 7.14 | About nullspaces | 443 |
| 7.14.1 | Pressure normalisation, nullspace | 443 |
| 7.14.2 | Removing rotational nullspace | 444 |
| 8 | The Discontinuous Galerkin Finite Element Method (DG-FEM) | 448 |
| 8.1 | First-order advection ODE in 1D | 449 |
| 8.2 | Steady state diffusion in 1D | 451 |
| 8.3 | Time-dependent diffusion PDE in 1D | 456 |
| 8.4 | Time-dependent advection PDE in 1D | 458 |
| 8.5 | Steady-state diffusion in 2D | 461 |
| 8.5.1 | The special case of linear rectangular elements | 471 |
| 8.5.2 | The special case of linear triangular elements | 474 |
| 8.6 | Time-dependent diffusion PDE in 2D | 476 |
| 8.7 | Stokes equations | 477 |
| 9 | Additional techniques, features, measurements | 478 |
| 9.1 | Dealing with a free surface (and mesh deformation) | 479 |
| 9.2 | Convergence criterion for nonlinear iterations | 491 |
| 9.3 | Strain weakening | 494 |
| 9.4 | Assigning values to quadrature points | 496 |
| 9.5 | Matrix (Sparse) storage | 499 |
| 9.5.1 | 2D domain - Q_1 - One degree of freedom per node | 499 |
| 9.5.2 | 2D domain - Q_1 - Symmetric matrix CSR storage | 501 |
| 9.5.3 | 2D domain - Q_1 - Two degrees of freedom per node | 501 |
| 9.5.4 | 2D domain - Q_2 - Two degrees of freedom per node | 504 |
| 9.5.5 | 3D domain - Q_1 - CSR storage - One degree of freedom | 507 |
| 9.5.6 | 3D domain - Q_2 - CSR storage - one degree of freedom | 508 |
| 9.5.7 | Matrix Storage in fieldstone | 509 |
| 9.5.8 | About Sparse Matrix-Vector multiplication | 509 |
| 9.5.9 | SpMV and SpMV-T with the CSR format - a concrete example | 510 |
| 9.6 | Mesh generation | 514 |

| | | |
|--------|-------------------------------------------------------------------------------|-----|
| 9.6.1 | Quadrilateral-based meshes | 514 |
| 9.6.2 | Delaunay triangulation and Voronoi cells, and triangle-based meshes | 516 |
| 9.6.3 | Tetrahedra | 520 |
| 9.6.4 | Hexahedra | 521 |
| 9.6.5 | Adaptive Mesh Refinement | 521 |
| 9.6.6 | Conformal Mesh Refinement | 525 |
| 9.6.7 | Stretching the mesh | 527 |
| 9.6.8 | Meshes in an annulus | 529 |
| 9.6.9 | Meshes in/on a hollow sphere | 529 |
| 9.7 | Pressure smoothing/filtering/recovery for $Q_1 \times P_0$ elements | 535 |
| 9.7.1 | Scheme 1 | 535 |
| 9.7.2 | Schemes 2,3 | 536 |
| 9.7.3 | Scheme 4 | 537 |
| 9.7.4 | Scheme 5 - Least squares | 538 |
| 9.7.5 | Scheme 6 - Consistent pressure recovery | 541 |
| 9.7.6 | Scheme 7 | 543 |
| 9.7.7 | Scheme 8 - bilinear interpolation | 543 |
| 9.8 | The value of the timestep | 545 |
| 9.9 | Exporting data to vtk/vtu format | 546 |
| 9.10 | Runge-Kutta methods | 552 |
| 9.11 | Am I in or not? - finding reduced coordinates | 556 |
| 9.12 | Error measurements and convergence rates | 561 |
| 9.13 | The initial temperature field | 564 |
| 9.14 | The consistent boundary flux (CBF) | 569 |
| 9.15 | Computing gradients - the recovery process | 573 |
| 9.16 | Tracking materials and/or interfaces | 574 |
| 9.17 | Static condensation | 587 |
| 9.18 | Measuring incompressibility | 588 |
| 9.19 | Picard and Newton | 589 |
| 9.19.1 | Defect correction formulation | 590 |
| 9.20 | Parallel or not? | 592 |
| 9.21 | Corner flow | 595 |
| 9.22 | Surface processes | 597 |
| 9.23 | Geometric multigrid | 599 |
| 9.24 | Algebraic multigrid | 601 |
| 9.25 | Computing depth | 602 |
| 9.26 | The Geoid | 603 |
| 9.27 | Mixing and stirring, the Lyapunov time/exponent | 606 |
| 9.27.1 | The Lyapunov exponent | 607 |
| 9.27.2 | configurational 'Shannon' entropy | 609 |
| 9.27.3 | Literature to sort out | 609 |
| 9.28 | Phase transitions | 610 |
| 9.29 | Implementation of an elasto-viscous rheology | 612 |
| 9.30 | Interpolation inside an element | 616 |
| 9.31 | Conservative Velocity Interpolation (CVI) | 621 |
| 9.31.1 | A few remarks about Wang <i>et al.</i> (2015) | 621 |
| 9.31.2 | In 2D with Q_1 basis functions - Naive approach | 622 |
| 9.31.3 | In 2D with Q_1 basis functions - better approach | 625 |
| 9.31.4 | Comparison with Wang <i>et al.</i> (2015) for 2D | 626 |
| 9.31.5 | In 3D with Q_1 basis functions - Naive approach | 628 |

| | | |
|-----------|-------------------------------------------------------------------------------------------------|------------|
| 9.31.6 | In 3D with Q_1 basis functions - better approach | 631 |
| 9.31.7 | Comparison with Wang <i>et al.</i> (2015) for 3D | 639 |
| 9.31.8 | In 2D with P_1 basis functions - what about triangles? | 643 |
| 9.31.9 | In 2D with Q_2 basis functions - Naive approach | 647 |
| 9.32 | Computing field derivatives -WIP | 655 |
| 9.33 | Iterative solvers | 658 |
| 9.34 | Weak seeds in extension modelling | 662 |
| 9.35 | Computing the volume of a hexahedron | 664 |
| 9.36 | Bandwidth reduction, matrix reordering | 666 |
| 9.37 | Scaling between dimensioned and dimensionless quantities | 668 |
| 9.38 | Spectral methods | 670 |
| 10 | Geodynamics GEO3-1313 syllabus (Utrecht University) | 671 |
| 10.1 | Introduction | 671 |
| 10.2 | Global internal structure and temperature of the Earth | 673 |
| 10.3 | The moment of inertia of a spherically symmetric density distribution | 676 |
| 10.4 | Density, gravity and pressure in the Earth | 680 |
| 10.5 | The gravitational potential for spherical problems | 689 |
| 10.5.1 | The gravity and pressure field for parameterized density models with self-gravitation | 690 |
| 10.5.2 | The pressure effect on density | 693 |
| 10.5.3 | Adiabatic density distribution | 698 |
| 10.5.4 | Earth's chemical composition | 701 |
| 10.5.5 | Phase transitions as anchor points of the geotherm | 703 |
| 10.6 | geostationary orbit | 706 |
| 10.7 | Programming exercises - February 2024 | 709 |
| 10.8 | Exam - February 2020 | 714 |
| 10.9 | Exam - March 2021 | 718 |
| 10.10 | WORK in PROGRESS. DUH. | 722 |
| 10.11 | Gravity benchmarks | 726 |
| 10.11.1 | Buried sphere (3D) | 726 |
| 10.11.2 | Buried horizontal cylinder (3D) | 727 |
| 10.11.3 | Buried column (2D) | 727 |
| 10.11.4 | Buried columns (2D) | 727 |
| 10.11.5 | Uniform layer of rock | 727 |
| 10.11.6 | A constant density shell (Root <i>et al.</i> , 2021) | 728 |
| 10.11.7 | The WINTERC mono-layer benchmark (Root <i>et al.</i> , 2021) | 730 |
| 10.11.8 | Moho benchmark (Root <i>et al.</i> , 2021) | 732 |
| 10.11.9 | Gravity potential and gravity field of a two-layer spherically symmetric planet | 733 |
| 10.12 | Gravity forward calculations in practice | 738 |
| 10.13 | Instruments to measure gravity | 746 |
| 10.14 | Gravity anomalies | 748 |
| 10.15 | Gravity reductions | 749 |
| 10.16 | How not to think about gravity (or Earth Sciences) | 750 |
| 11 | Mantle Dynamics GEO4-1416 syllabus (Utrecht University) | 753 |
| 11.0.1 | The continuity equation | 754 |
| 11.1 | Review of some essentials of continuum mechanics | 754 |
| 11.1.1 | Stress | 755 |
| 11.1.2 | Force balance equation of a continuum at rest | 755 |

| | | |
|-----------|-----------------------------------------------------------------------------------------------|------------|
| 11.1.3 | The material derivative | 756 |
| 11.1.4 | The material derivative of a material volume integral | 757 |
| 11.1.5 | Diffusion processes | 758 |
| 11.2 | The basic equations of continuum mechanics | 758 |
| 11.2.1 | The continuity equation | 758 |
| 11.2.2 | The general equation of motion (momentum equation) | 759 |
| 11.2.3 | Velocity gradient, strain rate, and rotation rate | 760 |
| 11.2.4 | Pressure and stress | 761 |
| 11.3 | Constitutive equations | 762 |
| 11.3.1 | Linear rheology | 762 |
| 11.3.2 | Non-linear rheology | 764 |
| 11.4 | The Navier-Stokes equation | 765 |
| 11.5 | Density perturbations as a driving force for mantle convection | 766 |
| 11.6 | Two-dimensional formulation for incompressible fluids: the stream function approach | 767 |
| 11.6.1 | Application of the stream function approach: Post-glacial rebound | 768 |
| 11.7 | The energy equation | 772 |
| 11.8 | The equation of state | 778 |
| 11.8.1 | The complete set of perturbation equations | 779 |
| 11.9 | Scaling of equations | 779 |
| 11.10 | The Boussinesq approximation | 781 |
| 11.10.1 | The Rayleigh-Bénard convection | 781 |
| 11.10.2 | Linear stability analysis (the onset of convection problem) | 784 |
| 11.11 | Video resources | 791 |
| 11.12 | Computer practicals | 792 |
| 11.12.1 | Introduction | 792 |
| 11.12.2 | Reminder of the governing model equations | 792 |
| 11.12.3 | Numerical solution of the equations | 793 |
| 11.12.4 | Two-dimensional convection in a unit box | 794 |
| 11.12.5 | Obtaining the python code | 795 |
| 11.12.6 | Experiments | 795 |
| 12 | Manufactured solutions & numerical benchmarks | 798 |
| 12.1 | The method of manufactured solutions | 798 |
| 12.1.1 | The repository | 798 |
| 12.1.2 | Manufactured solution in Donea and Huerta [341] (book) | 800 |
| 12.1.3 | Manufactured solution in Dohrmann and Bochev [336] (2004) | 803 |
| 12.1.4 | Analytical benchmark III - "DB3D" | 805 |
| 12.1.5 | Analytical benchmark IV - "Bercovier & Engelman" | 806 |
| 12.1.6 | Analytical benchmark VI - "Ilinca & Pelletier" | 807 |
| 12.1.7 | Analytical benchmark VII - "grooves" | 807 |
| 12.1.8 | Analytical benchmark VIII - "Kovaszny" | 810 |
| 12.1.9 | Analytical benchmark IX - "VJ2" | 811 |
| 12.1.10 | Manufactured solution in John, Linke, Merdon, Neilan, and Rebholz [655] | 811 |
| 12.1.11 | Manufactured solution in Lamichhane [741] (2017) | 814 |
| 12.1.12 | Manufactured solution in Mu and Ye [911] | 816 |
| 12.1.13 | Manufactured solution in Boffi, Cavallini, Gardini, and Gastaldi [110] (2012) | 817 |
| 12.1.14 | Manufactured solution in John [650] (book) | 819 |
| 12.1.15 | Manufactured solution in John, Kaiser, and Novo [652] | 822 |
| 12.1.16 | Manufactured solution in John [649] (1998) on a disc | 823 |
| 12.1.17 | Annulus with kinematical b.c. - pure rotation | 823 |

| | | |
|---------|---------------------------------------------------------------------------|-----|
| 12.1.18 | Viscous beam under extension | 825 |
| 12.1.19 | Channel flow with Herschel-Bulkley rheology | 826 |
| 12.1.20 | Flow in a square using Stream Functions | 835 |
| 12.1.21 | One-dimensional advection-diffusion equation | 836 |
| 12.1.22 | Annulus with kinematical b.c. - shear flow | 836 |
| 12.1.23 | Generic framework for 3D solution in Cartesian coordinates | 838 |
| 12.1.24 | 2D Analytical benchmark XII | 843 |
| 12.1.25 | 2D analytical benchmark from Burman & Hansbo (2006) | 843 |
| 12.1.26 | 2D analytical benchmark from Cioncolini & Boffi (2019) | 844 |
| 12.1.27 | 3D Analytical benchmark XIII | 844 |
| 12.1.28 | 2D Analytical benchmark XIV | 844 |
| 12.1.29 | Poisson equation on 3D shell | 846 |
| 12.1.30 | SolCx | 846 |
| 12.1.31 | SolKz | 847 |
| 12.1.32 | SolVi | 848 |
| 12.1.33 | Simple shear heating | 849 |
| 12.1.34 | 2D solution with nontrivial interface jump | 850 |
| 12.1.35 | 3D solution with nontrivial interface jump | 850 |
| 12.2 | Geodynamical benchmarks | 851 |
| 12.2.1 | Poiseuille flow | 855 |
| 12.2.2 | Relaxation of sinusoidal topography | 856 |
| 12.2.3 | the plastic brick | 857 |
| 12.2.4 | Time-dependent benchmark in an annulus | 860 |
| 12.2.5 | Convection in 2D-box | 861 |
| 12.2.6 | The sinker problem | 862 |
| 12.2.7 | The hot blob problem | 864 |
| 12.2.8 | The punch/indenter problem in 2D | 865 |
| 12.2.9 | Driven cavity with analytical solution | 866 |
| 12.2.10 | Viscous flow around a cylinder in 2D and 3D | 869 |
| 12.2.11 | Heat flow around a cylinder | 869 |
| 12.2.12 | Thermal diffusion of half-cooling space | 870 |
| 12.2.13 | Laplace equation on a semi infinite plate | 870 |
| 12.2.14 | Slab detachment benchmark | 873 |
| 12.2.15 | Layered flow with viscosity contrast | 873 |
| 12.2.16 | The annulus convection benchmark # 1 | 876 |
| 12.2.17 | The annulus convection benchmark #2 - no slip bc | 887 |
| 12.2.18 | Rayleigh-Bénard convection for silicon oil | 894 |
| 12.2.19 | Rayleigh-Taylor experiment of van Keken <i>et al.</i> (1997) | 895 |
| 12.2.20 | (Instantaneous) Sinking block (2D) | 913 |
| 12.2.21 | (Instantaneous) Stokes sphere (3D) | 918 |
| 12.2.22 | (Instantaneous) Stokes sphere (2D) | 925 |
| 12.2.23 | Stokes sphere (2D) in fluid with deformable free surface | 932 |
| 12.2.24 | Relaxation of topography (Cramer <i>et al.</i> , 2012) | 944 |
| 12.2.25 | 3D spherical shell convection benchmark | 945 |
| 12.2.26 | 2D convection benchmark ('Blankenbach <i>et al.</i> benchmark') | 946 |
| 12.2.27 | Subduction 'benchmark' of Schmeling <i>et al.</i> (2008) | 948 |
| 12.2.28 | Thin layer entrainment | 949 |

| | |
|---------------------------------------------------------------------------------------------------------------------------------------------|-------------|
| 13 Vorticity-stream function approach | 953 |
| 13.1 Vorticity-stream function formulation of the isoviscous Navier-Stokes equation | 953 |
| 13.2 Vorticity-stream function formulation of the isoviscous Stokes equation in 2d | 955 |
| 13.3 Vorticity-stream function formulation of the non-isoviscous Stokes equation in 2d . . . | 956 |
| 13.4 Boundary conditions | 958 |
| 13.4.1 Free slip, 'stress free' boundary conditions | 958 |
| 13.4.2 No slip | 959 |
| 13.4.3 Stress b.c. | 961 |
| 13.4.4 Line of symmetry | 961 |
| 13.5 Pressure Poisson Equation for the isoviscous N-S equation | 961 |
| 13.6 Vorticity-stream function formulation in 3D | 963 |
| 13.6.1 Cartesian domain (1) | 964 |
| 13.6.2 Cartesian domain (2) | 964 |
| 13.6.3 Cartesian domain (3) | 965 |
| 13.6.4 Spherical shell | 966 |
| 13.7 Vorticity-stream function formulation in polar/cylindrical coordinates | 966 |
| 13.8 Vorticity-stream function formulation in spherical coordinates for three-dimensional incompressible flow with axisymmetry | 966 |
| 13.9 Numerical approach | 967 |
| 13.9.1 Finite differences | 968 |
| 13.9.2 Finite elements | 969 |
| 13.10 Algorithm for stream function-vorticity formulation | 971 |
| 13.11 The nondimensional equations | 971 |
| 13.11.1 Isoviscous case | 971 |
| 13.12 Incorporation of phase changes | 973 |
| 13.13 The energy equation | 973 |
| 13.14 Remark, misc | 973 |
| 13.15 Literature | 973 |
| 14 Heat Transfer & convection in a porous medium | 979 |
| 14.0.1 Darcy's law for groundwater movement | 979 |
| 14.0.2 The equations of non-isothermal fluid flow in a porous medium | 980 |
| 14.0.3 Weak form and discretisation | 981 |
| 14.0.4 The equations in dimensionless form | 983 |
| 15 Adjoint methods | 986 |
| 16 Elasticity: physics, formulations and FEM | 987 |
| 16.1 Basic equations | 988 |
| 16.2 Plane strain | 992 |
| 16.3 Plane stress | 994 |
| 16.4 The axisymmetric case | 995 |
| 16.5 FEM: Incompressible formulation from Zienkiewicz & Taylor book | 999 |
| 16.6 Elastic parameter values for Earth materials | 1002 |
| 16.7 Benchmarks and analytical solutions | 1002 |
| 17 Visco-elasticity: physics, formulations and FEM | 1003 |
| 17.1 A remark | 1003 |
| 17.2 Analytical Benchmarks | 1003 |
| 17.2.1 the 1D solution | 1003 |
| 17.2.2 Pure shear | 1005 |

| | | |
|-----------|---------------------------------------------------------------------------------------------------|-------------|
| 17.2.3 | simple shear | 1005 |
| 17.2.4 | Rayleigh-Taylor instability | 1005 |
| 17.2.5 | stress build-up inside an elastic inclusion in a viscous matrix (Beuchert & Podlachikov | 1006 |
| 17.2.6 | Viscoelastic flow past a cylinder in a channel (Beuchert & Podlachikov) | 1007 |
| 17.2.7 | Elastic Simple Shear, von Tscharner and Schmalholz [1328] (2015) | 1009 |
| 17.2.8 | Response to load from ice sheet - Nakiboglu and Lambeck (1982) | 1010 |
| 17.3 | Numerical Benchmarks | 1010 |
| 17.3.1 | Bending of elastic slab (Gerya's book) | 1010 |
| 17.3.2 | Flexure of elastic plate (Choi et al) | 1011 |
| 17.3.3 | Elastic beam in viscous matrix - von Tscharner and Schmalholz | 1012 |
| 17.3.4 | Elastic beam in viscous matrix - Keller, May, and Kaus | 1013 |
| 17.3.5 | Boxcar load on an incompressible viscoelastic lithosphere - Wu (1992) | 1013 |
| 17.3.6 | Boxcar load on an incompressible viscoelastic lithosphere - Hampel et al. (2019) | 1014 |
| 17.3.7 | Kusznir and Bott (1977) experiment | 1015 |
| 17.3.8 | Parallel-Plate Viscometer Problem - SNAC manual | 1016 |
| 17.3.9 | Relaxation after extention - Hassani syllabus | 1017 |
| 17.3.10 | Role of elasticity in slab bending - Fourel, Goes, and Morra [406] (2014) | 1018 |
| 17.3.11 | Shear test in 2D - Farrington, Moresi, and Capitanio [387] (2014) | 1018 |
| 17.3.12 | Tortion test in 3D - Farrington, Moresi, and Capitanio [387] (2014) | 1019 |
| 17.3.13 | Cylindrical tunnel - Segall book (?) | 1020 |
| 17.3.14 | Relevant literature & various notes | 1020 |
| 18 | Geophysical data | 1022 |
| 18.1 | The PREM model | 1023 |
| 18.2 | From 1D tomography to density/temperature | 1027 |
| 18.3 | Earth radial viscosity profile | 1029 |
| 18.4 | Earth radial temperature profile | 1033 |
| 18.5 | Earth radial thermal expansion profile | 1041 |
| 18.6 | Earth radial density profile | 1042 |
| 18.7 | Earth radial thermal conductivity profile | 1045 |
| A | Matrix properties | 1046 |
| A.0.1 | Symmetric matrices | 1046 |
| A.0.2 | Eigenvalues of positive definite matrix | 1046 |
| A.0.3 | Schur complement | 1047 |
| B | Don't be a hero - unless you have to | 1048 |
| C | Some useful Python commands | 1050 |
| C.0.1 | Sparse matrices | 1050 |
| C.0.2 | condition number | 1050 |
| C.0.3 | Weird stuff | 1050 |
| C.0.4 | Making simple 2D plots | 1051 |
| C.0.5 | Making simple 3D plots of scatter | 1051 |
| C.0.6 | How to debug your code | 1052 |
| C.0.7 | Optional arguments | 1052 |
| C.0.8 | drawing and filling quadrilaterals | 1053 |

| | | |
|----------|--------------------------------------------------------------------|-------------|
| D | Some useful maths | 1054 |
| D.0.1 | Inverse of a 3x3 matrix | 1054 |
| D.0.2 | Inverse of a 3x3 matrix | 1054 |
| E | Elemental matrices for simple geometries | 1056 |
| E.0.1 | 1D segments | 1056 |
| E.0.2 | Quadrilaterals: rectangular linear elements | 1061 |
| E.0.3 | Quadrilaterals: rectangular quadratic elements | 1076 |
| E.0.4 | Hexahedra: cuboid elements | 1078 |
| E.0.5 | Triangles: linear elements | 1079 |
| F | Finite element terminology in various languages | 1087 |
| G | Fun modelling | 1088 |
| H | Beautiful/interesting images from computational geodynamics | 1091 |
| I | Working with Git from the terminal | 1095 |
| I.0.1 | Contributing to ASPECT | 1095 |
| I.0.2 | Contributing a cookbook in ASPECT | 1100 |
| I.0.3 | Contributing to fieldstone | 1101 |
| J | Writing a report as homework | 1107 |
| J.0.1 | Computational Geodynamics Report | 1111 |
| K | Analytical expressions for \mathbb{G}_{el} | 1113 |
| K.0.1 | $Q_1 \times P_0$ element - 2D | 1113 |
| K.0.2 | $Q_1 \times P_0$ element - 3D | 1115 |
| K.0.3 | $Q_1 \times Q_1$ element | 1116 |
| K.0.4 | $Q_1^+ \times Q_1$ element | 1121 |
| K.0.5 | $Q_1^+ \times Q_1$ element in 3D | 1128 |
| K.0.6 | $Q_2 \times Q_1$ element | 1132 |
| L | Computational Geophysics GEO4-1427 - projects | 1135 |
| L.0.1 | Convection in a box * | 1135 |
| L.0.2 | Corner flow subduction * | 1135 |
| L.0.3 | From 2D to 3D ** | 1136 |
| L.0.4 | Triangular linear elements */** | 1136 |
| L.0.5 | Triangular linear elements *** | 1136 |
| L.0.6 | Diffusion of topography **** | 1137 |
| L.0.7 | An example of a hand-built triangular mesh | 1137 |
| L.0.8 | How to visualise data on a triangular mesh with Paraview | 1139 |
| M | Using prisms in forward gravity modelling | 1140 |
| M.0.1 | Basic formulas | 1140 |
| M.0.2 | The gravitational potential | 1141 |
| M.0.3 | The gravity vector \vec{g} | 1145 |
| M.0.4 | The gravity gradient tensor | 1146 |
| M.0.5 | Revisiting Poisson's equation | 1149 |
| M.0.6 | Better numerical stability | 1150 |

| | |
|----------------------------------------------------------------|-------------|
| N Solutions to exercises of GEO3-1313 | 1151 |
| N.0.1 Problem 1 | 1151 |
| N.0.2 Problem 2 | 1153 |
| N.0.3 Problem 3 | 1154 |
| N.0.4 Problem 4 | 1156 |
| N.0.5 Problem 5 | 1156 |
| N.0.6 Problem 6 | 1156 |
| N.0.7 Problem 7 | 1156 |
| N.0.8 Problem 8 | 1156 |
| N.0.9 Problem 9 | 1157 |
| N.0.10 Problem 10 | 1157 |
| N.0.11 Problem 11 | 1157 |
| N.0.12 Problem 12 | 1158 |
| N.0.13 Problem 13 | 1158 |
| N.0.14 Problem 14 | 1158 |
| N.0.15 Problem 15 | 1160 |
| O A quick guide to Paraview and gnuplot | 1161 |
| O.0.1 Paraview | 1161 |
| O.0.2 gnuplot | 1165 |
| P A few L^AT_EX features | 1174 |
| Q Linux how to | 1176 |
| R on using Fortran | 1181 |
| R.0.1 Full matrix multiplications in fortran | 1181 |
| R.0.2 A simple example of an Interface | 1182 |
| S mineral parameters | 1184 |
| S.0.1 Olivine | 1184 |
| S.0.2 Quartz | 1185 |
| S.0.3 Plagioclase | 1185 |
| S.0.4 Peridotite | 1185 |
| S.0.5 Diabase | 1185 |
| S.0.6 Gabbro | 1186 |
| S.0.7 Serpentine | 1186 |
| T Invariants | 1187 |
| U The Γ tensor in plasticity | 1200 |
| U.0.1 Computing the Γ matrix | 1200 |
| V Using gmsh | 1206 |
| W Directional derivative, total and material derivative | 1209 |
| W.0.1 Directional derivative | 1209 |
| W.0.2 Total differential | 1210 |
| W.0.3 The material derivative | 1211 |
| W.0.4 Material derivative of a volume integral | 1212 |
| General index | 1303 |
| Contributors | 1307 |

Chapter 1

Introduction

chapter1.tex

1.1 Philosophy

philosophy.tex

This document was written with my students in mind, i.e. 3rd and 4th year Geology/Geophysics students at Utrecht University. I have chosen to use as little jargon as possible unless it is a term that is commonly found in the geodynamics literature (methods paper as well as application papers). There is no mathematical proof of any theorem that may be mentioned but I will try to refer to the appropriate sources, i.e. generic Numerical Analysis, Finite Element and Linear Algebra books. If you find that this book lacks references to Sobolev spaces, Hilbert spaces, and other spaces, this book is just not for you.

The codes I provide here are by no means optimised as I have chosen code readability over code efficiency. I have also chosen to avoid resorting to multiple code files or even functions in order to favour a sequential reading of the codes. These codes are not designed to form the basis of a real life application: Existing open source highly optimised codes should be preferred, such as ASPECT [732, 560], CITCOMS [1414, 1412], LAMEM [683], PTATIN [848, 845], PYLITH [1], ... (see Appendix ??).

Concerning figures I have consciously decided not to place them inside *figure* \LaTeX environments since it does not allow for complete control over where they end up. Instead they are inserted when they are needed in the text.

All kinds of feedback is welcome on the text (grammar, typos, ...), on the text, the equations or on the code(s). You will have my eternal gratitude if you wish to contribute an example, a benchmark, a cookbook.

All the python scripts and tex files are freely available at

<https://github.com/cedrict/fieldstone>

This document is available at:

<https://cedrict.github.io/>

Disclaimer: there are many things in this huge document I probably do not fully understand, or that I am simply wrong about. I sometimes write open questions in the text about such things. My commitment is to revisit this document time and time again, until it is 99% correct. This is not a book, it has not been edited by anybody. It is not perfect in any way. I nevertheless hope it will be useful to many in the long run.

1.2 ambition & motivation

motivation.tex

I wish to provide the community with:

- a ginormous bibliography data base - simply search the pdf for keywords. The L^AT_EX bib file¹ is also available next to the manual.tex file on github;
- a go-to document for anybody who wants to know more about a particular topic in computational geodynamics;
- a useful teaching tool for researchers, teachers, students and PhD students alike;
- small, readable, educative codes.

1.3 Acknowledgements

acknowledgments.tex

I have benefitted from many discussions, lectures, tutorials, coffee machine discussions, debugging sessions, conference poster sessions, etc ... over the years. I wish to name these instrumental people in particular and in alphabetic order: Wolfgang Bangerth, Jean Braun, Taco Broerse, Rens Elbertsen, Philippe Fullsack, Menno Fraters, Anne Glerum, Timo Heister, Dave May, Robert Myhill, John Naliboff, E. Gerry Puckett, Melchior Schuh-Senlis, Michael Tetley, Lukas van de Wiel, Arie van den Berg, Eric van den Hoogen, Tom Weir, and the whole ASPECT family/team.

1.4 About the author

I have BSc in mathematics, and an MSc diploma in physics (with a specialization in musical acoustics [329]). I did my PhD at the university of Groningen (The Netherlands) titled *Thermodynamically consistent fluid particle modelling of phase separating mixtures*². Although half of the thesis deals with the re-derivation of the Navier-Stokes equations for such systems[382], the second half is concerned with the implementation of these equations with the Smoothed Particle Hydrodynamics method [1264, 1265, 1263].

I then taught physics and programming at the University of Rennes (France) for a year, after which I did a 2-year post-doc with Prof. J. Braun³ in the Geosciences department. I then did a 4-year post-doc with prof. R. Huismans⁴ at the University of Bergen (Norway), followed by a 3-year post-doc with profs. T. Torsvik and W. Spakman at the Utrecht University (The Netherlands). Since June 2015 I am assistant professor there in the 'Mantle dynamics & theoretical geophysics' group.

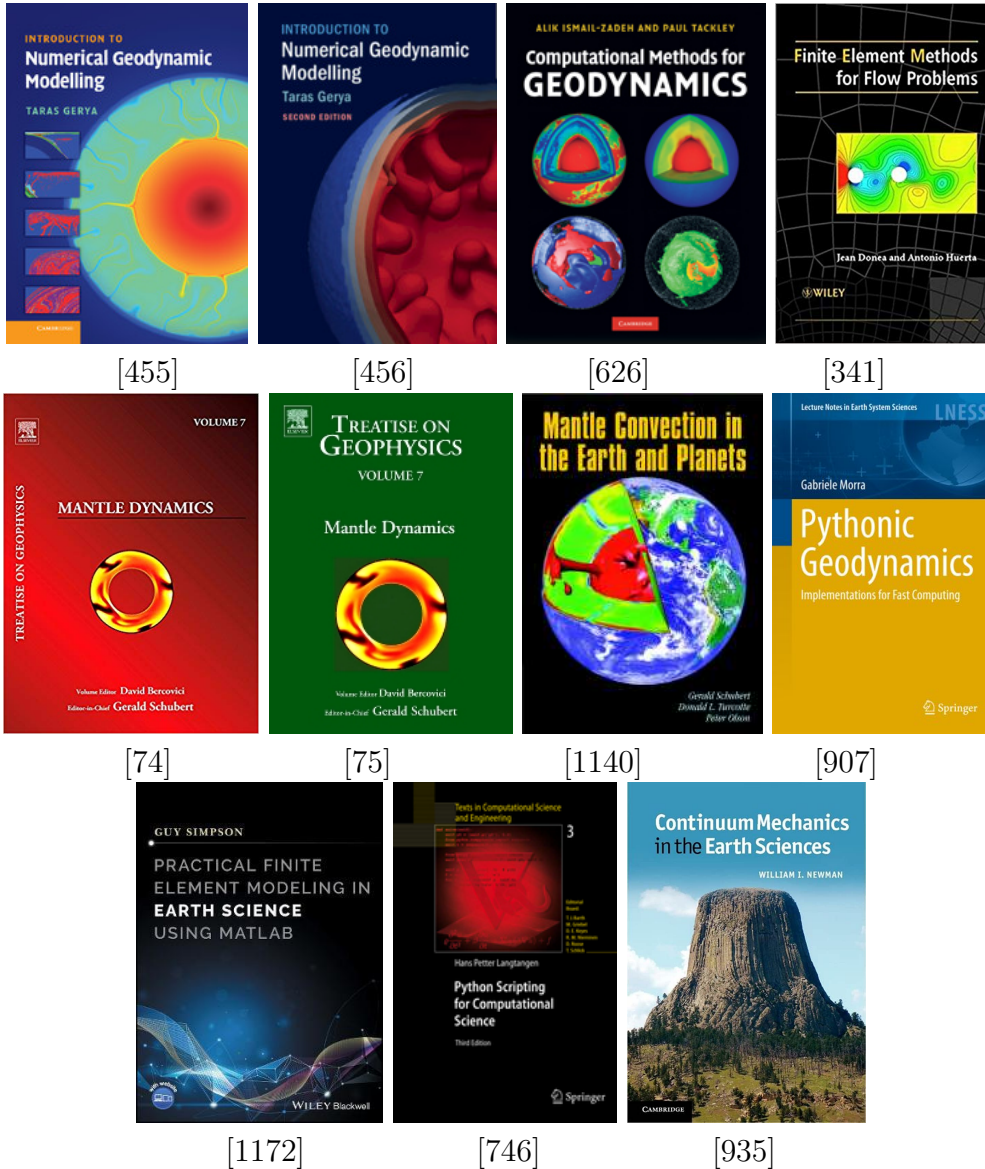
¹https://github.com/cedrict/fieldstone/blob/master/biblio_geosciences.bib

²<http://cedricthieulot.net/thesis.html>

³<https://www.gfz-potsdam.de/en/staff/jean-braun/>

⁴<https://folk.uib.no/huismans/>

1.5 Essential/relevant literature



- *Numerical modeling of Earth Systems* by Thorsten W. Becker and Boris J. P. Kaus, <http://www-udc.ig.utexas.edu/external/becker/teaching-557.html>
- *Myths & Methods in Modeling* by M. Spiegelman, <https://www.ldeo.columbia.edu/~mspieg/mmm/>
- *Computational Science I* by Matthew G. Knepley, <https://cse.buffalo.edu/~knepley/classes/caam519/Syllabus.html>
- *Introduction to Numerical Methods for Variational Problems* by Hans Petter Langtangen and Kent-Andre Mardal, <https://hplgit.github.io/fem-book/doc/pub/book/pdf/fem-book-4print.pdf>

1.6 Installing packages

Python

If numpy, scipy or matplotlib are not installed on your machine, here is how you can install them:

```
sudo apt install python3-numpy
sudo apt install python3-scipy
```

To install the umfpack solver (check?):

```
pip install --upgrade scikit-umfpack --user
```

If you need to install pip:

```
sudo apt install python3-pip
```

Julia

In order to have vim supporting the Julia language, do

```
git clone git@github.com:JuliaEditorSupport/julia-vim.git
```

and copy the content of the julia-vim folder in the .vim folder. That's it.

LaTeX

To install siunitx package:

```
sudo apt -y install texlive-science
```

To install additional fonts:

```
sudo apt-get install texlive-fonts-extra
```

To install biber package:

```
sudo apt install biber
```

1.7 What is a (real) fieldstone?

whatisafieldstone.tex



Taken from <https://en.wikipedia.org/wiki/Fieldstone>

Simply put, it is a stone collected from the surface of fields where it occurs naturally. It also stands for the bad acronym: *finite element deformation of stones* which echoes the primary application of these codes: geodynamic modelling.

1.8 Why the Finite Element method?

why.tex

The Finite Element Method (FEM) is by no means the only method to solve PDEs in geodynamics, nor it is necessarily always the best one. Other methods are employed very successfully, such as the Finite Difference Method (FDM), the Finite Volume Method (FVM), and to a lesser extent the Discrete Element Method (DEM) [1235, 360, 361, 430, 647], the Lattice-Boltzmann method [601], the Rigid Element Method [752], or the Element Free Galerkin Method (EFGM) [530]. I have been using FEM since 2008 and I do not have real experience to speak of in FVM or FDM (except for chapter 11) so I concentrate in this book on what I know best.

1.9 Notations

notations.tex

Scalars such as temperature, density, pressure, etc ... are simply obtained in L^AT_EX by using the math mode, e.g. T , ρ , p . Although it is common to lump vectors and matrices/tensors together by using bold fonts, I have decided in the interest of clarity to distinguish between those: vectors are denoted by an arrow atop the quantity, e.g. \vec{v} , \vec{g} , while matrices and tensors are in bold \mathbf{M} , $\boldsymbol{\sigma}$, etc ...

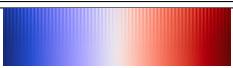











Also I use the \cdot notation between two vectors to denote a dot product $\vec{u} \cdot \vec{v} = u_i v_i$ or a matrix-vector multiplication $\mathbf{M} \cdot \vec{a} = M_{ij} a_j$. If there is no \cdot between vectors, it means that the result $\vec{a}\vec{b} = a_i b_j$ is a matrix (it is a dyadic product⁵). Case in point, $\vec{\nabla} \cdot \vec{v}$ is the velocity divergence while $\vec{\nabla}\vec{v}$ is the velocity gradient tensor.

1.10 Colour maps for visualisation

colourscale.tex

In an attempt to homogenise the figures obtained with ParaView, I have decided to use a fixed colour scale for each field throughout this document. These colour scales were obtained from <https://peterkovesi.com/projects/colourmaps> and are Perceptually Uniform Colour Maps [727].

⁵<https://en.wikipedia.org/wiki/Dyadics>

| Field | colour code | |
|-----------------------|-------------|-------------------------------------------------------------------------------------|
| Velocity/displacement | CET-D01A |  |
| Pressure | CET-L17 |  |
| Velocity divergence | CET-L01 |  |
| Density | CET-D03 |  |
| Strain rate | CET-R2 |  |
| Viscosity | CET-R3 |  |
| Temperature | CET-D09 |  |
| stress | CET-L18 |  |
| Spin tensor | CET-R1 |  |
| Composition field | CET-CBD1 |  |
| Gravity acceleration | vik |  |
| Gravity potential | roma |  |
| Vorticity | CET-L12 | |
| Stream function | CET-D02 | |

vik and roma are available at <http://www.fabiocrameri.ch/colourmaps.php>. See also Crameri *et al.* (2020) for a discussion about the misuse of colour in science communication.

1.11 How my bibliography works

mybib.tex

There is a single (large) bibliography file for this document:

`biblio_geosciences.bib`

If the paper is a single-author paper, say by Garfield⁶, published in 1978⁷, its code in my bibliography file is *garf78* (i.e. the first four letters of the name, followed by the two digits of the publication year).

If the paper was written by two authors, say Garfield and Odie, in 1987, its code will be *gaod87*, i.e. the first two letters of the first author followed by the two first letters of the second author followed by two digits.

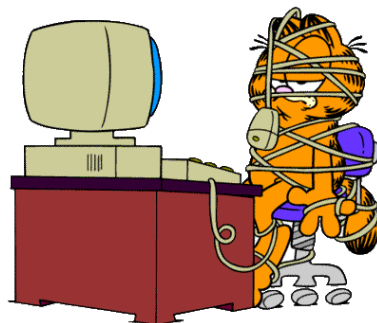
If the paper was written by three or more authors, say Garfield, Odie, John and Irene in 2003, its code will be *gaoj03*, i.e. the first two letters of the first author followed by the first letter of the second author, the first letter of the third author and the year.

If multiple papers are published the same year by the same authors, I simply append a,b,c... to the above rules.

⁶This is just an example

⁷May be not, after all, since Garfield the cat was born in 1978

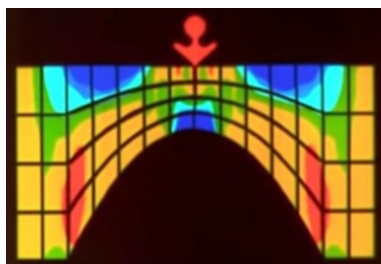
Remark. Dutch names such as 'van Hunen' or 'van den Berg' are classified under letter 'v', not 'h' or 'd' nor 'b'.



1.12 Youtube resources

youtube.tex

- https://youtu.be/aLJMDn_2-d8 [10min]



- https://youtu.be/j2_dJY_mIys [10min] Smarter Every Day channel



- <https://youtu.be/X4zd4Qpsbs8> [2min] Reversible Stokes flow (cylinder + dye)
- <https://youtu.be/wzcVT0oZJkg> [12min] (Boring Through The Earth's Crust)
- <https://youtu.be/GHjopp47vvQ> [18min] Understanding the Finite Element Method

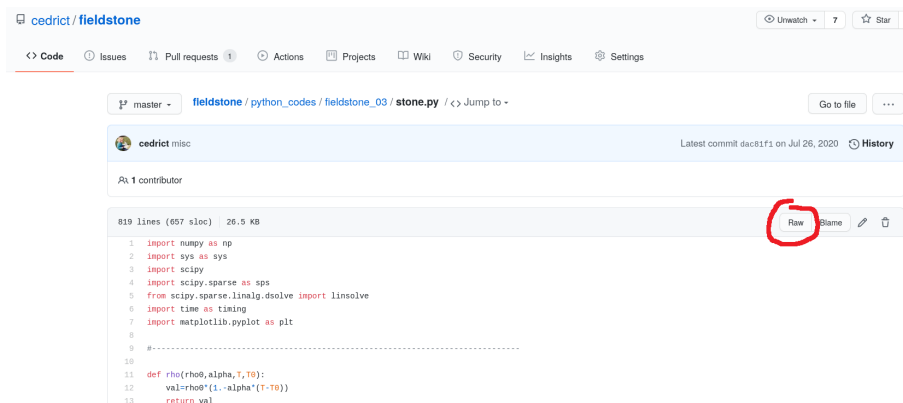


- <https://youtu.be/aPuLqiXci14> [35min] Plate Tectonics: Linking Surface Geology to Earth's Deep Interior by Clint Conrad
- <https://youtu.be/olbSuf6EGPM> [1h30] Models of mantle convection by Clint Conrad
- <https://youtu.be/aTQ-1Vpncjw> [1h30] Mantle flow for the present day by Clint Conrad
- https://youtu.be/OG5qDon-3_w [54min] Mantle flow for Earth history by Clint Conrad
- <https://youtu.be/4UAdEwbGKiM> [24min] 50 years of plate tectonics. But what is the driving force? by Clint Conrad
- https://www.esa.int/Applications/Observing_the_Earth/GOCE/Gravity_mission_still_unearthing_hidden_secrets [3min] GOCE helps create new model of crust and upper mantle
- https://youtu.be/_5q8hzF9VVE [12min] Continental drift (Wegener theory)
- <https://youtu.be/ZTRu620bIsE> [12min] Plate tectonics
- https://youtu.be/V_zsD8vXyik [5min] Heat transfer
- <https://youtu.be/q6503qA0-n4> [4min] What is sea level? (geoid)

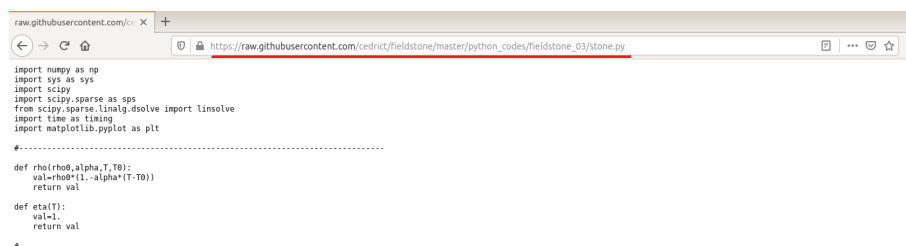
1.13 How to download a single stone

Say you wish to only download a single python program, for example stone 3. You then go to <https://github.com/cedrict/fieldstone> and click on `python_codes`, then on `fieldstone_03` and then on `stone.py`.

Then click on Raw:



And then copy the address in the address bar:



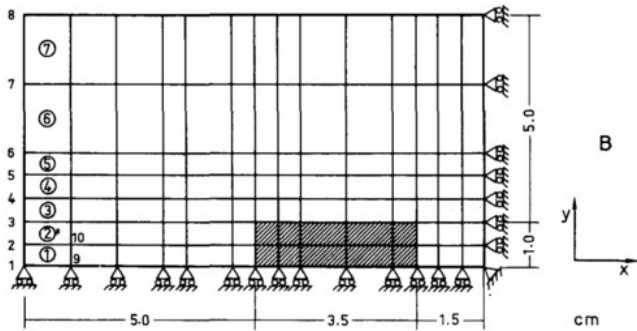
Finally, in the terminal of your Linux/Apple computer type

`wget https://raw.githubusercontent.com/cedrict/fieldstone/master/python_codes/fieldstone_03/stone.py`

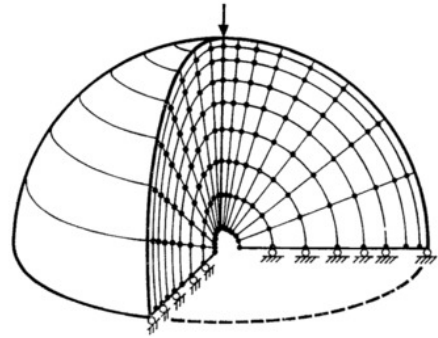
1.14 Oldies but goodies

oldies.tex

I hereunder show a few figures taken from early-ish geodynamics FEM papers.

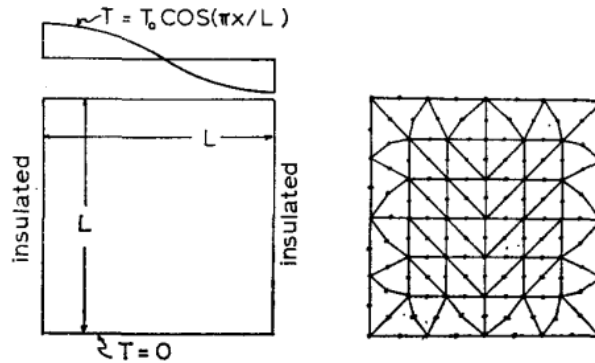


Model a boudinage structure - Stephansson and Berner [1208] (1971)

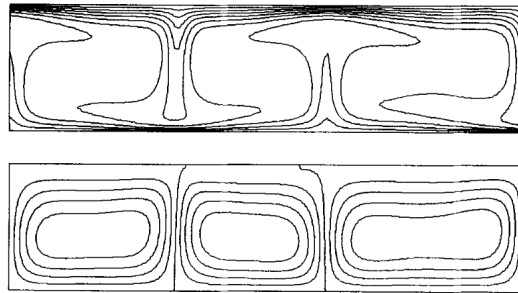


Crustal Structure from Surface Load Tilts - Beaumont and Lambert

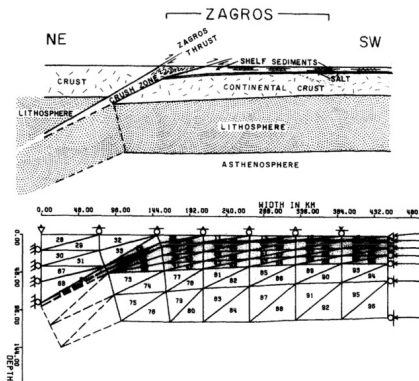
[62] (1972)



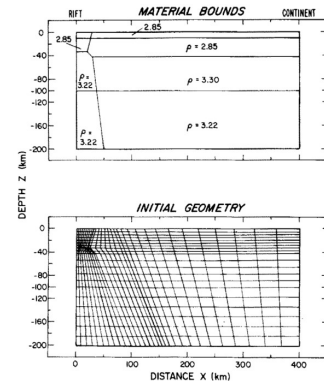
Mantle convection in a square domain - Sato and Thompson [1114] (1976)



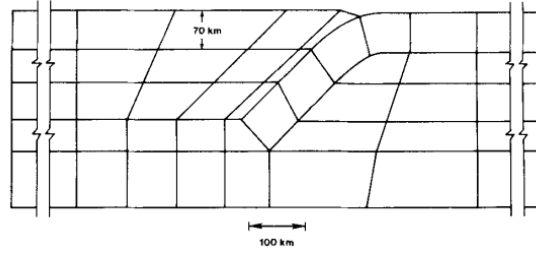
Mantle convection in a rectangular domain - Lux, Davies, and Thomas [817] (1979)



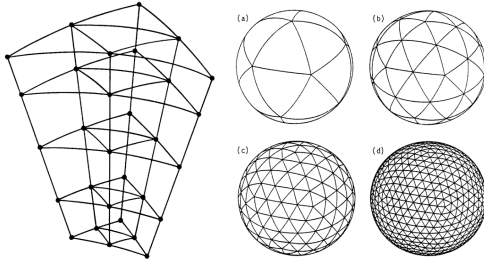
Finite element modelling of lithosphere deformation: the Zagros collision orogeny - Bird [91] (1978)



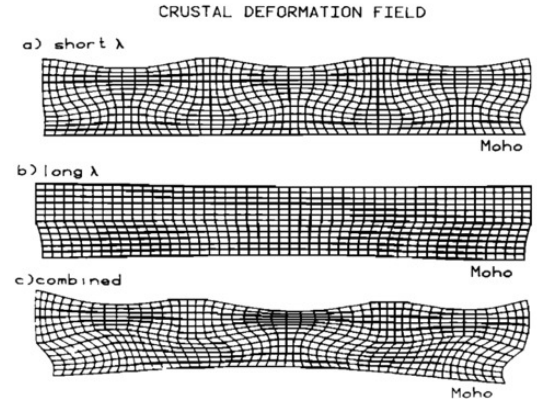
Thermal regimes, mantle diapirs and crustal stresses of continental rifts - Bridwell and Potzick [152] (1981)



Numerical models of subduction and forearc deformation - Tharp [1252] (1985)



Three-Dimensional Treatment of Convective Flow in the Earth's Mantle -
Baumgardner [57] (1985)



Lithospheric necking: a dynamic model for rift morphology - Zuber,
Parmentier, and Fletcher [1443] (1986)

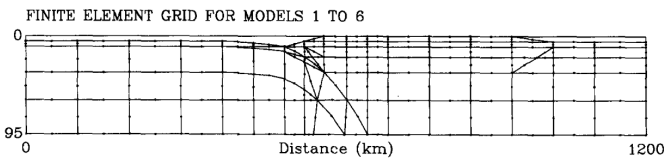
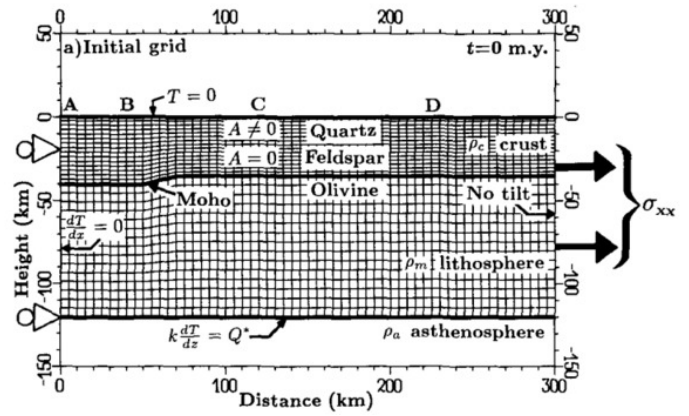
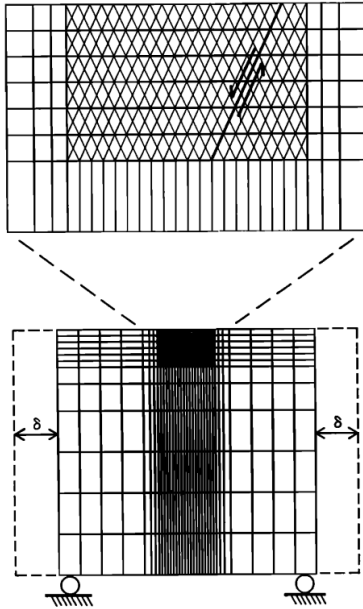


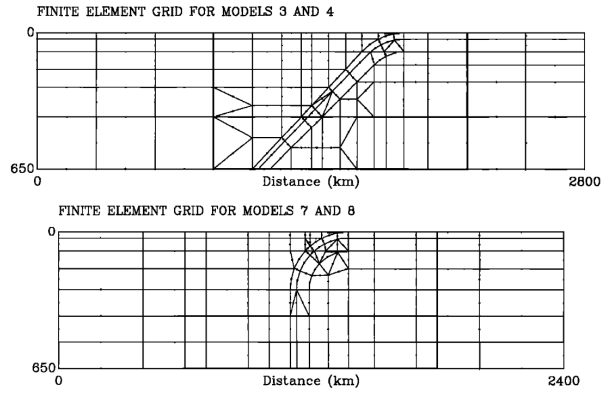
Plate boundary forces at subduction zones and trench-arc compression -
Bott, Waghorn, and Whittaker [117] (1989)



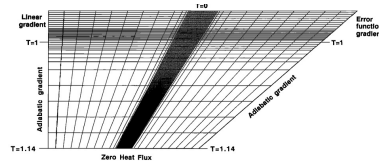
Relation between flank uplifts and the breakup unconformity at rifted
continental margins - Braun and Beaumont [141] (1989)



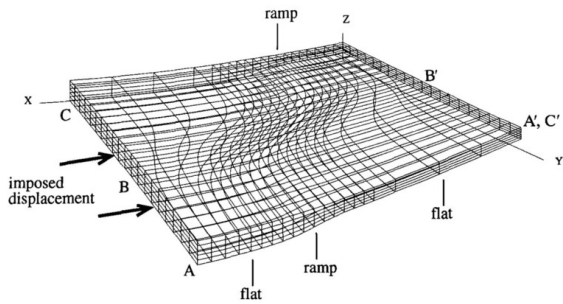
Mechanics of graben formation in crustal rocks -
Melosh and Williams Jr [863] (1989)



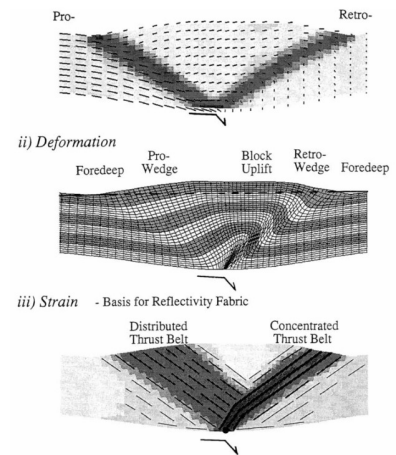
Stresses and plate boundary forces associated with subduction plate margins - [1353]
(1992)



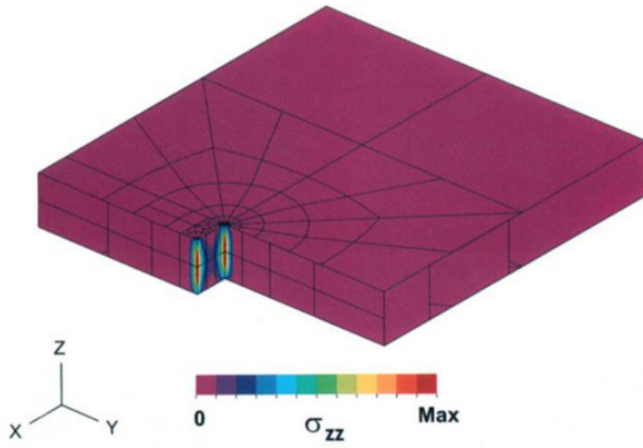
Temperature field in subduction zones - Davies and Stevenson [317] (1992)



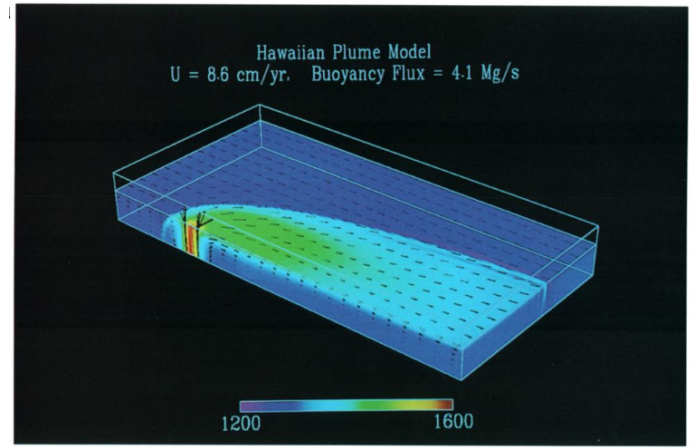
3D numerical modeling of compressional orogenies: Thrust geometry
and oblique convergence - Braun [139] (1993)



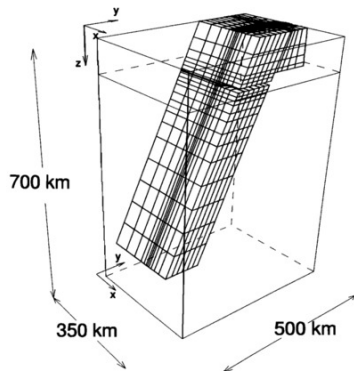
Crustal-scale compressional orogens - Beaumont and Quinlan [63]
(1994)



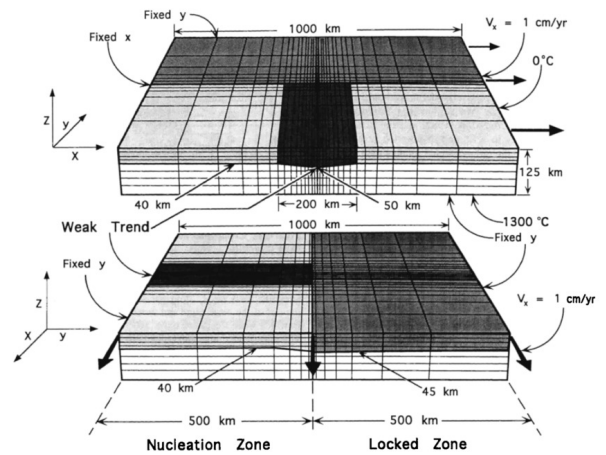
Modeling of pull-apart basins - Katzman, Brink, and Lin [677] (1995)



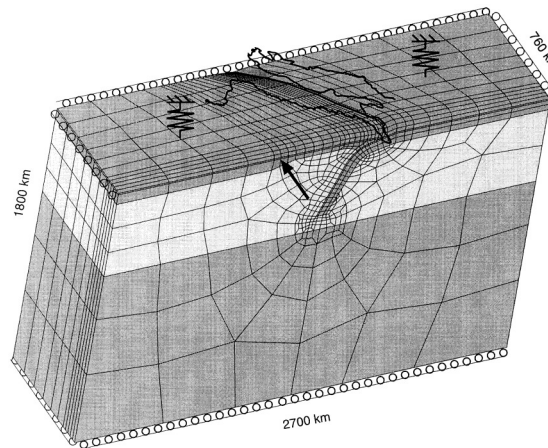
Plume-lithosphere interaction - Ribe and Christensen [1065] (1994)




3D numerical modeling of detachment of subducted lithosphere - Yoshioka and Wortel [1390] (1995)



3D dynamical model of continental rift propagation and margin plateau formation - Dunbar and Sawyer [351] (1996)



Model geometry, boundary conditions and 3-D finite element mesh used in the calculations. The circles denote a free-slip condition. The arrow denotes the velocity applied in some calculations to the southern boundary of the Tyrrhenian domain to simulate the motion of the African plate. The springs represent the buoyant restoring force applied at the surface - Negredo, Sabadini, Bianco, and Fernandez [932] (1999)

 **Relevant Literature:** Gartling [437] (1978), Anderson and Bridwell [21] (1980), Melosh and Raefsky [862] (1980), Bridwell and Anderson [151] (1980), England [374] (1982), Tharp [1252] (1985), Schubert and Anderson [1141] (1985), England and Houseman [375] (1986), Moretti and Froidevaux [905] (1986), Zuber and Parmentier [1442] (1986), Bott, Waghorn, and Whittaker [117] (1989).

Chapter 2

Physics and a bit of mathematics

chapter3.tex

2.1 Some maths

maths.tex

2.1.1 About vectors

Remark. *In this document I have chosen to (when possible) use the notation \vec{a} to denote a vector and \mathbf{a} to denote a tensor/matrix. More often than not the same notation \mathbf{a} is used for both in the literature.*

In mathematics, physics and engineering, a Euclidean vector or simply a vector is a geometric object that has magnitude (or length) and direction. Many algebraic operations on real numbers such as addition, subtraction, multiplication, and negation have close analogues for vectors.

Let \vec{v} be a vector in 3D space. Its Euclidean norm (or magnitude) is given in a coordinate-free way by

$$|\vec{v}| := \sqrt{\vec{v} \cdot \vec{v}}$$

This definition makes use of the dot product, see next section. The Euclidean norm is also called the L_2 -norm, or 2-norm. It is also sometimes noted $\|\cdot\|_2$.

In Cartesian coordinates the vector \vec{v} is given by

$$\vec{v} = \begin{pmatrix} v_x \\ v_y \\ v_z \end{pmatrix} = v_x \vec{e}_x + v_y \vec{e}_y + v_z \vec{e}_z \quad \text{with} \quad \vec{e}_x = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \vec{e}_y = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad \vec{e}_z = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

Its norm then simply writes

$$|\vec{v}| = \sqrt{v_x^2 + v_y^2 + v_z^2}$$

A unit vector is any vector with a length of one. A vector of arbitrary length can be divided by its length to create a unit vector. If \vec{a} is a vector, the corresponding unit vector is often denoted

$$\vec{e}_a = \frac{\vec{a}}{|\vec{a}|}$$

2.1.2 dot products, cross products and dyadic products

The **dot product** (or sometimes called inner product, or even scalar product) of two vectors is denoted by $\vec{a} \cdot \vec{b}$ and is defined as:

$$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta$$

where θ is the measure of the angle between \vec{a} and \vec{b} .

FIGURE

In Cartesian coordinates the dot product can also be defined as the sum of the products of the components of each vector as

$$\vec{a} \cdot \vec{b} = a_x b_x + a_y b_y + a_z b_z$$

The dot product can also be interpreted as an answer to the question “how similar are vectors \vec{a} and \vec{b} in magnitude and direction?” Indeed, if $\vec{a} = \vec{b}$ then $\theta = 0$ and $\cos \theta = 1$, while if \vec{a} is perpendicular to \vec{b} , then $\theta = \pi/2$, $\cos \theta = 0$ and $\vec{a} \cdot \vec{b} = 0$.

In Cartesian coordinates, we find that

$$\vec{v} \cdot \vec{e}_x = (v_x \vec{e}_x + v_y \vec{e}_y + v_z \vec{e}_z) \cdot \vec{e}_x = v_x \underbrace{\vec{e}_x \cdot \vec{e}_x}_{=1} + v_y \underbrace{\vec{e}_y \cdot \vec{e}_x}_{=0} + v_z \underbrace{\vec{e}_z \cdot \vec{e}_x}_{=0} = v_x$$

In this case the interpretation of $\vec{v} \cdot \vec{e}_x$ could be “how much of \vec{v} is in the direction \vec{e}_x ”.

The **cross product** (also called the vector product or outer product) of two vectors is also a vector. It is denoted $\vec{a} \times \vec{b}$ and defined as

$$\vec{c} = \vec{a} \times \vec{b} = |\vec{a}| |\vec{b}| \sin \theta \vec{n}$$

where θ is the measure of the angle between \vec{a} and \vec{b} and \vec{n} is a unit vector perpendicular to both \vec{a} and \vec{b} which completes a right-handed system.

FIGURE

The norm of the cross product, say $|\vec{c}| = |\vec{a} \times \vec{b}|$, is actually the area of the parallelogram having \vec{a} and \vec{b} as sides.

Also note that $\vec{a} \times \vec{b} = -\vec{b} \times \vec{a}$ (think about the direction of the normal vector in each case). In Cartesian coordinates the cross product can be written as

$$\vec{a} \times \vec{b} = (a_y b_z - a_z b_y) \vec{e}_x + (a_z b_x - a_x b_z) \vec{e}_y + (a_x b_y - a_y b_x) \vec{e}_z$$

Finally, let us look at the **dyadic product** of two vectors \vec{a} and \vec{b} which denoted by $\vec{a} \vec{b}^T$ (juxtaposed; no symbols, multiplication signs, crosses, dots, etc...). The result is a tensor:

$$\vec{a} = \begin{pmatrix} a_x \\ a_y \\ a_z \end{pmatrix}, \quad \vec{b} = \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix}, \quad \vec{a} \vec{b}^T = \begin{pmatrix} a_x \\ a_y \\ a_z \end{pmatrix} (b_x \ b_y \ b_z) = \begin{pmatrix} a_x b_x & a_x b_y & a_x b_z \\ a_y b_x & a_y b_y & a_y b_z \\ a_z b_x & a_z b_y & a_z b_z \end{pmatrix}$$

In conclusion the dot product yields a scalar, the cross product yields a vector and the dyadic product yields a tensor.

2.1.3 Rotation matrix

After much confusion, <https://mathworld.wolfram.com/RotationMatrix.html> is a source of clarity: one must be careful when speaking of ‘rotation matrix’. Indeed, there are two possible conventions: rotation of the axes, and rotation of the object relative to fixed axes.

We consider in \mathbb{R}^2 the matrix \mathbf{R} that rotates a given vector \vec{v} by a counterclockwise angle θ in a fixed coordinate system. It writes

$$\mathbf{R} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

with $\vec{v}' = \mathbf{R} \cdot \vec{v}$.

On the other hand, consider the matrix that rotates the coordinate system through a counterclockwise angle θ . The coordinates of the fixed vector \vec{v} in the rotated coordinate system are now given by a rotation matrix which is the transpose of the fixed-axis matrix and, as can be seen in the above diagram, is equivalent to rotating the vector by a counterclockwise angle of θ relative to a fixed set of axes, giving

$$\mathbf{R} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

In the following example we start from $\vec{v} = (2, 1)$. If we rotate the vector by 90° , the rotation matrix is given by

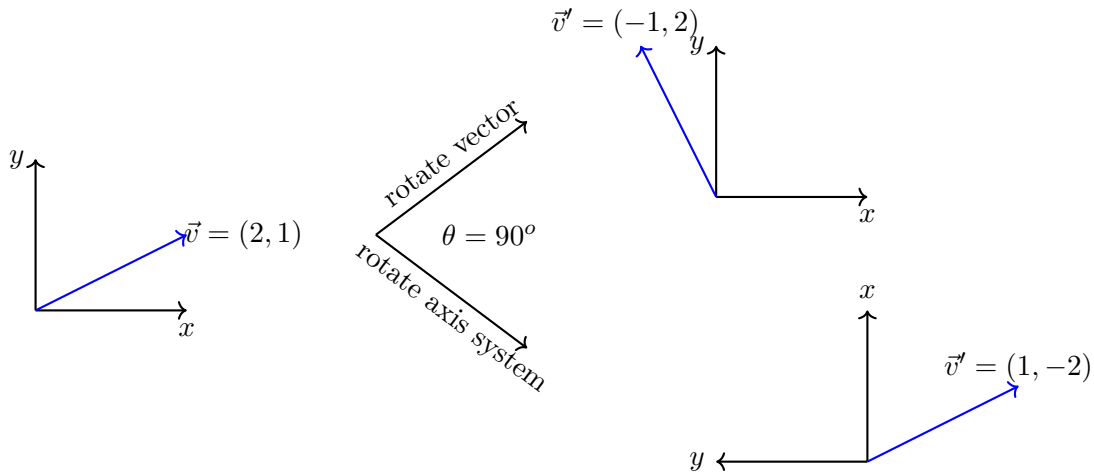
$$\mathbf{R} = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$$

so that $\vec{v}' = (-1, 2)$. If we rotate the axis by 90° , the rotation matrix is given by

$$\mathbf{R} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

and the coordinates of the resulting vector are $\vec{v}' = (1, -2)$.

(rotation_matrix.tex)

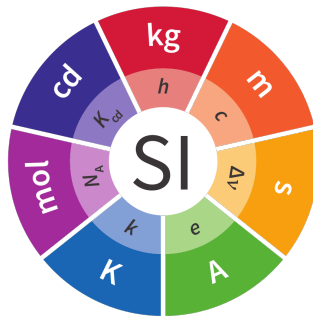


2.2 Units

nomenclature.tex

| Symbol | meaning | unit |
|----------------------------------------------------------------|------------------------------------|----------------------------------|
| t | Time | s |
| x, y, z | Cartesian coordinates | m |
| r, θ | Polar coordinates | m, - |
| r, θ, z | Cylindrical coordinates | m, -, m |
| r, θ, ϕ | Spherical coordinates | m, -, - |
| $\vec{v} = (u, v, w)$ | velocity vector ⁽¹⁾ | m s^{-1} |
| $\vec{v} = (\mathbf{v}_r, \mathbf{v}_\theta, \mathbf{v}_z)$ | velocity vector ⁽²⁾ | m s^{-1} |
| $\vec{v} = (\mathbf{v}_r, \mathbf{v}_\theta, \mathbf{v}_\phi)$ | velocity vector ⁽³⁾ | m s^{-1} |
| $\vec{u} = (\mathbf{u}_x, \mathbf{u}_y, \mathbf{u}_z)$ | displacement vector | m |
| ρ | mass density | kg m^{-3} |
| η | dynamic viscosity | Pa s |
| λ | penalty parameter | Pa s |
| T | temperature | K |
| $\vec{\nabla}$ | gradient operator | m^{-1} |
| $\vec{\nabla} \cdot$ | divergence operator | m^{-1} |
| p | pressure | Pa |
| $\dot{\epsilon}(\vec{v})$ | strain rate tensor | s^{-1} |
| $\dot{\epsilon}^d(\vec{v})$ | deviatoric strain rate tensor | s^{-1} |
| α | thermal expansion coefficient | K^{-1} |
| k | thermal conductivity | $\text{W m}^{-1} \text{K}^{-1}$ |
| C_p | Heat capacity at constant pressure | $\text{J kg}^{-1} \text{K}^{-1}$ |
| H | intrinsic specific heat production | W kg^{-1} |
| β_T | isothermal compressibility | Pa^{-1} |
| $\boldsymbol{\tau}$ | deviatoric stress tensor | Pa |
| $\boldsymbol{\sigma}$ | full stress tensor | Pa |
| θ_L | Lodé angle | - |
| λ | bulk modulus | Pa |
| μ | shear modulus | Pa |
| ν | Poisson ratio | - |
| E | Young's modulus | Pa |

(1) Cartesian coordinates; (2) Cylindrical coordinates; (3) Spherical coordinates.



Taken from Wikipedia¹. The SI logo, produced by the BIPM (International Bureau of Weights and Measures), showing the seven SI base units and the seven defining constants.

A quick note about units and \LaTeX . This document relies on the `siunitx` package². For instance, $\rho = 3300 \text{ kg m}^{-3}$ is obtained with

```
\rho = 3300~\si{\kg\per\cubic\metre}
```

¹https://en.wikipedia.org/wiki/International_System_of_Units

²<https://ctan.org/pkg/siunitx>

or

`\rho = \SI{3300}{\kg\per\cubic\metre}`

Note that the `si` command can be used outside of the math environment.

2.3 Coordinate systems

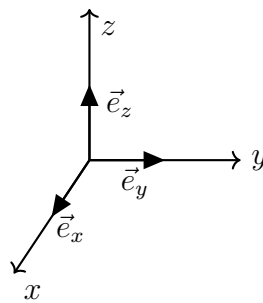
coordinate_systems.tex



2.3.1 Cartesian coordinates

The unit vectors along the x , y and z axis are \vec{e}_x , \vec{e}_y and \vec{e}_z respectively.

(tikz_cartesian_coordinates.tex)



Any vector can then be written

$$\vec{V} = V_x \vec{e}_x + V_y \vec{e}_y + V_z \vec{e}_z$$

'How much of \vec{V} is there in the x -direction' is obtained with $\vec{V} \cdot \vec{e}_x = V_x$. The gradient of a function f is

$$\vec{\nabla} f = \text{grad } f = \frac{\partial f}{\partial x} \vec{e}_x + \frac{\partial f}{\partial y} \vec{e}_y + \frac{\partial f}{\partial z} \vec{e}_z,$$

the divergence of a vector \vec{V} is

$$\vec{\nabla} \cdot \vec{V} = \frac{\partial V_x}{\partial x} + \frac{\partial V_y}{\partial y} + \frac{\partial V_z}{\partial z}$$

and the Laplace operator of a function f is:

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + \frac{\partial^2 f}{\partial z^2}$$

Finally the path increment is

$$d\vec{r} = dx \vec{e}_x + dy \vec{e}_y + dz \vec{e}_z$$

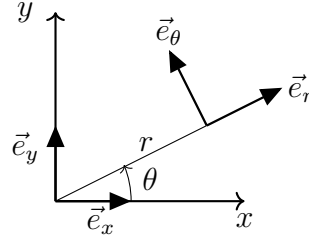
and the volume element is

$$dV = dx \, dy \, dz.$$

2.3.2 Polar coordinates

We have $r > 0$ and $\theta = [0, 2\pi[$, defined in the (x, y) -plane.

(tikz-polar-coordinates.tex)



The relation between the unit vector in Cartesian and Polar/Cylindrical coordinates is given by:

$$\begin{pmatrix} \vec{e}_r \\ \vec{e}_\theta \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \cdot \begin{pmatrix} \vec{e}_x \\ \vec{e}_y \end{pmatrix}$$

which should be read:

$$\begin{aligned} \vec{e}_r &= \cos \theta \vec{e}_x + \sin \theta \vec{e}_y \\ \vec{e}_\theta &= -\sin \theta \vec{e}_x + \cos \theta \vec{e}_y \end{aligned} \quad (2.1)$$

Obviously for $\theta = 0$ we find $\vec{e}_r = \vec{e}_x$ and $\vec{e}_\theta = \vec{e}_y$, while for $\theta = \pi/2$ then $\vec{e}_r = \vec{e}_y$ and $\vec{e}_\theta = -\vec{e}_x$.

Note that this 2×2 matrix is a rotation matrix³ corresponding to an angle $-\theta$. The inverse of this matrix always exists (we can always counter-rotate) and it then yields

$$\begin{pmatrix} \vec{e}_x \\ \vec{e}_y \end{pmatrix} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \cdot \begin{pmatrix} \vec{e}_r \\ \vec{e}_\theta \end{pmatrix}$$

so that for any vector \vec{V}

$$\begin{aligned} \vec{V} &= V_x \vec{e}_x + V_y \vec{e}_y \\ &= V_x [(\cos \theta) \vec{e}_r - (\sin \theta) \vec{e}_\theta] + V_y [(\sin \theta) \vec{e}_r + (\cos \theta) \vec{e}_\theta] \\ &= [V_x (\cos \theta) + V_y (\sin \theta)] \vec{e}_r + [-V_x (\sin \theta) + V_y (\cos \theta)] \vec{e}_\theta \\ &= V_r \vec{e}_r + V_\theta \vec{e}_\theta \end{aligned}$$

with

$$\begin{aligned} V_r &= V_x \cos \theta + V_y \sin \theta \\ V_\theta &= -V_x \sin \theta + V_y \cos \theta \end{aligned}$$

Finally the path increment is

$$d\vec{r} = dr \vec{e}_r + r \sin \theta d\theta \vec{e}_\theta$$

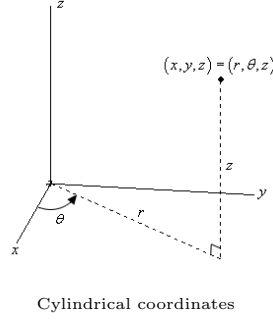
and the volume element is

$$dV = r dr d\theta.$$

The gradient, divergence and Laplacian formulae are given in the following section about the cylindrical coordinates.

³https://en.wikipedia.org/wiki/Rotation_matrix

2.3.3 Cylindrical coordinates



$$\vec{V} = V_r \vec{e}_r + V_\theta \vec{e}_\theta + V_z \vec{e}_z$$

We have

$$\begin{aligned} x &= r \cos \theta \\ y &= r \sin \theta \\ r &= \sqrt{x^2 + y^2} \end{aligned}$$

Let $f(r, \theta)$ be a function of the spatial coordinates. Its gradient is then

$$\vec{\nabla} f = \frac{\partial f}{\partial r} \vec{e}_r + \frac{1}{r} \frac{\partial f}{\partial \theta} \vec{e}_\theta + \frac{\partial f}{\partial z} \vec{e}_z$$

The divergence of a vector field \vec{V} is

$$\vec{\nabla} \cdot \vec{V} = \frac{1}{r} \frac{\partial}{\partial r} (r V_r) + \frac{1}{r} \frac{\partial V_\theta}{\partial \theta} + \frac{\partial V_z}{\partial z}$$

and the Laplacian of f is

$$\Delta f = \frac{1}{r} \frac{\partial}{\partial r} \left(r \frac{\partial f}{\partial r} \right) + \frac{1}{r^2} \frac{\partial^2 f}{\partial \theta^2} + \frac{\partial^2 f}{\partial z^2}$$

Finally the path increment is

$$d\vec{r} = dr \vec{e}_r + r \sin \theta d\theta \vec{e}_\theta + dz \vec{e}_z$$

and the volume element is

$$dV = r dr d\theta dz$$

Remark. Cylindrical coordinates can also be denoted by (ρ, θ) , (r, ϕ) or even (ρ, ϕ) . They are sometimes called "cylindrical polar coordinates" or "polar cylindrical coordinates".

The divergence of the second order tensor field \mathbf{S} in cylindrical polar coordinates is given by

$$\begin{aligned} \vec{\nabla} \cdot \mathbf{S} &= \frac{\partial S_{rr}}{\partial r} \vec{e}_r + \frac{\partial S_{r\theta}}{\partial r} \vec{e}_\theta + \frac{\partial S_{rz}}{\partial r} \vec{e}_z \\ &+ \frac{1}{r} \left[\frac{\partial S_{r\theta}}{\partial \theta} + (S_{rr} - S_{\theta\theta}) \right] \vec{e}_r + \frac{1}{r} \left[\frac{\partial S_{\theta\theta}}{\partial \theta} + (S_{r\theta} + S_{\theta r}) \right] \vec{e}_\theta + \frac{1}{r} \left[\frac{\partial S_{\theta z}}{\partial \theta} + S_{rz} \right] \vec{e}_z \\ &+ \frac{\partial S_{zr}}{\partial z} \vec{e}_r + \frac{\partial S_{z\theta}}{\partial z} \vec{e}_\theta + \frac{\partial S_{zz}}{\partial z} \vec{e}_z \end{aligned}$$

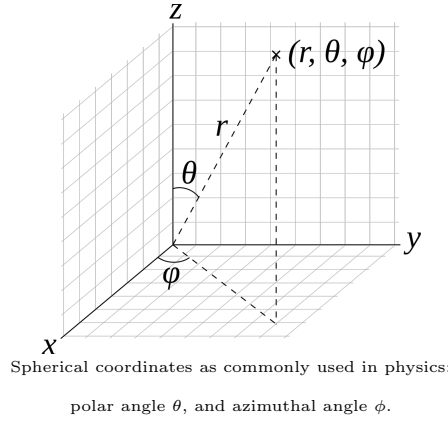
In the case of polar coordinates then all quantities featuring z (or ∂_z) are removed:

$$\begin{aligned}\vec{\nabla} \cdot \mathbf{S} &= \frac{\partial S_{rr}}{\partial r} \vec{e}_r + \frac{\partial S_{r\theta}}{\partial r} \vec{e}_\theta + \frac{1}{r} \left[\frac{\partial S_{r\theta}}{\partial \theta} + (S_{rr} - S_{\theta\theta}) \right] \vec{e}_r + \frac{1}{r} \left[\frac{\partial S_{\theta\theta}}{\partial \theta} + (S_{r\theta} + S_{\theta r}) \right] \vec{e}_\theta \\ &= \left[\frac{\partial S_{rr}}{\partial r} + \frac{1}{r} \frac{\partial S_{r\theta}}{\partial \theta} + \frac{1}{r} (S_{rr} - S_{\theta\theta}) \right] \vec{e}_r + \left[\frac{\partial S_{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial S_{\theta\theta}}{\partial \theta} + \frac{2}{r} S_{r\theta} \right] \vec{e}_\theta\end{aligned}\quad (2.2)$$

where we have assumed that the tensor \mathbf{S} is symmetric (i.e. $S_{r\theta} = S_{\theta r}$).

2.3.4 Spherical coordinates

On the following figure are represented the three Cartesian axis, a point and its spherical coordinates r, θ, ϕ :



In this case $\theta \in [0 : \pi]$ and $\phi \in]-\pi : \pi]$ and we have the following relationships:

$$r = \sqrt{x^2 + y^2 + z^2} \quad (2.3)$$

$$\theta = \arccos(z/r) \quad (2.4)$$

$$\phi = \arctan(y/x) \quad (2.5)$$

$$x = r \sin \theta \cos \phi \quad (2.6)$$

$$y = r \sin \theta \sin \phi \quad (2.7)$$

$$z = r \cos \theta \quad (2.8)$$

The inverse tangent used to compute ϕ must be suitably defined, taking into account the correct quadrant of (x, y) , which is why the `atan2` intrinsic function is used in FORTRAN for example. This is often written as follows:

$$\theta = \arctan\left(\sqrt{x^2 + y^2}, z\right) \quad (2.9)$$

$$\phi = \arctan(y, x) \quad (2.10)$$

where we formally take advantage of the two argument `arctan` function to eliminate quadrant confusion.

The path increment is expressed as:

$$d\vec{r} = dr \vec{e}_r + r d\theta \vec{e}_\theta + r \sin \theta d\phi \vec{e}_\phi \quad (2.11)$$

The gradient of a function $f(r, \theta, \phi)$ is

$$\vec{\nabla} f = \frac{\partial f}{\partial r} \vec{e}_r + \frac{1}{r} \frac{\partial f}{\partial \theta} \vec{e}_\theta + \frac{1}{r \sin \theta} \frac{\partial f}{\partial \phi} \vec{e}_\phi \quad (2.12)$$

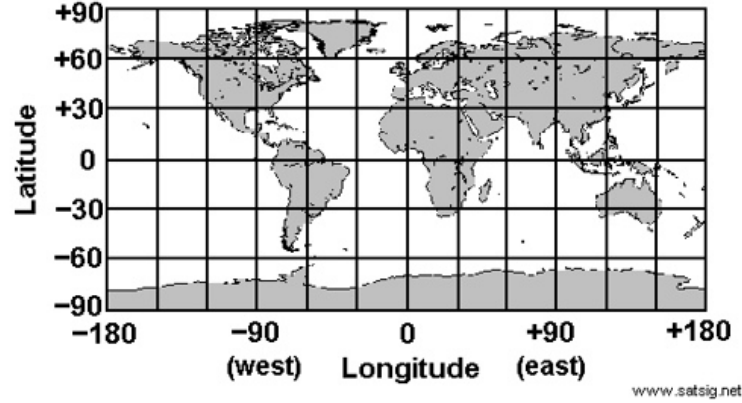
The divergence of a vector \vec{V} is

$$\vec{\nabla} \cdot \vec{V} = \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 V_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (V_\theta \sin \theta) + \frac{1}{r \sin \theta} \frac{\partial V_\phi}{\partial \phi} = 0 \quad (2.13)$$

The Laplacian of function f is given by:

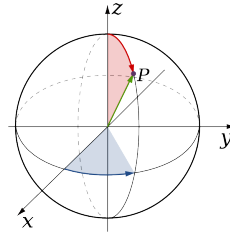
$$\Delta f = \vec{\nabla}^2 f = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial f}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial f}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2 f}{\partial \phi^2} \quad (2.14)$$

In geography one uses latitude and longitude, represented hereunder:



- Latitude $\in [-90 : 90]$, or $\in [-\pi/2 : \pi/2]$
- Longitude $\in] - 180 : 180]$, or $\in] - \pi : \pi]$

Since the colatitude is the complementary angle of the latitude, i.e. the difference between 90 and the latitude, where southern latitudes are denoted with a minus sign, θ as shown above is actually is the colatitude. The colatitude is shown in red on the following figure:



The volume of a sphere of radius R is easily obtained by computing

$$\begin{aligned} V_{sphere} &= \iiint_{sphere} dV \\ &= \int_0^R r^2 dr \int_0^\pi \sin \theta d\theta \int_0^{2\pi} d\phi \\ &= \frac{1}{3} R^3 \cdot 2 \cdot 2\pi \\ &= \frac{4}{3} \pi R^3 \end{aligned} \quad (2.15)$$

The volume of a spherical shell of inner radius R_i and outer radius R_o is equally easily obtained by computing

$$\begin{aligned}
V_{shell} &= \iiint_{shell} dV \\
&= \int_{R_i}^{R_o} r^2 dr \int_0^\pi \sin \theta d\theta \int_0^{2\pi} d\phi \\
&= \frac{1}{3} (R_o^3 - R_i^3) \cdot 2 \cdot 2\pi \\
&= \frac{4}{3} \pi (R_o^3 - R_i^3)
\end{aligned} \tag{2.16}$$

The spherical unit vectors are related to the Cartesian unit vectors by:

$$\begin{pmatrix} \vec{e}_r \\ \vec{e}_\theta \\ \vec{e}_\phi \end{pmatrix} = \begin{pmatrix} \sin \theta \cos \phi & \sin \theta \sin \phi & \cos \theta \\ \cos \theta \cos \phi & \cos \theta \sin \phi & -\sin \theta \\ -\sin \phi & \cos \phi & 0 \end{pmatrix} \begin{pmatrix} \vec{e}_x \\ \vec{e}_y \\ \vec{e}_z \end{pmatrix}$$

and the Cartesian unit vectors are related to the spherical unit vectors by

$$\begin{pmatrix} \vec{e}_x \\ \vec{e}_y \\ \vec{e}_z \end{pmatrix} = \begin{pmatrix} \sin \theta \cos \phi & \cos \theta \cos \phi & -\sin \phi \\ \sin \theta \sin \phi & \cos \theta \sin \phi & \cos \phi \\ \cos \theta & -\sin \theta & 0 \end{pmatrix} \begin{pmatrix} \vec{e}_r \\ \vec{e}_\theta \\ \vec{e}_\phi \end{pmatrix}$$

Finally, the velocity vector \vec{v} then becomes

$$\begin{aligned}
\vec{v} &= u \vec{e}_x + v \vec{e}_y + w \vec{e}_z \\
&= u (\sin \theta \cos \phi \vec{e}_r + \cos \theta \cos \phi \vec{e}_\theta - \sin \phi \vec{e}_\phi) \\
&\quad + v (\sin \theta \sin \phi \vec{e}_r + \cos \theta \sin \phi \vec{e}_\theta + \cos \phi \vec{e}_\phi) \\
&\quad + w (\cos \theta \vec{e}_r - \sin \theta \vec{e}_\theta) \\
&= v_r \vec{e}_r + v_\theta \vec{e}_\theta + v_\phi \vec{e}_\phi
\end{aligned} \tag{2.17}$$

with

$$\begin{aligned}
v_r &= u \sin \theta \cos \phi + v \sin \theta \sin \phi + w \cos \theta \\
v_\theta &= u \cos \theta \cos \phi + v \cos \theta \sin \phi - w \sin \theta \\
v_\phi &= -u \sin \phi + v \cos \phi
\end{aligned} \tag{2.18}$$

2.3.5 Converting tensors between Cartesian and Cylindrical bases

$$\begin{aligned}
\mathbf{T}_{Cyl} &= \begin{pmatrix} T_{rr} & T_{r\theta} & T_{rz} \\ T_{\theta r} & T_{\theta\theta} & T_{\theta z} \\ T_{zr} & T_{z\theta} & T_{zz} \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} T_{xx} & T_{xy} & T_{xz} \\ T_{yx} & T_{yy} & T_{yz} \\ T_{zx} & T_{zy} & T_{zz} \end{pmatrix} \cdot \begin{pmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix} \\
\mathbf{T}_{Cart} &= \begin{pmatrix} T_{xx} & T_{xy} & T_{xz} \\ T_{yx} & T_{yy} & T_{yz} \\ T_{zx} & T_{zy} & T_{zz} \end{pmatrix} = \begin{pmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} T_{rr} & T_{r\theta} & T_{rz} \\ T_{\theta r} & T_{\theta\theta} & T_{\theta z} \\ T_{zr} & T_{z\theta} & T_{zz} \end{pmatrix} \cdot \begin{pmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{pmatrix}
\end{aligned}$$

2.3.6 Converting tensors between Cartesian and Spherical bases

Let \mathbf{T} be a tensor

$$\mathbf{T} = \begin{pmatrix} T_{xx} & T_{xy} & T_{xz} \\ T_{yx} & T_{yy} & T_{yz} \\ T_{zx} & T_{zy} & T_{zz} \end{pmatrix} \quad \mathbf{T} = \begin{pmatrix} T_{rr} & T_{r\theta} & T_{r\phi} \\ T_{\theta r} & T_{\theta\theta} & T_{\theta\phi} \\ T_{\phi r} & T_{\phi\theta} & T_{\phi\phi} \end{pmatrix}$$

in the Cartesian basis (left) and the spherical basis (right).

The two sets of components are related by

$$\begin{pmatrix} T_{xx} & T_{xy} & T_{xz} \\ T_{yx} & T_{yy} & T_{yz} \\ T_{zx} & T_{zy} & T_{zz} \end{pmatrix} = \begin{pmatrix} \sin \theta \cos \phi & \cos \theta \cos \phi & -\sin \phi \\ \sin \theta \sin \phi & \cos \theta \sin \phi & \cos \phi \\ \cos \theta & -\sin \theta & 0 \end{pmatrix} \cdot \begin{pmatrix} T_{rr} & T_{r\theta} & T_{r\phi} \\ T_{\theta r} & T_{\theta\theta} & T_{\theta\phi} \\ T_{\phi r} & T_{\phi\theta} & T_{\phi\phi} \end{pmatrix} \cdot \begin{pmatrix} \sin \theta \cos \phi & \sin \theta \sin \phi \\ \cos \theta \cos \phi & \cos \theta \sin \phi \\ -\sin \phi & \cos \phi \end{pmatrix}$$

or

$$\begin{pmatrix} T_{rr} & T_{r\theta} & T_{r\phi} \\ T_{\theta r} & T_{\theta\theta} & T_{\theta\phi} \\ T_{\phi r} & T_{\phi\theta} & T_{\phi\phi} \end{pmatrix} = \begin{pmatrix} \sin \theta \cos \phi & \sin \theta \sin \phi & \cos \theta \\ \cos \theta \cos \phi & \cos \theta \sin \phi & -\sin \theta \\ -\sin \phi & \cos \phi & 0 \end{pmatrix} \cdot \begin{pmatrix} T_{xx} & T_{xy} & T_{xz} \\ T_{yx} & T_{yy} & T_{yz} \\ T_{zx} & T_{zy} & T_{zz} \end{pmatrix} \cdot \begin{pmatrix} \sin \theta \cos \phi & \cos \theta \cos \phi \\ \sin \theta \sin \phi & \cos \theta \sin \phi \\ \cos \theta & -\sin \theta \end{pmatrix}$$

If we now assume that the tensor \mathbf{T} is symmetric (e.g. stress tensor, strain rate tensor), then there are only 6 independent terms.

2.4 A continuum mechanics primer

continuum_mechanics.tex

Contains contributions by W. Spakman - Continuum mechanics course syllabus

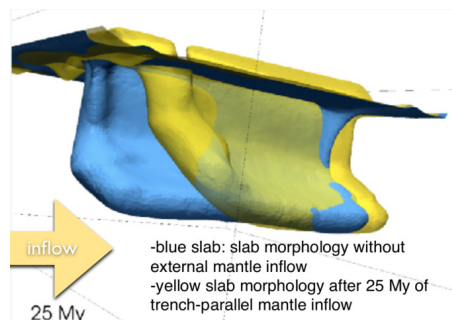
2.4.1 Forces

In continuum mechanics we make a distinction between two broad classes of forces:

- Body forces defined as force per unit volume (N m^{-3}): gravity, electro-magnetic forces
- Tractions: Surface forces defined as force per unit surface area (N m^{-2}): Contact forces, elastic forces per unit area, internal flow friction, pressure, ...

A traction is the surface average of all atomic forces exerted by atoms on the one side on atoms on the other side of the surface. For real-Earth processes, internal tractions are ultimately caused by the body forces, usually gravity.

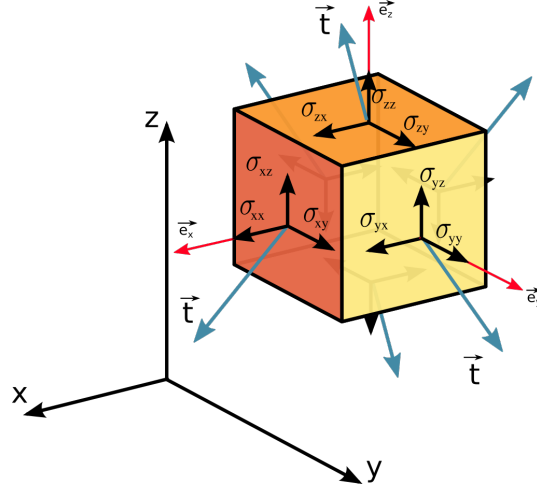
Existing mantle flow(i.e. flow that is forced elsewhere) can exert tractions (shear stresses) on the subducting slab or for instance at the base of lithosphere plates. In HPT-laboratory experiments external tractions (pressure, shear traction) are applied to a rock sample, which cause internal tractions to balance the exerted forces.



2.4.2 Stress tensor and tractions

The Cauchy stress tensor⁴ consists of nine components σ_{ij} that completely define the state of stress at a point inside a material. The tensor relates a unit-length direction vector \vec{n} to the so-called 'stress vector' (most commonly called 'traction') $\vec{t}(\vec{n})$ across an imaginary surface perpendicular to \vec{n} :

$$\vec{t}(\vec{n}) = \boldsymbol{\sigma} \cdot \vec{n}$$



Modified from original file on Wikipedia⁵

With respect to an orthonormal basis $\{\vec{e}_x, \vec{e}_y, \vec{e}_z\}$, the Cauchy stress tensor is given by:

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{yx} & \sigma_{yy} & \sigma_{yz} \\ \sigma_{zx} & \sigma_{zy} & \sigma_{zz} \end{pmatrix} \quad (2.19)$$

The three diagonal elements are called normal stresses while the off-diagonal terms are called shear stresses.

One can easily prove (see for instance Section 3.3.6 of [493]) that the balance of angular momentum leads reduces to the statement that the Cauchy stress tensor is symmetric, i.e. $\boldsymbol{\sigma} = \boldsymbol{\sigma}^T$. Therefore, the stress state of the medium at any point and instant can be specified by only six independent parameters, rather than nine:

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} \end{pmatrix} \quad \text{or sometimes} \quad \boldsymbol{\sigma} = \begin{pmatrix} \sigma_x & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \sigma_y & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \sigma_z \end{pmatrix} \quad (2.20)$$

where the elements σ_x , σ_y , σ_z are called the orthogonal normal stresses (relative to the chosen coordinate system), and τ_{xy} , τ_{xz} , τ_{yz} the orthogonal shear stresses. The left form is preferred in this document. As seen above, the SI units of both stress tensor and traction are N m^{-2} .

specify the underlying assumptions in what follows

⁴https://en.wikipedia.org/wiki/Cauchy_stress_tensor

⁵https://commons.wikimedia.org/wiki/File:Components_stress_tensor_cartesian.svg

In Cylindrical coordinates the stress tensor components are given by:

$$\sigma_{rr} = -p + 2\eta \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.21)$$

$$\sigma_{\theta\theta} = -p + 2\eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} + \frac{\mathbf{v}_r}{r} \right) \quad (2.22)$$

$$\sigma_{zz} = -p + 2\eta \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.23)$$

$$\sigma_{r\theta} = \eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial r} - \frac{\mathbf{v}_\theta}{r} \right) \quad (2.24)$$

$$\sigma_{rz} = \eta \left(\frac{\partial \mathbf{v}_r}{\partial z} + \frac{\partial \mathbf{v}_z}{\partial r} \right) \quad (2.25)$$

$$\sigma_{\theta z} = \eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_z}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial z} \right) \quad (2.26)$$

In Spherical coordinates the stress tensor components are given by:

$$\sigma_{rr} = -p + 2\eta \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.27)$$

$$\sigma_{\theta\theta} = -p + 2\eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} + \frac{\mathbf{v}_r}{r} \right) \quad (2.28)$$

$$\sigma_{\phi\phi} = -p + 2\eta \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} + \frac{\mathbf{v}_r}{r} + \frac{\mathbf{v}_\theta \cot \theta}{r} \right) \quad (2.29)$$

$$\sigma_{r\theta} = \eta \left(r \frac{\partial}{\partial r} \frac{\mathbf{v}_\theta}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \quad (2.30)$$

$$\sigma_{r\phi} = \eta \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_r}{\partial \phi} + r \frac{\partial}{\partial r} \frac{\mathbf{v}_\phi}{r} \right) \quad (2.31)$$

$$\sigma_{\theta\phi} = \eta \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\theta}{\partial \phi} + \frac{\sin \theta}{r} \frac{\partial}{\partial \theta} \frac{\mathbf{v}_\phi}{\sin \theta} \right) \quad (2.32)$$

2.4.3 Strain rate and spin tensor

The velocity gradient \mathbf{L} is given in Cartesian coordinates by:

$$\mathbf{L}(\vec{\mathbf{v}}) = \vec{\nabla} \vec{\mathbf{v}} = \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial v}{\partial x} & \frac{\partial w}{\partial x} \\ \frac{\partial u}{\partial y} & \frac{\partial v}{\partial y} & \frac{\partial w}{\partial y} \\ \frac{\partial u}{\partial z} & \frac{\partial v}{\partial z} & \frac{\partial w}{\partial z} \end{pmatrix} \quad (2.33)$$

It can be decomposed into its symmetric and skew-symmetric parts according to:

$$\vec{\nabla} \vec{\mathbf{v}} = (\vec{\nabla} \vec{\mathbf{v}})^s + (\vec{\nabla} \vec{\mathbf{v}})^w = \dot{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) + \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}) \quad (2.34)$$

The symmetric part is called the strain rate (or rate of deformation)⁶:

$$\dot{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) = \frac{1}{2} \left(\vec{\nabla} \vec{\mathbf{v}} + (\vec{\nabla} \vec{\mathbf{v}})^T \right) \quad (2.35)$$

⁶Note that often the dot is omitted and for example the ASPECT manual uses the $\boldsymbol{\varepsilon}$ notation.

The skew-symmetric tensor is called spin tensor (or vorticity tensor):

$$\dot{\omega}(\vec{v}) = \frac{1}{2} \left(\vec{\nabla} \vec{v} - (\vec{\nabla} \vec{v})^T \right) \quad (2.36)$$

Remark. In the mathematical literature a different notation for the strain rate tensor is often used, i.e. $\mathbf{D}(\vec{v})$ - or simply \mathbf{D} , such as for instance in Fullsack (1995) [426].

2.5 Viscous Newtonian rheology

physics.tex

The relationship between velocity-related stresses and velocity derivatives is such that the total stress tensor has the form [74]

$$\boldsymbol{\sigma} = -p\mathbf{1} + \mathbf{A} : \dot{\boldsymbol{\epsilon}}(\vec{v}) \quad (2.37)$$

where p is the thermodynamic pressure which is a function of the density ρ and the temperature T (an equation of state is then needed) and \mathbf{A} is the fourth-rank stiffness tensor.

Since both the stress and the strain tensors are symmetric and for isotropic fluids we have (see Malvern [831])

$$\mathbf{A} : \dot{\boldsymbol{\epsilon}}(\vec{v}) = \lambda(\vec{\nabla} \cdot \vec{v})\mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}(\vec{v}) \quad (2.38)$$

where λ is the bulk viscosity and η is the dynamic viscosity⁷. The stress tensor is then

$$\boldsymbol{\sigma} = (-p + \lambda(\vec{\nabla} \cdot \vec{v}))\mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}(\vec{v}) \quad (2.39)$$

By writing

$$\dot{\boldsymbol{\epsilon}}(\vec{v}) = \frac{1}{3}\text{tr}(\dot{\boldsymbol{\epsilon}}(\vec{v}))\mathbf{1} + \dot{\boldsymbol{\epsilon}}^d(\vec{v}) = \frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} + \dot{\boldsymbol{\epsilon}}^d(\vec{v})$$

where $\dot{\boldsymbol{\epsilon}}^d(\vec{v})$ is the deviatoric strain rate tensor and (in Cartesian coordinates)

$$\vec{\nabla} \cdot \vec{v} = \text{div}(\vec{v}) = \text{tr}(\dot{\boldsymbol{\epsilon}}(\vec{v})) = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} \quad (2.40)$$

where tr is the trace operator, we arrive at

$$\boldsymbol{\sigma} = (-p + \lambda(\vec{\nabla} \cdot \vec{v}))\mathbf{1} + 2\eta \left[\frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} + \dot{\boldsymbol{\epsilon}}^d \right] \quad (2.41)$$

$$= \left[-p + \left(\lambda + \frac{2}{3}\eta \right) (\vec{\nabla} \cdot \vec{v}) \right] \mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}^d \quad (2.42)$$

Introducing the second viscosity $\zeta = \lambda + \frac{2}{3}\eta$:

$$\boldsymbol{\sigma} = \left[-p + \zeta(\vec{\nabla} \cdot \vec{v}) \right] \mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}^d \quad (2.43)$$

$$= -\bar{p}\mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}^d \quad (2.44)$$

The effect of the volume viscosity ζ is that the mechanical pressure \bar{p} is not equivalent to the thermodynamic pressure p

$$\bar{p} = p - \zeta(\vec{\nabla} \cdot \vec{v}) \quad (2.45)$$

In other words: the isotropic average of the total stress is *not* the pressure term! This difference is usually neglected (and it is safe to do so, see [74, section 7.02.3.2.2]) by explicitly assuming $\zeta = 0$ (also called the Stokes assumption [1140, p256]), so that one can then refer to pressure as a single well-defined value. Note that in the case of an incompressible Newtonian Fluid, the strain rate tensor is deviatoric ($\text{tr}(\dot{\boldsymbol{\epsilon}}(\vec{v})) = \text{div}(\vec{v}) = 0$) and the above considerations vanish.

Finally, for both compressible and incompressible flow, the stress tensor becomes simply

⁷also sometimes called shear viscosity

$$\boldsymbol{\sigma} = -p\mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}^d(\vec{v}) = -p\mathbf{1} + \boldsymbol{\tau} \quad (2.46)$$

where $\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}^d(\vec{v})$ is the deviatoric stress tensor.

Remark. On page 256 of Schubert, Turcotte and Olson [1140], equation 6.5.3, the authors write $\tau_{ii}/3 = k_B e_{ii}$ while stating that τ is deviatoric in equation 6.4.2. This is an obvious conflict of notations.

2.6 The heat transport equation - energy conservation equation

physics.tex

2.7 The momentum conservation equations

physics.tex

As explained in Section 2.11, in Earth science applications the Navier-Stokes equations reduce to the Stokes equation:

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} = \vec{0} \quad (2.47)$$

Since

$$\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau} \quad (2.48)$$

it also writes

$$-\vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\tau} + \rho \vec{g} = \vec{0} \quad (2.49)$$

Using the relationship $\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}^d(\vec{v})$ we arrive at

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta\dot{\boldsymbol{\epsilon}}^d(\vec{v})) + \rho \vec{g} = \vec{0} \quad (2.50)$$

The divergence of a tensor field in cylindrical coordinates (r, θ, z) has been obtained in Section 2.3.3. The equations of motion (2.47) becomes⁸

$$\frac{\partial \sigma_{rr}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{r\theta}}{\partial \theta} + \frac{\partial \sigma_{rz}}{\partial z} + \frac{1}{r}(\sigma_{rr} - \sigma_{\theta\theta}) + \rho g_r = 0 \quad (2.51)$$

$$\frac{\partial \sigma_{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{\theta\theta}}{\partial \theta} + \frac{\partial \sigma_{\theta z}}{\partial z} + \frac{2}{r} \sigma_{r\theta} + \rho g_\theta = 0 \quad (2.52)$$

$$\frac{\partial \sigma_{rz}}{\partial r} + \frac{1}{r} \frac{\partial \sigma_{\theta z}}{\partial \theta} + \frac{\partial \sigma_{zz}}{\partial z} + \frac{1}{r} \sigma_{rz} + \rho g_z = 0 \quad (2.53)$$

2.8 The mass conservation equations

physics.tex

⁸https://en.wikipedia.org/wiki/Linear_elasticity

The mass conservation equation (often called continuity equation) is given by

$$\frac{D\rho}{Dt} + \rho \vec{\nabla} \cdot \vec{v} = 0$$

or, since

$$\frac{D\rho}{Dt} = \frac{\partial \rho}{\partial t} + \vec{v} \cdot \vec{\nabla} \rho$$

then

$$\frac{D\rho}{Dt} + \rho \vec{\nabla} \cdot \vec{v} = \frac{\partial \rho}{\partial t} + \vec{v} \cdot \vec{\nabla} \rho + \rho \vec{\nabla} \cdot \vec{v} = 0$$

and finally:

$$\frac{\partial \rho}{\partial t} + \vec{\nabla} \cdot (\rho \vec{v}) = 0 \quad (2.54)$$

In the case of an incompressible flow, then $\partial \rho / \partial t = 0$ and $\vec{\nabla} \rho = 0$, i.e. $D\rho/Dt = 0$ and the remaining equation is simply:

$$\vec{\nabla} \cdot \vec{v} = 0$$

A vector field that is divergence-free is also called solenoidal⁹.

In cylindrical coordinates (r, θ, ϕ) the continuity equation for an incompressible fluid is :

$$\frac{1}{r} \frac{\partial}{\partial r}(r v_r) + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} + \frac{\partial v_z}{\partial z} = 0$$

In spherical coordinates (r, θ, ϕ) the continuity equation for an incompressible fluid is :

$$\frac{1}{r^2} \frac{\partial}{\partial r}(r^2 v_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta}(v_\theta \sin \theta) + \frac{1}{r \sin \theta} \frac{\partial v_\phi}{\partial \phi} = 0 \quad (2.55)$$

2.9 The equations in Aspect manual

physics.tex

The following is lifted off the ASPECT manual. We focus on the system of equations in a $d = 2$ - or $d = 3$ -dimensional domain Ω that describes the motion of a highly viscous fluid driven by differences in the gravitational force due to a density that depends on the temperature. In the following, we largely follow the exposition of this material in Schubert, Turcotte and Olson [1140].

Specifically, we consider the following set of equations for velocity \vec{v} , pressure p and temperature T :

$$-\vec{\nabla} \cdot \left[2\eta \left(\dot{\epsilon}(\vec{v}) - \frac{1}{3}(\vec{\nabla} \cdot \vec{v}) \mathbf{1} \right) \right] + \vec{\nabla} p = \rho \vec{g} \quad \text{in } \Omega, \quad (2.56)$$

$$\vec{\nabla} \cdot (\rho \vec{v}) = 0 \quad \text{in } \Omega, \quad (2.57)$$

$$\begin{aligned} \rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T &= \rho H \\ &+ 2\eta \left(\dot{\epsilon}(\vec{v}) - \frac{1}{3}(\vec{\nabla} \cdot \vec{v}) \mathbf{1} \right) : \left(\dot{\epsilon}(\vec{v}) - \frac{1}{3}(\vec{\nabla} \cdot \vec{v}) \mathbf{1} \right) \\ &+ \alpha T \left(\vec{v} \cdot \vec{\nabla} p \right) \end{aligned} \quad \text{in } \Omega, \quad (2.58)$$

⁹https://en.wikipedia.org/wiki/Solenoidal_vector_field

where $\dot{\epsilon}(\vec{v}) = \frac{1}{2}(\vec{\nabla}\vec{v} + \vec{\nabla}\vec{v}^T)$ is the symmetric gradient of the velocity (often called the *strain rate* tensor).

In this set of equations, (2.56) and (2.57) represent the compressible Stokes equations in which $\vec{v} = \vec{v}(\mathbf{x}, t)$ is the velocity field and $p = p(\mathbf{x}, t)$ the pressure field. Both fields depend on space \mathbf{x} and time t . Fluid flow is driven by the gravity force that acts on the fluid and that is proportional to both the density of the fluid and the strength of the gravitational pull.

Coupled to this Stokes system is equation (2.58) for the temperature field $T = T(\mathbf{x}, t)$ that contains heat conduction terms as well as advection with the flow velocity \mathbf{v} . The right hand side terms of this equation correspond to

- internal heat production for example due to radioactive decay;
- friction (shear) heating;
- adiabatic compression of material;

In order to arrive at the set of equations in the ASPECT manual we need to

- neglect the $\partial p / \partial t$ _____
- neglect the $\partial \rho / \partial t$ in (2.54).

wrong
rephrase

from equations above. A partial answer is given in the next section.

Also, their definition of the shear heating term Φ is:

$$\Phi = k_B(\vec{\nabla} \cdot \vec{v})^2 + 2\eta\dot{\epsilon}^d : \dot{\epsilon}^d$$

For many fluids the bulk viscosity k_B is very small and is often taken to be zero, an assumption known as the Stokes assumption: $k_B = \lambda + 2\eta/3 = 0$. Note that η is the dynamic viscosity and λ the second viscosity. Also,

$$\boldsymbol{\tau} = 2\eta\dot{\epsilon} + \lambda(\nabla \cdot \vec{v})\mathbf{1}$$

but since $k_B = \lambda + 2\eta/3 = 0$, then $\lambda = -2\eta/3$ so

$$\boldsymbol{\tau} = 2\eta\dot{\epsilon} - \frac{2}{3}\eta(\nabla \cdot \vec{v})\mathbf{1} = 2\eta\dot{\epsilon}^d$$

2.10 Equations for thermal convection in an anelastic, compressible, self-gravitating spherical mantle

physics.tex

What follows is borrowed from Section 2.1 of Glišović *et al.* (2012) [468]. We start from the conservation mass, momentum and energy equations (the full Navier-Stokes equations):

$$\frac{\partial \rho}{\partial t} + \vec{\nabla} \cdot (\rho \vec{v}) = 0 \quad (2.59)$$

$$\rho \frac{D\vec{v}}{Dt} = \vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} \quad (2.60)$$

$$\rho C_p \frac{DT}{Dt} = \vec{\nabla} \cdot k \vec{\nabla} T + \alpha T \frac{Dp}{Dt} + \Phi + Q \quad (2.61)$$

In solving for the mantle flow field that satisfies the equation of momentum conservation, we incorporate all effects arising from self-gravitation and we must therefore explicitly consider the 3-D variation of gravity throughout Earth's interior. The gravitational acceleration is written as

$$\vec{g} = \vec{\nabla} \phi$$

where ϕ is Earth's gravitational potential which satisfies Poisson's equation

$$\Delta \phi = -4\pi \mathcal{G} \rho$$

The gravitational potential is expressed as

$$\phi = \phi_0(r) + \phi_1(r, \theta, \phi)$$

where the subscript 0 denotes a hydrostatic reference state, in which the structure of the mantle (density, gravity, pressure, temperature) varies with radius alone and the subscript 1 denotes all 3D perturbations arising from the thermal convection process. This decomposition makes sense in the context of a perfect sphere.

The total perturbed density and pressure fields in the mantle may similarly be expressed as

$$\rho = \rho_0(r) + \rho_1(r, \theta, \phi)$$

$$p = p_0(r) + p_1(r, \theta, \phi)$$

The equation of state relates the density perturbations to the temperature and pressure perturbations as follows

$$\rho_1 = \rho_0 [1 - \alpha(T - T_0(r)) + K_T^{-1}(p - p_0(r))]$$

where K_T is the bulk modulus and the term $T_0(r)$ represents the horizontally averaged temperature (i.e. the geotherm) which varies with radius only. The effects of compressibility on the density are found to be at least two orders of magnitude smaller than the effects of temperature variations. Therefore, the last term of this equation is often neglected. Note that this expression is a first order expansion of any Equation of State.

Also, this equation can be misleading if one forgets that the parameters α and K_T cannot be constant but must be related through Maxwell relations (for example, their definitions

$$\alpha = \frac{1}{V} \left(\frac{\partial V}{\partial T} \right)_P = -\frac{1}{\rho} \left(\frac{\partial \rho}{\partial T} \right)_P$$

and

$$K_T = -V \left(\frac{\partial P}{\partial V} \right)_T = \rho \left(\frac{\partial P}{\partial \rho} \right)_T$$

imply that

$$\frac{\partial(\alpha\rho)}{\partial P} = -\frac{\partial(\rho/K_T)}{\partial T}$$

Some models can be found in the geophysical literature in which assumptions made inconsistently about thermodynamic parameters (either constant or depth-dependent) violate the Maxwell rules.

Important simplifications are made assuming the anelastic-liquid approximation (e.g. Jarvis & McKenzie (1980) [637], Solheim & Peltier (1990) [1178]). This approximation is justified because the velocities associated with mantle convection are very slow compared to the local sound speed and hence acoustic waves cannot be generated by the slow changes in the mantle pressure field. We therefore neglect the time derivative of density, thereby eliminating sound waves:

$$\frac{\partial \rho}{\partial t} \simeq 0$$

For the same reason, the pressure distribution may be considered (to first-order accuracy) as the pressure of a fluid in hydrostatic equilibrium which yields

$$\frac{Dp}{Dt} = \frac{\partial p}{\partial t} + \vec{v} \cdot \vec{\nabla} p \simeq -u_r \rho_0(r) g_0(r)$$

The equations are then rewritten in terms of dimensionless variables according to the relations:

$$r' = \frac{r}{d} \tag{2.62}$$

$$\mathbf{v}' = \frac{\mathbf{v}}{U} \tag{2.63}$$

$$t' = \frac{U}{d/t} \tag{2.64}$$

$$T' = \frac{T}{\Delta T} \tag{2.65}$$

$$\rho' = \frac{\rho}{\rho_{0s}} \tag{2.66}$$

$$g' = \frac{g}{g_{0s}} \tag{2.67}$$

$$\phi' = \frac{\phi}{g_{0s}d} \tag{2.68}$$

$$\alpha' = \frac{\alpha}{\alpha_s} \tag{2.69}$$

$$p' = \frac{p}{\alpha_s \Delta T \rho_{0s} g_{0s} d} \tag{2.70}$$

$$\tau_{ij} = \frac{\tau_{ij}}{\alpha_s \Delta T \rho_{0s} g_{0s} d} \tag{2.71}$$

$$\eta' = \frac{\eta}{\eta_s} \tag{2.72}$$

$$k' = \frac{k}{k_s} \tag{2.73}$$

$$Q' = \frac{Q d^2}{k_s \Delta T} \tag{2.74}$$

$$U = \frac{\rho_{0s} g_{0s} \alpha_s \Delta T d^2}{\eta_s} \tag{2.75}$$

in which the primes represent the dimensionless variables, the subscript s means that we consider the surface value of the variable to which it is applied. The length scale d and temperature scale ΔT are respectively the depth of the mantle and the difference of temperature between the bottom and the top of the mantle.

Often one deals with dimensionless variables and the primes are dropped for notational convenience (this is the case in what follows).

It is a tedious but trivial exercise to show that the dimensionless equation of conservation of momentum is then written as follows:

$$\rho \frac{\text{Ra}_s}{\text{Pr}_s} \frac{D\vec{v}}{Dt} = \frac{\rho}{\alpha_s \Delta T} \vec{\nabla} \phi - \vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\tau}$$

in which we introduce the surface Rayleigh Ra_s and Prandtl Pr_s numbers defined, respectively, by

$$\text{Ra}_s = \frac{\rho_0^2 C_p g_0 \alpha_s \Delta T d^3}{k_s \eta_s} \quad \text{Pr}_s = \frac{\eta_s C_p}{k_s}$$

Because of the very high viscosity of mantle rocks, the left-hand term is smaller than the other terms by several orders of magnitude and may therefore be neglected. This important simplification is called the infinite Prandtl number approximation.

The equation of energy conservation may also be rewritten in terms of the surface Rayleigh number, as follows

$$\frac{DT}{Dt} = \frac{1}{\rho \text{Ra}_s} \left(\vec{\nabla} \cdot k \vec{\nabla} T + Q \right) + \frac{Di}{\rho} \left(-\alpha T \frac{Dp}{Dt} + \Phi \right)$$

where Di is the dissipation number (see Peltier (1972) [986]) which measures the importance of compression work and frictional heating, and it is defined as

$$Di = \frac{\alpha_s g_0 d}{C_p}$$

Di also measures the ratio of the depth of mantle convection (d) to the adiabatic scale height ($C_p/\alpha g_0$) and for whole-mantle convection is close to order 1 (see Jarvis & McKenzie (1980) [637]).

After simplifications, the dimensionless system of governing equations is written as

$$\vec{\nabla} \cdot (\rho_0 \vec{v}) = 0 \quad (2.76)$$

$$\frac{\rho}{\alpha_s \Delta T} \vec{\nabla} \phi - \vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\tau} = \vec{0} \quad (2.77)$$

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \frac{1}{\rho_0 \text{Ra}_s} \left(\vec{\nabla} \cdot k \vec{\nabla} T + Q \right) + \frac{Di}{\rho_0} (-\alpha T \rho_0 g_0 u_r + \Phi) \quad (2.78)$$

with

$$\Delta \phi = -4\pi \mathcal{G} \rho$$

$$\rho_1 = \rho_0 (1 - \alpha(T - T_0(r)))$$

VERIFY all this !

2.11 Non-dimensionalisation of the Navier-Stokes equations

2.11.1 Approach # 1 - isothermal flow

We define (see for instance Massimi *et al.* (2006) [839]) four reference quantities which are relevant for geodynamics¹⁰:

- a reference viscosity value $\underline{\eta} = 10^{20}$ Pa s
- a reference mass density $\underline{\rho} = 1000$ kg m⁻³
- a reference time $\underline{t} = 1\text{Myr} \simeq 3.15 \cdot 10^{13}$ s
- a reference length $\underline{l} = 1000$ m
- a reference gravity $\underline{g} = 9.81$ m s⁻²

It follows that a reference pressure can be obtained:

$$\underline{p} = \underline{\rho g l} = 9.81 \cdot 10^6 \text{ Pa}$$

Note that there is unfortunately no natural selection for the pressure scale. We could also have used $\underline{p} = \underline{\rho \mathbf{v}^2}$ where dynamic effects are dominant i.e. high velocity flows, or $\underline{p} = \underline{\eta \mathbf{v} / l}$ where viscous effects are dominant i.e. creeping flows (which is the case in geodynamics). The definition of a reference velocity is more straightforward:

$$\underline{\mathbf{v}} = \frac{\underline{l}}{\underline{t}} = 1 \text{ mm yr}^{-1}$$

We define dimensionless variables through:

$$\textcolor{teal}{x} = \frac{x}{\underline{l}} \quad \textcolor{teal}{y} = \frac{y}{\underline{l}} \quad \textcolor{teal}{z} = \frac{z}{\underline{l}} \quad \textcolor{teal}{\vec{v}'} = \frac{\vec{v}}{\underline{\mathbf{v}}} \quad \textcolor{teal}{t} = \frac{t}{\underline{t}} \quad \textcolor{teal}{\eta} = \frac{\eta}{\underline{\eta}} \quad \textcolor{teal}{g} = \frac{g}{\underline{g}}$$

where the teal color indicates dimensionless values.

Consequently, time and space derivatives will be rescaled as follows:

$$\textcolor{teal}{\vec{\nabla}} = \underline{l} \vec{\nabla} \quad \textcolor{teal}{\partial_t} = \underline{t} \partial_t$$

Using this scaling relations the Navier-Stokes equation become:

$$\frac{\underline{\rho l}}{\underline{t}^2} \frac{\partial \vec{v}}{\partial \textcolor{teal}{t}} + \frac{\underline{\rho l}}{\underline{t}^2} (\vec{v} \cdot \textcolor{teal}{\vec{\nabla}}) \vec{v} = -\underline{\rho g} \vec{\nabla} p + \frac{\underline{\eta}}{\underline{l t}} \textcolor{teal}{\vec{\nabla}} \cdot \eta (\textcolor{teal}{\vec{\nabla}} \vec{v} + \textcolor{teal}{\vec{\nabla}} \vec{v}^T) + \rho \vec{g}$$

I make $\textcolor{teal}{\rho} = \rho / \underline{\rho}$ appear in the left hand side:

$$\frac{\underline{\rho l}}{\underline{t}^2} \textcolor{teal}{\rho} \frac{\partial \vec{v}}{\partial \textcolor{teal}{t}} + \frac{\underline{\rho l}}{\underline{t}^2} \rho (\vec{v} \cdot \textcolor{teal}{\vec{\nabla}}) \vec{v} = -\underline{\rho g} \vec{\nabla} p + \frac{\underline{\eta}}{\underline{l t}} \textcolor{teal}{\vec{\nabla}} \cdot \eta (\textcolor{teal}{\vec{\nabla}} \vec{v} + \textcolor{teal}{\vec{\nabla}} \vec{v}^T) + \rho \vec{g}$$

which we can divide by $\underline{\rho l} / \underline{t}^2$ to obtain:

$$\textcolor{teal}{\rho} \left(\frac{\partial \vec{v}}{\partial \textcolor{teal}{t}} + (\vec{v} \cdot \textcolor{teal}{\vec{\nabla}}) \vec{v} \right) = -\frac{\underline{g t^2}}{\underline{l}} \vec{\nabla} p + \frac{\underline{\eta t}}{\underline{\rho l^2}} \textcolor{teal}{\vec{\nabla}} \cdot \eta (\textcolor{teal}{\vec{\nabla}} \vec{v} + \textcolor{teal}{\vec{\nabla}} \vec{v}^T) + \frac{\underline{t^2}}{\underline{l}} \rho \vec{g}$$

One can recognise in this equation the Reynolds and Froude non-dimensional numbers (the ratio between the inertial and viscous forces, and the ratio between buoyancy and inertial forces respectively).

$$Re = \frac{\underline{\rho l^2}}{\underline{\eta t}} \quad Fr = \frac{\underline{l}}{\underline{g t^2}}$$

¹⁰Note that in the paper the authors conflate ρ and $\tilde{\rho}$ which prevents them from non-dimensionalising all terms as we do here.

From this we conclude that inertial forces in the Earth's mantle are small compared to viscous forces. We can then write:

$$\rho \left(\frac{\partial \vec{v}}{\partial t} + (\vec{v} \cdot \nabla) \vec{v} \right) = -\frac{1}{Fr} \nabla p + \frac{1}{Re} \nabla \cdot \eta (\nabla \vec{v} + \nabla \vec{v}^T) + \frac{1}{Fr} \vec{g}$$

In our case, given the definitions taken above, we have:

$$Re \simeq 3.174 \cdot 10^{-24} \quad Fr \simeq 1.027 \cdot 10^{-25}$$

so that the inertial terms can be dropped from the momentum equation (thereby yielding the dimensionless Stokes equations):

$$\nabla \cdot \eta (\nabla \vec{v} + \nabla \vec{v}^T) - \frac{Re}{Fr} \nabla p + \frac{Re}{Fr} \vec{g} = 0$$

Note that in our case $Re/Fr \simeq 30.5$.

2.11.2 Approach # 2 - Temperature dependent

dimensionless_equations2.tex.tex

Let us now consider a box heated from below and cooled from above. We define 4 fundamental reference quantities:

- a length L_{ref} (m), (L)
- a temperature T_{ref} (K), (θ)
- a viscosity η_{ref} (Pa s), ($ML^{-1}T^{-1}$)
- a thermal diffusion coefficient κ_{ref} ($m^2 s^{-1}$), ($L^2 T^{-1}$)

From these reference quantities one can form secondary ones, such as

- a time $t_{ref} = L_{ref}^2 / \kappa_{ref}$ (aka the diffusion time)
- a velocity $v_{ref} = L_{ref} / t_{ref} = \kappa_{ref} / L_{ref}$
- an acceleration $g_{ref} = v_{ref} / t_{ref} = \kappa_{ref}^2 / L_{ref}^3$
- a strain rate $\dot{\epsilon}_{ref} = t_{ref}^{-1} = \kappa_{ref} / L_{ref}^2$
- a pressure $p_{ref} = \eta_{ref} \dot{\epsilon}_{ref} = \eta_{ref} t_{ref}^{-1}$
- a reference density $\rho_{ref} = \eta_{ref} L_{ref} t_{ref} / L_{ref}^3 = \eta_{ref} L_{ref}^{-2} t_{ref}$
- a reference mass $M_{ref} = \eta_{ref} L_{ref} t_{ref}$
- a reference energy $E_{ref} = \eta_{ref} L_{ref} t_{ref} \frac{L_{ref}^2}{t_{ref}^2} = \eta_{ref} \frac{L_{ref}^3}{t_{ref}}$
- a reference heat conductivity¹¹ $k_{ref} = E_{ref} / t_{ref} / L_{ref} / T_{ref} = \eta_{ref} L_{ref}^2 / t_{ref}^2 / T_{ref}$
- a reference heat capacity¹² $C_{ref} = E_{ref} / M_{ref} / T_{ref} = L_{ref}^2 / t_{ref}^2 / T_{ref}$
- a reference heat production coefficient¹³ $H_{ref} = E_{ref} / t_{ref} / M_{ref} = \frac{L_{ref}^2}{t_{ref}^3}$
- a reference heat flux¹⁴ $q_{ref} = \eta_{ref} L_{ref} t_{ref}^{-2}$

¹¹Units: W/m/K

¹²Units: J/kg/K

¹³Units: W/kg

¹⁴Units: W m⁻², or kg s⁻³

We define **dimensionless** quantities as follows:

$$x = \frac{x}{L_{ref}} \quad \vec{v} = \frac{\vec{v}}{v_{ref}} \quad t = \frac{t}{t_{ref}} \quad \eta = \frac{\eta}{\eta_{ref}} \quad g = \frac{g}{g_{ref}} \quad k = \frac{k}{k_{ref}} \quad C_p = \frac{C_p}{C_{ref}} \quad \rho = \frac{\rho}{\rho_{ref}} \quad (2.79)$$

$$H = \frac{H}{H_{ref}} \quad \vec{\nabla} = L_{ref} \vec{\nabla} \quad \partial_t = t_{ref} \partial_t \quad T = \frac{T}{T_{ref}} \quad \dot{\epsilon} = \dot{\epsilon} t_{ref} \quad (2.80)$$

We start from the standard Navier-Stokes equation¹⁵

$$\rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}) + \rho \vec{g}$$

and assume that the density is temperature-dependent (Boussinesq approximation) so that

$$\rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}) + \rho_0(1 - \alpha T) \vec{g}$$

and remove the hydrostatic pressure (although we keep using p for simplicity, p is now the dynamic pressure):

$$\rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}) - \rho_0 \alpha T \vec{g}$$

We divide this equation by $p_{ref} = \eta_{ref} \dot{\epsilon}_{ref}$:

$$\frac{1}{\eta_{ref} \dot{\epsilon}_{ref}} \rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2 \frac{\eta}{\eta_{ref}} \frac{\dot{\epsilon}}{\dot{\epsilon}_{ref}} - \frac{\rho_0 \alpha T \vec{g}}{\eta_{ref} \dot{\epsilon}_{ref}}$$

Let us call \vec{e} the positive vertical vector (\vec{e}_z in Cartesian coordinates, \vec{e}_r in spherical coordinates), then $\vec{g} = -g_0 \vec{e}$ and we can write (using $\dot{\epsilon}_{ref} = \kappa_{ref}/L_{ref}^2$)

$$\frac{1}{\eta_{ref} \dot{\epsilon}_{ref}} (\rho_{ref} \rho) \frac{D(v_{ref} \vec{v})}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \frac{\rho_0 \alpha T g_0}{\eta_{ref} (\kappa_{ref}/L_{ref}^2)} \vec{e}$$

Finally, dividing by L_{ref}^{-1} (i.e. multiplying by L_{ref}) yields

$$\frac{v_{ref} \rho_{ref} L_{ref}}{\eta_{ref} \dot{\epsilon}_{ref}} \rho \frac{D\vec{v}}{t_{ref} Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \frac{\rho_0 \alpha (T T_{ref}) g_0 L_{ref}^3}{\eta_{ref} \kappa_{ref}} \vec{e}$$

and finally (using $v_{ref} = L_{ref}/t_{ref}$)

$$\frac{\rho_{ref} \kappa_{ref}}{\eta_{ref}} \rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \frac{\rho_0 \alpha T_{ref} g_0 L_{ref}^3}{\eta_{ref} \kappa_{ref}} T \vec{e}$$

In the context of a system with a temperature difference ΔT between the bottom and top boundaries separated by a distance H , one would then take $T_{ref} = \Delta T$ and $L_{ref} = H$ so that the equation becomes:

$$\underbrace{\frac{\rho_{ref} \kappa_{ref}}{\eta_{ref}}}_{Pr^{-1}} \rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \underbrace{\frac{\rho_0 \alpha \Delta T g_0 H^3}{\eta_{ref} \kappa_{ref}}}_{Ra} T \vec{e}$$

and we obviously recover the classical definition of the Rayleigh number.

¹⁵https://en.wikipedia.org/wiki/Navier-Stokes_equations

$$\frac{1}{\text{Pr}} \rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \text{Ra} T \vec{e}$$

On the left side of the equation we recognize the (inverse of the) Prandtl number $\text{Pr} = \frac{\eta}{\rho\kappa}$. We can estimate the dimensionless number before the inertial term for Earth geodynamics:

$$\text{Pr} \simeq \frac{10^{20-23}}{3000 \cdot 10^{-6}} \gg 10^{23}$$

Its inverse is then extremely small and this is why we neglect the inertial terms in mantle modelling.

Note that if the fluid is isoviscous, one can then set $\eta_{ref} = \eta = \eta_0$ and then $\eta = 1$

Turning now to the continuity equation $\vec{\nabla} \cdot \vec{v} = 0$, it is trivial to show that $\vec{\nabla} \cdot \vec{v} = 0$. Finally, starting from the simple heat transport equation:

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T$$

We divide each side by T_{ref} so that

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T$$

We now divide each side by the reference velocity v_{ref} and we obtain

$$\frac{L_{ref}}{\kappa_{ref}} \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \frac{L_{ref}}{\kappa_{ref}} \kappa \Delta T$$

We multiply each side by L_{ref} and we finally get

$$\frac{L_{ref}^2}{\kappa_{ref}} \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T$$

and finally

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T$$

The set of dimensionless equations is then:

$$-\vec{\nabla} p + \vec{\nabla} \cdot 2\eta \dot{\epsilon} + \text{Ra} T \vec{e} = \vec{0} \quad (2.81)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (2.82)$$

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T \quad (2.83)$$

Looking now at the Extended Boussinesq Approximation (EBA), we have to consider two additional terms in the energy equation:

- the shear heating Φ (See Eq.(11.84)) $\Phi = 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d$
- the adiabatic heating $\alpha T \vec{v} \cdot \vec{\nabla} p$

We start this time from

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T = \rho H + 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \alpha T \vec{v} \cdot \vec{\nabla} p$$

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T = \rho H + 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \alpha T \vec{v} \cdot \vec{\nabla} p$$

$$\rho C_p \frac{T_{ref}}{t_{ref}} \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \frac{T_{ref}}{L_{ref}^2} \vec{\nabla} \cdot k \vec{\nabla} T = \rho_{ref} \rho H + \frac{\eta_{ref}}{t_{ref}^2} 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \frac{p_{ref}}{t_{ref}} \alpha T \vec{v} \cdot \vec{\nabla} p$$

we then use $p_{ref} = \eta_{ref} t_{ref}^{-1}$ and $\rho_{ref} = \eta_{ref} L_{ref}^{-2} t_{ref}$

$$\rho_{ref} C_{ref} \frac{T_{ref}}{t_{ref}} \rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \frac{T_{ref} k_{ref}}{L_{ref}^2} \vec{\nabla} \cdot k \vec{\nabla} T = \frac{\eta_{ref}}{L_{ref}^2} t_{ref} \frac{L_{ref}^2}{t_{ref}^3} \rho H + \frac{\eta_{ref}}{t_{ref}^2} 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \frac{\eta_{ref}}{t_{ref}^2} \alpha T \vec{v} \cdot \vec{\nabla} p$$

or, multiplying all by t_{ref}^2/η_{ref} :

$$\frac{t_{ref}^2}{\eta_{ref}} \rho_{ref} C_{ref} \frac{T_{ref}}{t_{ref}} \rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \frac{t_{ref}^2}{\eta_{ref}} \frac{T_{ref} k_{ref}}{L_{ref}^2} \vec{\nabla} \cdot k \vec{\nabla} T = \rho H + 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \alpha T \vec{v} \cdot \vec{\nabla} p$$

We then make use of $C_{ref} = L_{ref}^2/t_{ref}^2/T_{ref}$ and $k_{ref} = \eta_{ref} L_{ref}^2/t_{ref}^2/T_{ref}$ to arrive at

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T = \rho H + 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \alpha T \vec{v} \cdot \vec{\nabla} p$$

2.12 The Navier-Stokes equations in cylindrical coordinates

physics.tex

In cylindrical coordinates, (r, θ, z) , the continuity equation for an incompressible fluid is

$$\frac{1}{r} \frac{\partial}{\partial r}(r \mathbf{v}_r) + \frac{1}{r} \frac{\partial}{\partial \theta}(\mathbf{v}_\theta) + \frac{\partial \mathbf{v}_z}{\partial z} = 0 \quad (2.84)$$

or

$$\frac{\partial \mathbf{v}_r}{\partial r} + \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial}{\partial \theta}(\mathbf{v}_\theta) + \frac{\partial \mathbf{v}_z}{\partial z} = 0 \quad (2.85)$$

The Navier-Stokes equations of motion for an incompressible fluid with uniform viscosity are:

$$\begin{aligned} \rho \left(\frac{D \mathbf{v}_r}{Dt} - \frac{\mathbf{v}_\theta^2}{r} \right) &= -\frac{\partial p}{\partial r} + f_r + \eta \left(\Delta \mathbf{v}_r - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right) \\ \rho \left(\frac{D \mathbf{v}_\theta}{Dt} + \frac{\mathbf{v}_\theta \mathbf{v}_r}{r} \right) &= -\frac{1}{r} \frac{\partial p}{\partial \theta} + f_\theta + \eta \left(\Delta \mathbf{v}_\theta - \frac{\mathbf{v}_\theta}{r^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \\ \rho \frac{D \mathbf{v}_z}{Dt} &= -\frac{\partial p}{\partial z} + f_z + \eta \Delta \mathbf{v}_z \end{aligned} \quad (2.86)$$

where the Lagrangian or material derivative is

$$\frac{D}{Dt} = \frac{\partial}{\partial t} + \mathbf{v}_r \frac{\partial}{\partial r} + \frac{\mathbf{v}_\theta}{r} \frac{\partial}{\partial \theta} + \mathbf{v}_z \frac{\partial}{\partial z}$$

and the Laplacian operator is

$$\Delta = \frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} + \frac{\partial^2}{\partial z^2} \quad (2.87)$$

and for an incompressible, Newtonian fluid

$$\sigma_{rr} = -p + 2\eta \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.88)$$

$$\sigma_{\theta\theta} = -p + 2\eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} + \frac{\mathbf{v}_r}{r} \right) \quad (2.89)$$

$$\sigma_{zz} = -p + 2\eta \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.90)$$

$$\sigma_{rz} = \eta \left(\frac{\partial \mathbf{v}_r}{\partial z} + \frac{\partial \mathbf{v}_z}{\partial r} \right) \quad (2.91)$$

$$\sigma_{r\theta} = \eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial r} - \frac{\mathbf{v}_\theta}{r} \right) \quad (2.92)$$

$$\sigma_{\theta z} = \eta \left(\frac{1}{r} \frac{\partial \mathbf{v}_z}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial z} \right) \quad (2.93)$$

2.13 The Stokes equations in spherical coordinates

physics.tex

In spherical coordinates, (r, θ, ϕ) , the continuity equation for an incompressible fluid is

$$\frac{1}{r^2} \frac{\partial}{\partial r} (r^2 \mathbf{v}_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\mathbf{v}_\theta \sin \theta) + \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} = 0 \quad (2.94)$$

Concerning the momentum equation, we start from

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{f} = \vec{0} \quad (2.95)$$

The buoyancy force \vec{f} is nearly always given by $\vec{f} = \rho \vec{g} = -\rho g \vec{e}_r$ (with $g > 0$), i.e. $f_\phi = f_\theta = 0$, and then

$$-\vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\tau} - \rho g \vec{e}_r = \vec{0}$$

or,

$$\begin{aligned} -(\vec{\nabla} p)_r + (\vec{\nabla} \cdot \boldsymbol{\tau})_r &= \rho g \\ -(\vec{\nabla} p)_\theta + (\vec{\nabla} \cdot \boldsymbol{\tau})_\theta &= 0 \\ -(\vec{\nabla} p)_\phi + (\vec{\nabla} \cdot \boldsymbol{\tau})_\phi &= 0 \end{aligned}$$

The pressure gradient is simply given by:

$$\begin{aligned} (\vec{\nabla} p)_r &= \frac{\partial p}{\partial r} \\ (\vec{\nabla} p)_\theta &= \frac{1}{r} \frac{\partial p}{\partial \theta} \\ (\vec{\nabla} p)_\phi &= \frac{1}{r \sin \theta} \frac{\partial p}{\partial \phi} \end{aligned}$$

We now turn to the remaining three components of the divergence of deviatoric stress in spherical coordinates r, θ, ϕ , which are given by¹⁶

$$\begin{aligned} (\vec{\nabla} \cdot \boldsymbol{\tau})_r &= \frac{\partial \tau_{rr}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta r}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi r}}{\partial \phi} + \frac{2\tau_{rr} - \tau_{\theta\theta} - \tau_{\phi\phi}}{r} + \frac{\tau_{\theta r} \cot \theta}{r} \\ (\vec{\nabla} \cdot \boldsymbol{\tau})_\theta &= \frac{\partial \tau_{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta\theta}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi\theta}}{\partial \phi} + \frac{3\tau_{\theta r} + (\tau_{\theta\theta} - \tau_{\phi\phi}) \cot \theta}{r} \\ (\vec{\nabla} \cdot \boldsymbol{\tau})_\phi &= \frac{\partial \tau_{r\phi}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta\phi}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi\phi}}{\partial \phi} + \frac{3\tau_{r\phi} + 2\tau_{\phi\theta} \cot \theta}{r} \end{aligned} \quad (2.96)$$

And finally the momentum equation writes:

$$\begin{aligned} -\frac{\partial p}{\partial r} + \frac{\partial \tau_{rr}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta r}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi r}}{\partial \phi} + \frac{2\tau_{rr} - \tau_{\theta\theta} - \tau_{\phi\phi}}{r} + \frac{\tau_{\theta r} \cot \theta}{r} &= \rho g \\ -\frac{1}{r} \frac{\partial p}{\partial \theta} + \frac{\partial \tau_{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta\theta}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi\theta}}{\partial \phi} + \frac{3\tau_{\theta r} + (\tau_{\theta\theta} - \tau_{\phi\phi}) \cot \theta}{r} &= 0 \\ -\frac{1}{r \sin \theta} \frac{\partial p}{\partial \phi} + \frac{\partial \tau_{r\phi}}{\partial r} + \frac{1}{r} \frac{\partial \tau_{\theta\phi}}{\partial \theta} + \frac{1}{r \sin \theta} \frac{\partial \tau_{\phi\phi}}{\partial \phi} + \frac{3\tau_{r\phi} + 2\tau_{\phi\theta} \cot \theta}{r} &= 0 \end{aligned} \quad (2.97)$$

¹⁶Would be nice to have a ref here

The deviatoric stress tensor components are

$$\begin{aligned}\tau_{rr} &= 2\eta \left[\dot{\varepsilon}_{rr} - \frac{1}{3}(\dot{\varepsilon}_{rr} + \dot{\varepsilon}_{\theta\theta} + \dot{\varepsilon}_{\phi\phi}) \right] \\ \tau_{\theta\theta} &= 2\eta \left[\dot{\varepsilon}_{\theta\theta} - \frac{1}{3}(\dot{\varepsilon}_{rr} + \dot{\varepsilon}_{\theta\theta} + \dot{\varepsilon}_{\phi\phi}) \right] \\ \tau_{\phi\phi} &= 2\eta \left[\dot{\varepsilon}_{\phi\phi} - \frac{1}{3}(\dot{\varepsilon}_{rr} + \dot{\varepsilon}_{\theta\theta} + \dot{\varepsilon}_{\phi\phi}) \right] \\ \tau_{r\theta} &= 2\eta \dot{\varepsilon}_{r\theta} \end{aligned} \tag{2.98}$$

$$\tau_{r\phi} = 2\eta \dot{\varepsilon}_{r\phi} \tag{2.99}$$

$$\tau_{\theta\phi} = 2\eta \dot{\varepsilon}_{\theta\phi} \tag{2.100}$$

with

$$\begin{aligned}\dot{\varepsilon}_{rr} &= \frac{\partial \mathbf{v}_r}{\partial r} \\ \dot{\varepsilon}_{\theta\theta} &= \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \\ \dot{\varepsilon}_{\phi\phi} &= \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} + \frac{\mathbf{v}_r}{r} + \frac{\mathbf{v}_\theta \cot \theta}{r} \\ \dot{\varepsilon}_{\theta r} = \dot{\varepsilon}_{r\theta} &= \frac{1}{2} \left(r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\theta}{r} \right) + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \\ \dot{\varepsilon}_{\phi r} = \dot{\varepsilon}_{r\phi} &= \frac{1}{2} \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_r}{\partial \phi} + r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\phi}{r} \right) \right) \\ \dot{\varepsilon}_{\phi\theta} = \dot{\varepsilon}_{\theta\phi} &= \frac{1}{2} \left(\frac{\sin \theta}{r} \frac{\partial}{\partial \theta} \left(\frac{\mathbf{v}_\phi}{\sin \theta} \right) + \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\theta}{\partial \phi} \right)\end{aligned}$$

We go further by assuming the fluid to be incompressible (i.e. $\dot{\varepsilon}_{rr} + \dot{\varepsilon}_{\theta\theta} + \dot{\varepsilon}_{\phi\phi} = 0$) and then:

$$\begin{aligned}\tau_{rr} = 2\eta \dot{\varepsilon}_{rr} &= 2\eta \frac{\partial \mathbf{v}_r}{\partial r} \\ \tau_{\theta\theta} = 2\eta \dot{\varepsilon}_{\theta\theta} &= 2\eta \left(\frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right) \\ \tau_{\phi\phi} = 2\eta \dot{\varepsilon}_{\phi\phi} &= 2\eta \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} + \frac{\mathbf{v}_r}{r} + \frac{\mathbf{v}_\theta \cot \theta}{r} \right) \\ \tau_{r\theta} = 2\eta \dot{\varepsilon}_{r\theta} &= \eta \left(r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\theta}{r} \right) + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \\ \tau_{r\phi} = 2\eta \dot{\varepsilon}_{r\phi} &= \eta \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_r}{\partial \phi} + r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\phi}{r} \right) \right) \\ \tau_{\theta\phi} = 2\eta \dot{\varepsilon}_{\theta\phi} &= \eta \left(\frac{\sin \theta}{r} \frac{\partial}{\partial \theta} \left(\frac{\mathbf{v}_\phi}{\sin \theta} \right) + \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\theta}{\partial \phi} \right)\end{aligned}$$

Inserting these expressions in Eq. (2.97) is a cumbersome affair... Under the assumption that the

fluid is also isoviscous, we get¹⁷

$$\begin{aligned}
\rho g &= -\frac{\partial p}{\partial r} + \eta \left(\Delta \mathbf{v}_r - \frac{2\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} - \frac{2\mathbf{v}_\theta \cot \theta}{r^2} - \frac{2}{r^2 \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} \right) \\
0 &= -\frac{1}{r} \frac{\partial p}{\partial \theta} + \eta \left(\Delta \mathbf{v}_\theta + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2 \sin^2 \theta} - \frac{2 \cot \theta}{r^2 \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} \right) \\
0 &= -\frac{1}{r \sin \theta} \frac{\partial p}{\partial \phi} + \eta \left(\Delta \mathbf{v}_\phi + \frac{2}{r^2 \sin \theta} \frac{\partial \mathbf{v}_r}{\partial \phi} - \frac{\mathbf{v}_\phi}{r^2 \sin^2 \theta} + \frac{2 \cot \theta}{r^2 \sin \theta} \frac{\partial \mathbf{v}_\theta}{\partial \phi} \right)
\end{aligned} \tag{2.101}$$

and the Laplacian operator is

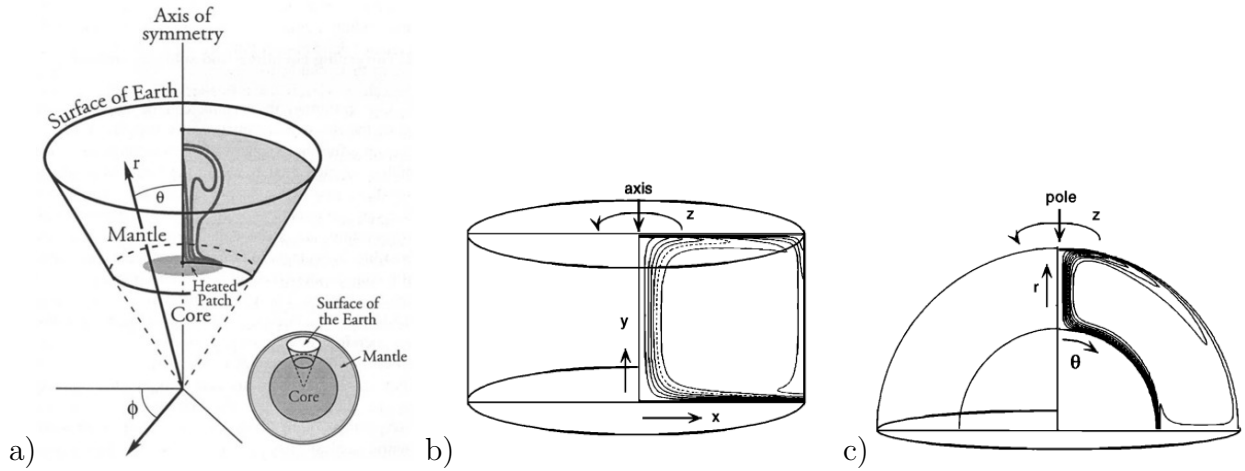
$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2}{\partial \phi^2}$$

Before being used these equations should be checked against multiple sources.

2.14 The equations for axisymmetric geometries

axisymmetric.eqs.tex

In what follows we are concerned with incompressible flow. In some cases the assumption can be made that the object we wish to study has an axisymmetric geometry, for example a plume:



a) Taken from Kellogg and King [691] (1997); b,c) Taken from Leitch, Steinbach, and Yuen [762] (1996).

Looking at the figure above we see that there are in fact two cases: axisymmetry in cylindrical coordinates (b) and axisymmetry in spherical coordinates (c).

As mentioned in Kellogg and King [691] (1997): "By imposing axisymmetry, we restrict the problem to two degrees of freedom, reducing the computational effort significantly over 3D calculations." However, [1052] (2004) also mention: "An important caveat of axisymmetric calculations is that there are no variations in the ϕ direction (i.e., there are no ϕ derivatives in the governing equations). Thus, as we get further from the pole, the results become increasingly less physical. Downwelling drips off the pole are actually downwelling doughnuts that follow the entire small circle. In a fully 3D calculation, this doughnut feature would in reality be a drip."

See Section 7.5.6 for the FE formulation of these equations.

¹⁷I have not thoroughly checked these equations yet

In cylindrical coordinates

The velocity vector is $\vec{\mathbf{v}} = (\mathbf{v}_r, \mathbf{v}_\theta, \mathbf{v}_z)$. Due to the symmetry we have $\mathbf{v}_\theta = 0$, $\partial_\theta \rightarrow 0$ and the Stokes equations then become ¹⁸

$$-\frac{\partial p}{\partial r} + \eta \left(\frac{1}{r} \frac{\partial}{\partial r} \left(r \frac{\partial \mathbf{v}_r}{\partial r} \right) + \frac{\partial^2 \mathbf{v}_r}{\partial z^2} - \frac{\mathbf{v}_r}{r^2} \right) + \rho g_r = 0 \quad (2.102)$$

$$-\frac{\partial p}{\partial z} + \eta \left(\frac{1}{r} \frac{\partial}{\partial r} \left(r \frac{\partial \mathbf{v}_z}{\partial r} \right) + \frac{\partial^2 \mathbf{v}_z}{\partial z^2} \right) + \rho g_z = 0 \quad (2.103)$$

$$\frac{1}{r} \frac{\partial}{\partial r} (r \mathbf{v}_r) + \frac{\partial \mathbf{v}_z}{\partial z} = 0 \quad (2.104)$$

The strain rate tensor in cylindrical coordinates is given by

$$\dot{\epsilon}_{rr} = \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.105)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \quad (2.106)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_\theta}{\partial r} - \frac{\mathbf{v}_\theta}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \quad (2.107)$$

$$\dot{\epsilon}_{zz} = \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.108)$$

$$\dot{\epsilon}_{rz} = \dot{\epsilon}_{zr} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_r}{\partial z} + \frac{\partial \mathbf{v}_z}{\partial r} \right) \quad (2.109)$$

$$\dot{\epsilon}_{\theta z} = \dot{\epsilon}_{z\theta} = \frac{1}{2} \left(\frac{1}{r} \frac{\partial \mathbf{v}_z}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial z} \right) \quad (2.110)$$

In the axisymmetric case, we have $\mathbf{v}_\theta = 0$ and $\partial_\theta \rightarrow 0$ so that

$$\dot{\epsilon}_{rr} = \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.111)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{\mathbf{v}_r}{r} \quad (2.112)$$

$$\dot{\epsilon}_{r\theta} = \dot{\epsilon}_{\theta r} = 0 \quad (2.113)$$


$$\dot{\epsilon}_{zz} = \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.114)$$

$$\dot{\epsilon}_{rz} = \dot{\epsilon}_{zr} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_r}{\partial z} + \frac{\partial \mathbf{v}_z}{\partial r} \right) \quad (2.115)$$

$$\dot{\epsilon}_{\theta z} = \dot{\epsilon}_{z\theta} = 0 \quad (2.116)$$

or,

$$\dot{\epsilon} = \begin{pmatrix} \dot{\epsilon}_{rr} & 0 & \dot{\epsilon}_{rz} \\ 0 & \dot{\epsilon}_{\theta\theta} & 0 \\ \dot{\epsilon}_{zr} & 0 & \dot{\epsilon}_{zz} \end{pmatrix}$$

 Relevant Literature: Daly & Raefsky (1985) [300], Kiefer & Hager (1992) [699].

This is implemented in **STONE** 36,63,90,91,92,96,106.

¹⁸https://en.wikipedia.org/wiki/Navier-Stokes_equations

In spherical coordinates

Assuming the flow velocity does not depend on ϕ ($\partial_\phi = 0$) and therefore also that $\mathbf{v}_\phi = 0$

$$0 = -\frac{\partial p}{\partial r} + f_r + \eta \left(\Delta v_r - \frac{2v_r}{r^2} - \frac{2}{r^2} \frac{\partial v_\theta}{\partial \theta} - \frac{2v_\theta \cot \theta}{r^2} \right)$$

$$0 = -\frac{1}{r} \frac{\partial p}{\partial \theta} + \eta \left(\Delta v_\theta + \frac{2}{r^2} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r^2 \sin^2 \theta} \right)$$

with

$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial}{\partial \theta} \right)$$

$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2}{\partial \phi^2}$$

THESE EQUATIONS SHOULD BE CHECKED and RE-CHECKED !!

From [1401]:

$$\frac{1}{r^2} \frac{\partial}{\partial r} (r^2 \mathbf{v}_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\mathbf{v}_\theta \sin \theta) + \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} = 0 \quad (2.117)$$

Pb with $1/r^2$??

$$0 = -\frac{\partial p}{\partial r} + (1 - \zeta) Ra \, r \, T + \frac{1}{r^2} \frac{\partial}{\partial r} \left(2\eta r^2 \frac{\partial \mathbf{v}_r}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\eta \sin \theta \frac{\partial \mathbf{v}_r}{\partial \theta} \right) + \frac{\partial}{\partial \theta} \left(\eta \frac{\partial}{\partial r} \frac{\mathbf{v}_\theta}{r} \right) \quad (118)$$

where $\zeta = R_i/R_o$

The dimensional form of the energy equation in a spherical axisymmetric geometry is given by (assuming the conductivity k to be constant):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \mathbf{v}_r \frac{\partial T}{\partial r} + \frac{\mathbf{v}_\theta}{r} \frac{\partial T}{\partial \theta} \right) = k \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial T}{\partial r} \right) + k \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial T}{\partial \theta} \right) \dots$$

THESE EQUATIONS SHOULD BE CHECKED and RE-CHECKED !!

2.15 The Boussinesq approximation

physics.tex

As nicely explained in Spiegel & Veronis [1186]: "In the study of problems of thermal convection it is a frequent practice to simplify the basic equations by introducing certain approximations which are attributed to Boussinesq (1903). The Boussinesq approximations can best be summarized by two statements:

1. The fluctuations in density which appear with the advent of motion result principally from thermal (as opposed to pressure) effects.
2. In the equations for the rate of change of momentum and mass, density variations may be neglected except when they are coupled to the gravitational acceleration in the buoyancy force."

Note that their paper examines the Boussinesq approximation for compressible fluids.

[from ASPECT manual] The Boussinesq approximation assumes that the density can be considered constant in all occurrences in the equations with the exception of the buoyancy term on the right hand side of (2.56). The primary result of this assumption is that the continuity equation (2.57) will now read $\vec{\nabla} \cdot \vec{v} = 0$. This implies that the strain rate tensor is deviatoric. Under the Boussinesq approximation, the equations are much simplified:

$$-\vec{\nabla} \cdot [2\eta\dot{\epsilon}(\vec{v})] + \vec{\nabla} p = \rho\vec{g} \quad \text{in } \Omega, \quad (2.119)$$

$$\vec{\nabla} \cdot (\rho\vec{v}) = 0 \quad \text{in } \Omega, \quad (2.120)$$

$$\rho_0 C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T = \rho H \quad \text{in } \Omega \quad (2.121)$$

Note that all terms on the rhs of the temperature equations have disappeared, with the exception of the source term.

2.16 The Extended Boussinesq approximation

physics.tex

Yuen *et al.* (2007) [1397] state that the background of the extended Boussinesq equations can be found described in Christensen and Yuen (1985) [251] and more completely in Matyska and Yuen (2007) [842].



Relevant Literature[543, 542]

2.17 Stokes equation for elastic medium

elastic_equations.tex

This will be moved to Section 16.7

What follows is mostly borrowed from Becker & Kaus lecture notes [66].

The strong form of the PDE that governs force balance in a medium is given by

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{f} = \vec{0}$$

where $\boldsymbol{\sigma}$ is the stress tensor and \vec{f} is a body force.

The stress tensor is related to the strain tensor through the generalised Hooke's law¹⁹:

$$\sigma_{ij} = \sum_{kl} C_{ijkl} \varepsilon_{kl} \quad \text{or} \quad \boldsymbol{\sigma} = \mathbf{C} : \boldsymbol{\varepsilon} \quad (2.122)$$

where \mathbf{C} is the fourth-order elastic tensor.

Due to the inherent symmetries of $\boldsymbol{\sigma}$, $\boldsymbol{\varepsilon}$, and \mathbf{C} , only 21 elastic coefficients of the latter are independent. For isotropic linear media (which have the same physical properties in any direction), \mathbf{C} can be reduced to only two independent numbers (for example the bulk modulus K and the shear modulus G that quantify the material's resistance to changes in volume and to shearing deformations, respectively). Thus

$$C_{ijkl} = \lambda \delta_{ij} \delta_{kl} + \mu (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk})$$

so that Eq. (2.122) becomes:

$$\sigma_{ij} = \lambda \varepsilon_{kk} \delta_{ij} + 2\mu \varepsilon_{ij}$$

or

$$\boldsymbol{\sigma} = \lambda (\vec{\nabla} \cdot \vec{u}) \mathbf{1} + 2\mu \boldsymbol{\varepsilon}(\vec{u}) \quad (2.123)$$

where λ is the Lamé parameter and μ is the shear modulus²⁰. The term $\vec{\nabla} \cdot \vec{u}$ is the isotropic dilation.

This can be re-written in the 6-dimensional stress/strain space as

$$\underbrace{\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix}}_{\vec{\sigma}} = \underbrace{\begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix}}_{\mathbf{C}} \cdot \underbrace{\begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{zz} \\ \varepsilon_{xy} \\ \varepsilon_{xz} \\ \varepsilon_{yz} \end{pmatrix}}_{\vec{\varepsilon}}$$

or, in terms of the compliance matrix \mathbf{C}^{-1} ,

$$\vec{\varepsilon} = \mathbf{C}^{-1} \cdot \vec{\sigma}$$

with

$$\mathbf{C}^{-1} = \frac{1}{\mu(3\lambda + 2\mu)} \begin{pmatrix} \lambda + \mu & -\lambda/2 & -\lambda/2 & 0 & 0 & 0 \\ -\lambda/2 & \lambda + \mu & -\lambda/2 & 0 & 0 & 0 \\ -\lambda/2 & -\lambda/2 & \lambda + \mu & 0 & 0 & 0 \\ 0 & 0 & 0 & 3\lambda + 2\mu & 0 & 0 \\ 0 & 0 & 0 & 0 & 3\lambda + 2\mu & 0 \\ 0 & 0 & 0 & 0 & 0 & 3\lambda + 2\mu \end{pmatrix}$$

¹⁹https://en.wikipedia.org/wiki/Hooke's_law

²⁰It is also sometimes written G

If we define the Young's modulus as $E = \mu(3\lambda + 2\mu)/(\lambda + \mu)$ and the Poisson's ratio as $\nu = \lambda(\lambda + \mu)/2$, then

$$\mathbf{C}^{-1} = \frac{1}{E} \begin{pmatrix} 1 & -\nu & -\nu & 0 & 0 & 0 \\ -\nu & 1 & -\nu & 0 & 0 & 0 \\ -\nu & -\nu & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2(1+\nu) & 0 & 0 \\ 0 & 0 & 0 & 0 & 2(1+\nu) & 0 \\ 0 & 0 & 0 & 0 & 0 & 2(1+\nu) \end{pmatrix}$$

Note that the determinant of \mathbf{C}^{-1} is $8(1+\nu)^5(1-2\nu)E^{-6}$, so that when $\nu \rightarrow 1/2$ (incompressible material), the compliance matrix is singular and the stress cannot be given as a function of strain [816].

The strain tensor is related to the displacement as follows:

$$\boldsymbol{\varepsilon}(\vec{u}) = \frac{1}{2}(\vec{\nabla}\vec{u} + (\vec{\nabla}\vec{u})^T)$$

The incompressibility (or bulk modulus) K is defined as $p = -K\vec{\nabla} \cdot \vec{u}$ where p is the pressure with

$$\begin{aligned} p &= -\frac{1}{3}\text{tr}(\boldsymbol{\sigma}) \\ &= -\frac{1}{3}[\lambda(\vec{\nabla} \cdot \vec{u})\text{tr}[\mathbf{1}] + 2\mu\text{tr}[\boldsymbol{\varepsilon}(\vec{u})]] \\ &= -\frac{1}{3}[\lambda(\vec{\nabla} \cdot \vec{u})3 + 2\mu(\vec{\nabla} \cdot \vec{u})] \\ &= -\left[\lambda + \frac{2}{3}\mu\right](\vec{\nabla} \cdot \vec{u}) \end{aligned} \tag{2.124}$$

so that


$$p = -K\vec{\nabla} \cdot \vec{u} \quad \text{with} \quad K = \lambda + \frac{2}{3}\mu$$

Remark. Eq. (2.122) and (2.123) are analogous to the ones that one has to solve in the context of viscous flow using the penalty method. In this case λ is the penalty coefficient, \vec{u} is the velocity, and μ is then the dynamic viscosity.

The Lamé parameter and the shear modulus are also linked to ν the poisson ratio, and E , Young's modulus:

$$\lambda = \mu \frac{2\nu}{1-2\nu} = \frac{\nu E}{(1+\nu)(1-2\nu)} \quad \text{with} \quad E = 2\mu(1+\nu)$$

The shear modulus, expressed often in GPa, describes the material's response to shear stress. The poisson ratio describes the response in the direction orthogonal to uniaxial stress. The Young modulus, expressed in GPa, describes the material's strain response to uniaxial stress in the direction of this stress.

 **Relevant Literature:** solvers for 3D Stokes and elasticity problems with heterogeneous coefficients [1106]

2.18 The strain rate tensor in all coordinate systems

strainrate_tensor.tex

The strain rate tensor $\dot{\epsilon}(\vec{v})$ is given by

$$\dot{\epsilon}(\vec{v}) = \frac{1}{2}(\vec{\nabla}\vec{v} + (\vec{\nabla}\vec{v})^T) \quad (2.125)$$

2.18.1 Cartesian coordinates

$$\dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} \quad (2.126)$$

$$\dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} \quad (2.127)$$

$$\dot{\epsilon}_{zz} = \frac{\partial w}{\partial z} \quad (2.128)$$

$$\dot{\epsilon}_{yx} = \dot{\epsilon}_{xy} = \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) \quad (2.129)$$

$$\dot{\epsilon}_{zx} = \dot{\epsilon}_{xz} = \frac{1}{2} \left(\frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \right) \quad (2.130)$$

$$\dot{\epsilon}_{zy} = \dot{\epsilon}_{yz} = \frac{1}{2} \left(\frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \right) \quad (2.131)$$

In the ASPECT manual there is an interesting discussion about the strain rate tensor in the case of 2D models: "The notion we adopt here is to think of two-dimensional models in the following way: We assume that the domain we want to solve on is a two-dimensional cross section (parameterized by x and y coordinates) that extends infinitely far in both negative and positive z direction. Further, we assume that the velocity is zero in z direction and that all variables have no variation in z direction. As a consequence, we ought to really think of these two-dimensional models as three-dimensional ones in which the z component of the velocity is zero and so are all z derivatives."

This of course makes sense but it means that when the deviatoric strain rate tensor needs to be computed, then it is given by

$$\dot{\epsilon}^d = \dot{\epsilon} - \frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & 0 \\ \dot{\epsilon}_{xy} & \dot{\epsilon}_{yy} & 0 \\ 0 & 0 & 0 \end{pmatrix} - \frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy})\mathbf{1} = \frac{1}{3} \begin{pmatrix} 2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} & 3\dot{\epsilon}_{xy} & 0 \\ 3\dot{\epsilon}_{xy} & -\dot{\epsilon}_{xx} + 2\dot{\epsilon}_{yy} & 0 \\ 0 & 0 & -\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} \end{pmatrix}$$

As a consequence the shear heating term Φ is given by

$$\begin{aligned} \Phi = 2\eta\dot{\epsilon}^d : \dot{\epsilon}^d &= 2\eta\frac{1}{9} [(2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + (-\dot{\epsilon}_{xx} + 2\dot{\epsilon}_{yy})^2 + 2 \cdot 9\dot{\epsilon}_{xy}^2 + (-\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2] \\ &= 2\eta\frac{1}{9} [4\dot{\epsilon}_{xx}^2 - 4\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + \dot{\epsilon}_{yy}^2 + \dot{\epsilon}_{xx}^2 - 4\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + 4\dot{\epsilon}_{yy}^2 + 18\dot{\epsilon}_{xy}^2 + \dot{\epsilon}_{xx}^2 + 2\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + \dot{\epsilon}_{yy}^2] \\ &= 2\eta\frac{1}{9} [6\dot{\epsilon}_{xx}^2 + 6\dot{\epsilon}_{yy}^2 - 6\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + 18\dot{\epsilon}_{xy}^2] \\ &= 2\eta \left[\frac{2}{3}\dot{\epsilon}_{xx}^2 + \frac{2}{3}\dot{\epsilon}_{yy}^2 - \frac{2}{3}\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{xy}^2 \right] \end{aligned} \quad (2.132)$$

2.18.2 Polar coordinates

$$\dot{\epsilon}_{rr} = \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.133)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \quad (2.134)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_\theta}{\partial r} - \frac{\mathbf{v}_\theta}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \quad (2.135)$$

2.18.3 Cylindrical coordinates

$$\dot{\epsilon}_{rr} = \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.136)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \quad (2.137)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_\theta}{\partial r} - \frac{\mathbf{v}_\theta}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \quad (2.138)$$

$$\dot{\epsilon}_{zz} = \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.139)$$

$$\dot{\epsilon}_{rz} = \dot{\epsilon}_{zr} = \frac{1}{2} \left(\frac{\partial \mathbf{v}_r}{\partial z} + \frac{\partial \mathbf{v}_z}{\partial r} \right) \quad (2.140)$$

$$\dot{\epsilon}_{\theta z} = \dot{\epsilon}_{z\theta} = \frac{1}{2} \left(\frac{1}{r} \frac{\partial \mathbf{v}_z}{\partial \theta} + \frac{\partial \mathbf{v}_\theta}{\partial z} \right) \quad (2.141)$$

The velocity divergence is given by

$$\vec{\nabla} \cdot \vec{\mathbf{v}} = \dot{\epsilon}_{rr} + \dot{\epsilon}_{\theta\theta} + \dot{\epsilon}_{zz} = \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r} \left(\frac{\partial \mathbf{v}_\theta}{\partial \theta} + \mathbf{v}_r \right) + \frac{\partial \mathbf{v}_z}{\partial z} \quad (2.142)$$

2.18.4 Spherical coordinates

$$\dot{\epsilon}_{rr} = \frac{\partial \mathbf{v}_r}{\partial r} \quad (2.143)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{\mathbf{v}_r}{r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \quad (2.144)$$

$$\dot{\epsilon}_{\phi\phi} = \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\phi}{\partial \phi} + \frac{\mathbf{v}_r}{r} + \frac{\mathbf{v}_\theta \cot \theta}{r} \quad (2.145)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\theta}{r} \right) + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial \theta} \right) \quad (2.146)$$

$$\dot{\epsilon}_{\phi r} = \dot{\epsilon}_{r\phi} = \frac{1}{2} \left(\frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_r}{\partial \phi} + r \frac{\partial}{\partial r} \left(\frac{\mathbf{v}_\phi}{r} \right) \right) \quad (2.147)$$

$$\dot{\epsilon}_{\phi\theta} = \dot{\epsilon}_{\theta\phi} = \frac{1}{2} \left(\frac{\sin \theta}{r} \frac{\partial}{\partial \theta} \left(\frac{\mathbf{v}_\phi}{\sin \theta} \right) + \frac{1}{r \sin \theta} \frac{\partial \mathbf{v}_\theta}{\partial \phi} \right) \quad (2.148)$$

2.18.5 Relationship between Cartesian and polar coordinates expressions

We can go from Cartesian to polar coordinates via the 2×2 transformation matrix:

$$\mathcal{P} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \quad (2.149)$$

The rows correspond to the components of \vec{e}_r and \vec{e}_θ in the Cartesian basis. A vector \vec{v} transforms from one orthonormal basis to another by multiplying it by the matrix \mathcal{P} . As we have seen before, this yields

$$\mathbf{v}_r = u \cos \theta + v \sin \theta \quad (2.150)$$

$$\mathbf{v}_\theta = -u \sin \theta + v \cos \theta \quad (2.151)$$

A second-order tensor \mathbf{a} in Cartesian coordinates transforms into \mathbf{a}^* in polar coordinates by

$$\mathbf{a}^* = \mathcal{P} \cdot \mathbf{a} \cdot \mathcal{P}^T$$

and obviously

$$\mathbf{a} = \mathcal{P}^T \cdot \mathbf{a}^* \cdot \mathcal{P}$$

We obtain for the strain rate tensor (or the stress tensor):

$$\begin{aligned} \dot{\epsilon}_{rr} &= \dot{\epsilon}_{xx} \cos^2 \theta + \dot{\epsilon}_{yy} \sin^2 \theta + 2\dot{\epsilon}_{xy} \sin \theta \cos \theta \\ \dot{\epsilon}_{\theta\theta} &= \dot{\epsilon}_{xx} \sin^2 \theta + \dot{\epsilon}_{yy} \cos^2 \theta - 2\dot{\epsilon}_{xy} \sin \theta \cos \theta \\ \dot{\epsilon}_{r\theta} &= \dot{\epsilon}_{xy}(\cos^2 \theta - \sin^2 \theta) + (\dot{\epsilon}_{yy} - \dot{\epsilon}_{xx}) \sin \theta \cos \theta \end{aligned}$$

Using the trigonometric identities $\sin 2\theta = 2 \sin \theta \cos \theta$ and $\cos^2 \theta - \sin^2 \theta = \cos 2\theta$, then we obtain

$$\begin{aligned} \dot{\epsilon}_{rr} &= \dot{\epsilon}_{xx} \cos^2 \theta + \dot{\epsilon}_{yy} \sin^2 \theta + \dot{\epsilon}_{xy} \sin 2\theta \\ \dot{\epsilon}_{\theta\theta} &= \dot{\epsilon}_{xx} \sin^2 \theta + \dot{\epsilon}_{yy} \cos^2 \theta - \dot{\epsilon}_{xy} \sin 2\theta \\ \dot{\epsilon}_{r\theta} &= \dot{\epsilon}_{xy} \cos 2\theta + \frac{1}{2}(\dot{\epsilon}_{yy} - \dot{\epsilon}_{xx}) \sin 2\theta \end{aligned}$$

and likewise:

$$\dot{\epsilon}_{xx} = \dot{\epsilon}_{rr} \cos^2 \theta + \dot{\epsilon}_{\theta\theta} \sin^2 \theta - 2\dot{\epsilon}_{r\theta} \sin \theta \cos \theta \quad (2.152)$$

$$\dot{\epsilon}_{yy} = \dot{\epsilon}_{rr} \sin^2 \theta + \dot{\epsilon}_{\theta\theta} \cos^2 \theta + 2\dot{\epsilon}_{r\theta} \sin \theta \cos \theta \quad (2.153)$$

$$\dot{\epsilon}_{xy} = \dot{\epsilon}_{r\theta}(\cos^2 \theta - \sin^2 \theta) + (\dot{\epsilon}_{rr} - \dot{\epsilon}_{\theta\theta}) \sin \theta \cos \theta \quad (2.154)$$

2.19 Boundary conditions

physics.tex

In mathematics, the Dirichlet (or first-type) boundary condition is a type of boundary condition, named after Peter Gustav Lejeune Dirichlet. When imposed on an ODE or PDE, it specifies the values that a solution needs to take along the boundary of the domain. Note that a Dirichlet boundary condition may also be referred to as a fixed boundary condition.

The Neumann (or second-type) boundary condition is a type of boundary condition, named after Carl Neumann. When imposed on an ordinary or a partial differential equation, the condition specifies the values in which the derivative of a solution is applied within the boundary of the domain.

It is possible to describe the problem using other boundary conditions: a Dirichlet boundary condition specifies the values of the solution itself (as opposed to its derivative) on the boundary, whereas the Cauchy boundary condition, mixed boundary condition and Robin boundary condition are all different types of combinations of the Neumann and Dirichlet boundary conditions.

2.19.1 The Stokes equations

You may find the following terms in the computational geodynamics literature:

- free surface: this means that no force is acting on the surface, i.e. $\boldsymbol{\sigma} \cdot \vec{n} = \vec{0}$. It is usually used on the top boundary of the domain and allows for topography evolution.
- free slip: $\vec{v} \cdot \vec{n} = 0$ and $(\boldsymbol{\sigma} \cdot \vec{n}) \times \vec{n} = \vec{0}$. This condition ensures a frictionless flow parallel to the boundary where it is prescribed.
- no slip: this means that the velocity (or displacement) is exactly zero on the boundary, i.e. $\vec{v} = \vec{0}$.
- prescribed velocity: $\vec{v} = \vec{v}_{bc}$
- stress b.c.:
- open .b.c.: see [STONE](#) 29.

2.19.2 The heat transport equation

There are two types of boundary conditions for this equation: temperature boundary conditions (Dirichlet boundary conditions) and heat flux boundary conditions (Neumann boundary conditions).

2.20 Meaningful physical quantities

physics.tex

- **Velocity** \vec{v} (m/s): This is a vector quantity and both magnitude and direction are needed to define it. It is the rate of change of position with respect to a frame of reference.
- **Root mean square velocity** v_{rms} (m/s):

$$v_{rms} = \left(\frac{\int_{\Omega} |\vec{v}|^2 dV}{\int_{\Omega} dV} \right)^{1/2} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} |\vec{v}|^2 dV \right)^{1/2} \quad (2.155)$$

Remark. V_{Ω} is usually computed numerically at the same time that v_{rms} is computed.

In Cartesian coordinates, for a cuboid domain of size $Lx \times Ly \times Lz$, the v_{rms} is simply given by:

$$v_{rms} = \left(\frac{1}{LxLyLz} \int_0^{Lx} \int_0^{Ly} \int_0^{Lz} (u^2 + v^2 + w^2) dx dy dz \right)^{1/2} \quad (2.156)$$

In the case of an annulus domain, although calculations are carried out in Cartesian coordinates, it makes sense to look at the radial velocity component v_r and the tangential velocity component v_{θ} , and their respective root mean square averages:

$$v_r|_{rms} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} v_r^2 d\Omega \right)^{1/2} \quad (2.157)$$

$$v_{\theta}|_{rms} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} v_{\theta}^2 d\Omega \right)^{1/2} \quad (2.158)$$

- **Pressure** p (Pa):
- **Stress tensor** σ (Pa):
- **Strain tensor** ϵ (dimensionless):
- **Strain rate tensor** $\dot{\epsilon}$ (s^{-1}):
- **Argand Number**: Non-dimensional number (Ar) representing the ratio of the stress arising from crustal thickness contrasts (vertical stress) to the stress required to deform the material at ambient strain rates (horizontal stress) It is commonly used in mountain building dynamics as a measure of the tendency of an orogen to collapse under its own gravitational potential energy. See England & McKenzie [376], Houseman & England [596].
- **(Thermal) Rayleigh number** Ra (or Ra_T) (X): It is a dimensionless number that expresses the ratio of the driving forces to the opposing forces. The buoyancy force comes from the volumetric thermal expansion while the viscous forces and the heat diffusivity oppose convection (the latter one smoothes out thermal gradients).

The Rayleigh number for convection driven by a constant temperature hot base and a cold surface in a domain of thickness D is:

$$Ra = \frac{\rho_0 g \alpha D^3}{\eta \kappa} \cdot \Delta T = \frac{\rho_0^2 C_p g \alpha D^3 \Delta T}{\eta k}$$

The Rayleigh number for convection driven by a hot base (constant basal heat flow q_b) and a colder surface is:

$$\text{Ra} = \frac{\rho_0 g \alpha D^3}{\eta \kappa} \cdot \frac{q_b D}{k}$$

The Rayleigh number for convection driven by internal heating H (production per cubic meter) is:

$$\text{Ra} = \frac{\rho_0 g \alpha D^3}{\eta \kappa} \cdot \frac{H D^2}{k}$$

The Rayleigh number for convection driven by both basal heat flow and internal heating is:

$$\text{Ra} = \frac{\rho_0 g \alpha D^3}{\eta \kappa} \cdot \frac{q_b D + H D^2}{k}$$

For convection to occur, the Rayleigh number must be larger than the so-called critical Rayleigh number, which ranges from 600 to 3000 (it depends on the boundary conditions and the geometry).

- **Compositional Rayleigh Number** Ra_C :

$$\text{Ra}_C = \frac{\Delta \rho_C g D^3}{\kappa \eta_0}$$

where $\Delta \rho_C$ is the difference in density between the distinct material compositions (when compared at identical temperatures). See for instance Trim *et al.* (2020) [1281].

- **Prandtl number** Pr (X): It is named after the German physicist Ludwig Prandtl²¹ and is defined as the ratio of momentum diffusivity to thermal diffusivity. It is given as:

$$\text{Pr} = \frac{\text{momentum diffusivity}}{\text{thermal diffusivity}} = \frac{\eta / \rho}{k / (\rho C_p)} = \frac{\eta C_p}{k}$$

For Earth materials, we have $\text{Pr} \sim (10^{21} 1000) / 3 \gg 1$, which means that momentum diffusivity dominates.

- **Nusselt number** Nu (X): the Nusselt number (Nu) is the ratio of convective to conductive heat transfer across (normal to) the boundary. The conductive component is measured under the same conditions as the heat convection but with a (hypothetically) stagnant (or motionless) fluid.

In practice the Nusselt number Nu of a layer (typically the mantle of a planet) is defined as follows:

$$\text{Nu} = \frac{q}{q_c} \tag{2.159}$$

where q is the heat transferred by convection while $q_c = k \Delta T / D$ is the amount of heat that would be conducted through a layer of thickness D with a temperature difference ΔT across it with k being the thermal conductivity.

For 2D Cartesian systems of size (L_x, L_y) the Nu is computed [95]

$$\text{Nu} = \frac{\frac{1}{L_x} \int_0^{L_x} k \frac{\partial T}{\partial y}(x, y = L_y) dx}{-\frac{1}{L_x} \int_0^{L_x} k T(x, y = 0) / L_y dx} = -L_y \frac{\int_0^{L_x} \frac{\partial T}{\partial y}(x, y = L_y) dx}{\int_0^{L_x} T(x, y = 0) dx}$$

i.e. it is the mean surface temperature gradient over the mean bottom temperature.

²¹https://en.wikipedia.org/wiki/Ludwig_Prandtl

Note that in the case when no convection takes place then the measured heat flux at the top is the one obtained from a purely conductive profile which yields $\text{Nu}=1$.

Note that a relationship $\text{Ra} \propto \text{Nu}^\alpha$ exists between the Rayleigh number Ra and the Nusselt number Nu in convective systems, see [1368] and references therein.

Turning now to cylindrical geometries with inner radius R_1 and outer radius R_2 , we define $f = R_1/R_2$. A small value of f corresponds to a high degree of curvature. We assume now that $R_2 - R_1 = 1$, so that $R_2 = 1/(1 - f)$ and $R_1 = f/(1 - f)$. Following [635], the Nusselt number at the inner and outer boundaries are:

$$\text{Nu}_{inner} = \frac{f \ln f}{1 - f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T}{\partial r} \right)_{r=R_1} d\theta \quad (2.160)$$

$$\text{Nu}_{outer} = \frac{\ln f}{1 - f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T}{\partial r} \right)_{r=R_2} d\theta \quad (2.161)$$

Note that a conductive geotherm in such an annulus between temperatures T_1 and T_2 is given by

$$T_c(r) = \frac{\ln(r/R_2)}{\ln(R_1/R_2)} = \frac{\ln(r(1 - f))}{\ln f}$$

so that

$$\frac{\partial T_c}{\partial r} = \frac{1}{r} \frac{1}{\ln f}$$

We then find:

$$\text{Nu}_{inner} = \frac{f \ln f}{1 - f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T_c}{\partial r} \right)_{r=R_1} d\theta = \frac{f \ln f}{1 - f} \frac{1}{R_1} \frac{1}{\ln f} = 1 \quad (2.162)$$

$$\text{Nu}_{outer} = \frac{\ln f}{1 - f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T_c}{\partial r} \right)_{r=R_2} d\theta = \frac{\ln f}{1 - f} \frac{1}{R_2} \frac{1}{\ln f} = 1 \quad (2.163)$$

As expected, the recovered Nusselt number at both boundaries is exactly 1 when the temperature field is given by a steady state conductive geotherm.

derive formula for Earth size R1 and R2

Relevant Literature[800]

- **Temperature** (K):
- **(Dynamic) Viscosity** (Pas): For air it is roughly 10^{-5} Pas, about 10^{-3} Pas for water, about 10^{10} Pas for ice and about 10^{17} Pas for salt.
- **Entropy** S (J K^{-1})
- **(mass) Density** ρ (kg m^{-3}):
- **Heat capacity** C_p (J K^{-1}): It is the measure of the heat/energy required to increase the temperature of a unit quantity of a substance by unit degree. Note that the *specific* heat capacity c_p of a substance is the heat capacity of a sample of the substance divided by the mass of the sample, with units $\text{J K}^{-1} \text{kg}^{-1}$.

“Different substances respond to heat in different ways. If a metal chair sits in the bright sun on a hot day, it may become quite hot to the touch. An equal mass of water under the same

sun exposure will not become nearly as hot. This means that water has a high heat capacity (the amount of heat required to raise the temperature of an object by 1 °C). Water is very resistant to changes in temperature, while metals generally are not.” ²²

- **Heat conductivity**, or thermal conductivity k ($\text{W m}^{-1} \text{K}^{-1}$). It is the property of a material that indicates its ability to conduct heat. It appears primarily in Fourier’s Law for heat conduction. Note that it is a function of temperature, especially in mantle convection settings, see Bonneville & Capolsini (1999) [115] and refs therein, Miyauchi & Kameyama (2013) [885], Hofmeister & Yuen (2007) [583]. Note also that it can be a tensorial quantity in anisotropic context. The heat conductivity of many rocks was determined in [23].
- **Heat diffusivity**: $\kappa = k/(\rho C_p)$ ($\text{m}^2 \text{s}^{-1}$). Substances with high thermal diffusivity rapidly adjust their temperature to that of their surroundings, because they conduct heat quickly in comparison to their volumetric heat capacity or ‘thermal bulk’.
- **thermal expansion** α (K^{-1}): it is the tendency of a matter to change in volume in response to a change in temperature. Note that it is a function of temperature, especially in mantle convection settings [885].

$$\alpha = \frac{1}{V} \left(\frac{\partial V}{\partial T} \right)_P$$

- **Urey Ratio**: mantle heat production divided by heat loss. It is a key constraint for thermal history models. Recent Urey ratio estimates are in the range of 0.21-0.49. [768]
- **Shear modulus**: modulus of rigidity, usually expressed in GPa. It describes the material response to shear stress.
- **Poisson ratio**: response in the direction orthogonal to uniaxial stress.
- **Young’s modulus**: describes the material strain response to uniaxial stress in the direction of this stress, usually expressed in GPa.
- **Average viscosity**: following Christensen (1983) [242], one can compute the averaged viscosity in a domain as follows:

$$\langle \eta \rangle = \frac{\int_V \eta \dot{\epsilon}_e^2 dV}{\int_V \dot{\epsilon}_e^2 dV} \quad (2.164)$$

check aspect manual The 2D cylindrical shell benchmarks by Davies *et al.* 5.4.12

²²[https://chem.libretexts.org/Bookshelves/Introductory_Chemistry/Introductory_Chemistry_\(CK-12\)/17%3A_Thermochemistry/17.04%3A_Heat_Capacity_and_Specific_Heat](https://chem.libretexts.org/Bookshelves/Introductory_Chemistry/Introductory_Chemistry_(CK-12)/17%3A_Thermochemistry/17.04%3A_Heat_Capacity_and_Specific_Heat)

2.21 Principal stress and principal invariants

physics.tex

As seen before (see Section 2.4.2) the stress tensor is a symmetric 3×3 real matrix, and linear algebra tells us that it therefore has three mutually orthogonal unit-length eigenvectors $\vec{n}_1, \vec{n}_2, \vec{n}_3$ and three real eigenvalues $\lambda_1, \lambda_2, \lambda_3$ such that $\boldsymbol{\sigma} \cdot \vec{n}_i = \lambda_i \vec{n}_i$.

As a consequence, in a coordinate system with axes $\vec{n}_1, \vec{n}_2, \vec{n}_3$, the stress tensor is a diagonal matrix, and has only the three normal components $\lambda_1, \lambda_2, \lambda_3$ i.e. the principal stresses. If the three eigenvalues are equal, the stress is an isotropic compression or tension, always perpendicular to any surface, there is no shear stress, and the tensor is a diagonal matrix in any coordinate frame.

2.21.1 In two dimensions

We are looking for the stress tensor eigenvector vector $\vec{n} = (n_x, n_y)$ associated to the eigenvalue λ such that

$$\begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \end{pmatrix} = \lambda \begin{pmatrix} n_x \\ n_y \end{pmatrix}$$

or,

$$\begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \end{pmatrix} = \vec{0}$$

i.e.,

$$\begin{pmatrix} \sigma_{xx} - \lambda & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} - \lambda \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \end{pmatrix} = \vec{0}$$

which yields

$$(\sigma_{xx} - \lambda)(\sigma_{yy} - \lambda) - \sigma_{xy}^2 = 0$$

or,

$$\lambda^2 - (\sigma_{xx} + \sigma_{yy})\lambda + (\sigma_{xx}\sigma_{yy} - \sigma_{xy}^2) = 0$$

The discriminant Δ is

$$\begin{aligned} \Delta &= (\sigma_{xx} + \sigma_{yy})^2 - 4(\sigma_{xx}\sigma_{yy} - \sigma_{xy}^2) \\ &= (\sigma_{xx} - \sigma_{yy})^2 + 4\sigma_{xy}^2 \end{aligned}$$

The roots are given by:

$$\begin{aligned} \lambda_{\pm} &= \frac{(\sigma_{xx} + \sigma_{yy}) \pm \sqrt{(\sigma_{xx} - \sigma_{yy})^2 + 4\sigma_{xy}^2}}{2} \\ &= \frac{\sigma_{xx} + \sigma_{yy}}{2} \pm \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2} \end{aligned}$$

The two principal stresses are then:

$$\begin{aligned} \sigma_1 &= \frac{\sigma_{xx} + \sigma_{yy}}{2} + \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2} \\ \sigma_2 &= \frac{\sigma_{xx} + \sigma_{yy}}{2} - \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2} \end{aligned} \tag{2.165}$$

with the convention $\sigma_1 > \sigma_2$. The maximum shear stress is defined as one-half the difference between the two principal stresses

$$\tau_{max} = \frac{\sigma_1 - \sigma_2}{2} = \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2} \quad (2.166)$$

The eigenvector \vec{n}_1 corresponding to σ_1 is obtained by solving

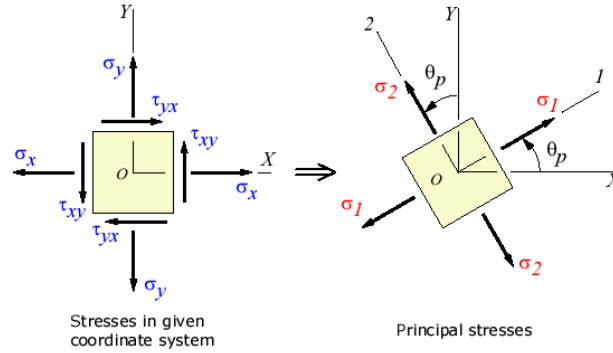
$$\boldsymbol{\sigma} \cdot \vec{n}_1 = \sigma_1 \vec{n}_1$$

and same for the other eigenvalue/vector:

$$\boldsymbol{\sigma} \cdot \vec{n}_2 = \sigma_2 \vec{n}_2$$

Each is a system of two equations with two unknowns. These are not difficult to solve, but can prove cumbersome. Note that linear algebra tells us that $\vec{n}_1 \cdot \vec{n}_2 = 0$, i.e. the eigenvectors form a basis of \mathbb{R}^2 .

This is the reason why often people go another route. One can ask the question: what is the value of the angle θ_p which, if used to perform a rotation of the axis system, yields a stress tensor that is diagonal, with the principal stresses on the diagonal?



Taken from https://www.efunda.com/formulae/solid_mechanics/mat_mechanics/plane_stress_principal.cfm

The rotation matrix is

$$\mathbf{R} = \begin{pmatrix} \cos \theta_p & -\sin \theta_p \\ \sin \theta_p & \cos \theta_p \end{pmatrix}$$

and the image of $\boldsymbol{\sigma}$ by means of the axis rotation is $\boldsymbol{\sigma}' = \mathbf{R} \cdot \boldsymbol{\sigma} \cdot \mathbf{R}^{-1}$, i.e.

$$\begin{aligned} \boldsymbol{\sigma}' &= \begin{pmatrix} \cos \theta_p & -\sin \theta_p \\ \sin \theta_p & \cos \theta_p \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} \end{pmatrix} \cdot \begin{pmatrix} \cos \theta_p & \sin \theta_p \\ -\sin \theta_p & \cos \theta_p \end{pmatrix} \\ &= \begin{pmatrix} \cos \theta_p & -\sin \theta_p \\ \sin \theta_p & \cos \theta_p \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \cos \theta_p - \sigma_{xy} \sin \theta_p & \sigma_{xx} \sin \theta_p + \sigma_{xy} \cos \theta_p \\ \sigma_{xy} \cos \theta_p - \sigma_{yy} \sin \theta_p & \sigma_{xy} \sin \theta_p + \sigma_{yy} \cos \theta_p \end{pmatrix} \\ &= \begin{pmatrix} \dots & \cos \theta_p(\sigma_{xx} \sin \theta_p + \sigma_{xy} \cos \theta_p) - \sin \theta_p(\sigma_{xy} \sin \theta_p + \sigma_{yy} \cos \theta_p) \\ \dots & \dots \end{pmatrix} \end{aligned}$$

In the matrix above I have only computed the off diagonal term since we are actually looking for θ_p such that $\sigma'_{xy} = 0$, or

$$\begin{aligned} \cos \theta_p(\sigma_{xx} \sin \theta_p + \sigma_{xy} \cos \theta_p) - \sin \theta_p(\sigma_{xy} \sin \theta_p + \sigma_{yy} \cos \theta_p) &= 0 \\ \sin \theta_p \cos \theta_p(\sigma_{xx} - \sigma_{yy}) + (\cos^2 \theta_p - \sin^2 \theta_p)\sigma_{xy} &= 0 \end{aligned}$$

and then

$$\frac{\sin \theta_p \cos \theta_p}{\cos^2 \theta_p - \sin^2 \theta_p} = \frac{\sigma_{xy}}{\sigma_{xx} - \sigma_{yy}}$$

The left hand term is actually a trigonometric identity²³:

$$\frac{\sin \theta_p \cos \theta_p}{\cos^2 \theta_p - \sin^2 \theta_p} = \frac{\frac{1}{2} \sin 2\theta_p}{\cos 2\theta_p} = \frac{1}{2} \tan 2\theta_p$$

and finally:

$$\tan 2\theta_p = \frac{2\sigma_{xy}}{\sigma_{xx} - \sigma_{yy}} \quad \text{or} \quad \boxed{\theta_p = \frac{1}{2} \tan^{-1} \frac{2\sigma_{xy}}{\sigma_{xx} - \sigma_{yy}}}$$

Once θ_p has been found the other direction is given by $\theta_p + \pi/2$.

Example: Let us assume a diagonal stress tensor of the form

$$\boldsymbol{\sigma} = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}$$

then $\tan 2\theta_p = 0$, and then $\theta_p = 0$. The principal directions are the horizontal and vertical directions, i.e. the Cartesian axis system, which is consistent.

add a remark that 2D does not exist and that plane strain incompressible actually is what is going on above

2.21.2 In three dimensions

We are looking for the stress tensor eigenvector vector $\vec{n} = (n_x, n_y, n_z)$ associated to the eigenvalue λ such that

$$\begin{pmatrix} \sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \\ n_z \end{pmatrix} = \lambda \begin{pmatrix} n_x \\ n_y \\ n_z \end{pmatrix}$$

or,

$$\begin{pmatrix} \sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \\ n_z \end{pmatrix} - \begin{pmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \\ n_z \end{pmatrix} = \vec{0}$$

$$\begin{pmatrix} \sigma_{xx} - \lambda & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} - \lambda & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} - \lambda \end{pmatrix} \cdot \begin{pmatrix} n_x \\ n_y \\ n_z \end{pmatrix} = \vec{0}$$

Non-trivial solutions of this equation require

$$\begin{vmatrix} \sigma_{xx} - \lambda & \sigma_{xy} & \sigma_{xz} \\ \sigma_{xy} & \sigma_{yy} - \lambda & \sigma_{yz} \\ \sigma_{xz} & \sigma_{yz} & \sigma_{zz} - \lambda \end{vmatrix} = 0$$

Expanding the determinant results in the following cubic equation:

$$\begin{aligned} 0 &= (\sigma_{xx} - \lambda)[(\sigma_{yy} - \lambda)(\sigma_{zz} - \lambda) - \sigma_{yz}^2] - \sigma_{xy}[\sigma_{xy}(\sigma_{zz} - \lambda) - \sigma_{yz}\sigma_{xz}] + \sigma_{xz}[\sigma_{xy}\sigma_{yz} - (\sigma_{yy} - \lambda)\sigma_{xz}] \\ &= (\sigma_{xx} - \lambda)[\sigma_{yy}\sigma_{zz} - \lambda(\sigma_{yy} + \sigma_{zz}) + \lambda^2 - \sigma_{yz}^2] - \sigma_{xy}[\sigma_{xy}(\sigma_{zz} - \lambda) - \sigma_{yz}\sigma_{xz}] + \sigma_{xz}[\sigma_{xy}\sigma_{yz} - (\sigma_{yy} - \lambda)\sigma_{xz}] \\ &= -\lambda^3 + (\sigma_{xx} + \sigma_{yy} + \sigma_{zz})\lambda^2 + (-\sigma_{yy}\sigma_{zz} - \sigma_{xx}\sigma_{yy} - \sigma_{xx}\sigma_{zz} + \sigma_{yz}^2 + \sigma_{xy}^2 + \sigma_{xz}^2)\lambda + \det(\boldsymbol{\sigma}) \end{aligned}$$

²³https://en.wikipedia.org/wiki/List_of_trigonometric_identities

or, after multiplying the last line by -1,

$$\lambda^3 - \mathcal{K}_1(\boldsymbol{\sigma})\lambda^2 + \mathcal{K}_2(\boldsymbol{\sigma})\lambda - \mathcal{K}_3(\boldsymbol{\sigma}) = 0 \quad (2.167)$$

with²⁴:

$$\begin{aligned} \mathcal{K}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\ \mathcal{K}_2(\boldsymbol{\sigma}) &= \sigma_{xx}\sigma_{yy} + \sigma_{yy}\sigma_{zz} + \sigma_{xx}\sigma_{zz} - \sigma_{xy}^2 - \sigma_{xz}^2 - \sigma_{yz}^2 \\ \mathcal{K}_3(\boldsymbol{\sigma}) &= \det(\boldsymbol{\sigma}) \\ &= \sigma_{xx}\sigma_{yy}\sigma_{zz} - \sigma_{xx}\sigma_{yz}^2 - \sigma_{xy}^2\sigma_{zz} + \sigma_{xy}\sigma_{yz}\sigma_{xz} + \sigma_{xz}\sigma_{xy}\sigma_{yz} - \sigma_{xz}^2\sigma_{yy} \\ &= \sigma_{xx}\sigma_{yy}\sigma_{zz} + 2\sigma_{xy}\sigma_{yz}\sigma_{xz} - (\sigma_{xx}\sigma_{yz}^2 + \sigma_{zz}\sigma_{xy}^2 + \sigma_{yy}\sigma_{xz}^2) \end{aligned} \quad (2.168)$$

\mathcal{K}_1 , \mathcal{K}_2 and \mathcal{K}_3 are called **principal invariants**²⁵ (see also Appendix A.1 of Zienkiewicz & Taylor [1431] or Eq. (6.4) of Freudenthal & Geiringer [418]). These invariants can be written in a coordinate-free manner²⁶:

$$\begin{aligned} \mathcal{K}_1(\boldsymbol{\sigma}) &= \text{tr}(\boldsymbol{\sigma}) \\ \mathcal{K}_2(\boldsymbol{\sigma}) &= \frac{1}{2}(\text{tr}(\boldsymbol{\sigma})^2 - \text{tr}(\boldsymbol{\sigma}^2)) \\ \mathcal{K}_3(\boldsymbol{\sigma}) &= \det(\boldsymbol{\sigma}) \end{aligned}$$

and if the stress tensor is diagonal, we have

$$\begin{aligned} \mathcal{K}_1(\boldsymbol{\sigma}) &= \sigma_1 + \sigma_2 + \sigma_3 \\ \mathcal{K}_2(\boldsymbol{\sigma}) &= \sigma_1\sigma_2 + \sigma_2\sigma_3 + \sigma_1\sigma_3 \\ \mathcal{K}_3(\boldsymbol{\sigma}) &= \sigma_1\sigma_2\sigma_3 \end{aligned}$$

The principal invariants $\mathcal{K}_{\{1,2,3\}}$ are related to the **moment invariants** $\mathcal{I}_{\{1,2,3\}}$ (see Section 2.22) as follows (Appendix A.2 of Zienkiewicz & Taylor [1431]):

$$\mathcal{I}_1(\boldsymbol{\sigma}) = \mathcal{K}_1(\boldsymbol{\sigma}) \quad (2.169)$$

$$\mathcal{I}_2(\boldsymbol{\sigma}) = \frac{1}{2}\mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma}) \quad (2.170)$$

$$\mathcal{I}_3(\boldsymbol{\sigma}) = \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma}) \quad (2.171)$$

write proofs in appendix

Very often we will find ourselves interested in the principal components of the deviatoric stress tensor $\boldsymbol{\tau}$ so that we now have the following determinant to compute:

$$\begin{vmatrix} \tau_{xx} - \lambda & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} - \lambda & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} - \lambda \end{vmatrix} = 0$$

and therefore obtain the following cubic equation

$$\lambda^3 - \mathcal{K}_1(\boldsymbol{\tau})\lambda^2 + \mathcal{K}_2(\boldsymbol{\tau})\lambda - \mathcal{K}_3(\boldsymbol{\tau}) = 0 \quad (2.172)$$

²⁴Note that in the equation (2.167) there is often a plus sign in front of \mathcal{K}_2 but not always. Be careful when reading literature!

²⁵https://en.wikipedia.org/wiki/Invariants_of_tensors

²⁶Proofs are in Appendix T

By definition of a deviatoric tensor we have $\mathcal{K}_1(\boldsymbol{\tau}) = 0$ and then Eqs. (2.170) and (2.171) become

$$\mathcal{I}_2(\boldsymbol{\tau}) = -\mathcal{K}_2(\boldsymbol{\tau}) \quad (2.173)$$

$$\mathcal{I}_3(\boldsymbol{\tau}) = \mathcal{K}_3(\boldsymbol{\tau}) \quad (2.174)$$

so that the cubic equation becomes

$$\lambda^3 - \mathcal{I}_2(\boldsymbol{\tau})\lambda - \mathcal{I}_3(\boldsymbol{\tau}) = 0 \quad (2.175)$$

Noting the trigonometric identity²⁷

$$\sin 3\theta = 3 \sin \theta - 4 \sin^3 \theta \quad \text{or,} \quad \sin^3 \theta - \frac{3}{4} \sin \theta + \frac{1}{4} \sin 3\theta = 0 \quad (2.176)$$

and substituting $\lambda = r \sin \theta$ into (2.175) we have²⁸

$$\sin^3 \theta - \frac{\mathcal{I}_2(\boldsymbol{\tau})}{r^2} \sin \theta - \frac{\mathcal{I}_3(\boldsymbol{\tau})}{r^3} = 0 \quad (2.177)$$

Comparing (2.176) and (2.177) gives

$$r = \frac{2}{\sqrt{3}} \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad (2.178)$$

$$\sin 3\theta = -\frac{4\mathcal{I}_3(\boldsymbol{\tau})}{r^3} = -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \quad (2.179)$$

The so-called Lode angle [1423] is then given by

$$\theta_L = \frac{1}{3} \sin^{-1} \left(-\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \quad (2.180)$$

with $-\pi/6 < \theta_L < \pi/6$. The very same equation is also found in Willett (1992) [1359] for instance.

The first root of (2.179) with θ_L determined for $3\theta_L$ in the range $\pm\pi/2$ is a convenient alternative to the third invariant, $\mathcal{I}_3(\boldsymbol{\tau})$. By noting the cyclic nature of $\sin(3\theta_L + 2n\pi)$ we have immediatly the three (and only three) possible values of $\sin \theta_L$ which define the three principal stresses. The deviatoric principal stresses are given by $\lambda = r \sin \theta_L$ on substitution of the three values of $\sin \theta_L$ in turn.

We then obtain

$$\begin{Bmatrix} \tau_1 \\ \tau_2 \\ \tau_3 \end{Bmatrix} = \frac{2}{\sqrt{3}} \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \begin{Bmatrix} \sin(\theta_L + 2\pi/3) \\ \sin \theta_L \\ \sin(\theta_L + 4\pi/3) \end{Bmatrix} \quad (2.181)$$

with $\tau_1 > \tau_2 > \tau_3$ and $-\pi/6 \leq \theta_L \leq \pi/6$. It is indeed easy to verify that for $-\pi/6 \leq \theta_L \leq \pi/6$ we have $\sin(\theta_L + 2\pi/3) > \sin \theta_L > \sin(\theta_L + 4\pi/3)$.

Finally, we wish to compute the principal stresses of the full stress tensor $\boldsymbol{\sigma}$. In the right coordinate system both stress and deviatoric stress tensors are diagonal and $\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau}$ writes:

$$\begin{pmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{pmatrix} = \begin{pmatrix} -p & 0 & 0 \\ 0 & -p & 0 \\ 0 & 0 & -p \end{pmatrix} + \begin{pmatrix} \tau_1 & 0 & 0 \\ 0 & \tau_2 & 0 \\ 0 & 0 & \tau_3 \end{pmatrix}$$

²⁷see section 7.4 of Owen & Hinton [967]

²⁸Note that r and θ have nothing to do with polar, cylindrical or spherical coordinates.

so that (since $p = -\frac{1}{3}\text{tr}(\boldsymbol{\sigma}) = -\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})$)

$$\sigma_1 = \tau_1 - p = \tau_1 + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \quad (2.182)$$

$$\sigma_2 = \tau_2 - p = \tau_2 + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \quad (2.183)$$

$$\sigma_3 = \tau_3 - p = \tau_3 + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \quad (2.184)$$

and finally the total principal stresses are

$$\begin{pmatrix} \sigma_1 \\ \sigma_2 \\ \sigma_3 \end{pmatrix} = \frac{2}{\sqrt{3}}\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \begin{pmatrix} \sin(\theta_L + 2\pi/3) \\ \sin \theta_L \\ \sin(\theta_L + 4\pi/3) \end{pmatrix} + \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad (2.185)$$

with $\sigma_1 > \sigma_2 > \sigma_3$ and $-\pi/6 \leq \theta_L \leq \pi/6$. We have

$$\begin{aligned} \sin(\theta_L + 2\pi/3) &= \sin \theta_L \cos 2\pi/3 + \cos \theta_L \sin 2\pi/3 \\ &= -\frac{1}{2} \sin \theta_L + \cos \theta_L \frac{\sqrt{3}}{2} \end{aligned} \quad (2.186)$$

$$\begin{aligned} \sin(\theta_L + 4\pi/3) &= \sin \theta_L \cos 4\pi/3 + \cos \theta_L \sin 4\pi/3 \\ &= -\frac{1}{2} \sin \theta_L - \cos \theta_L \frac{\sqrt{3}}{2} \end{aligned} \quad (2.187)$$

so that

$$\begin{pmatrix} \sigma_1 \\ \sigma_2 \\ \sigma_3 \end{pmatrix} = \frac{2}{\sqrt{3}}\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \begin{pmatrix} -\frac{1}{2} \sin \theta_L + \cos \theta_L \frac{\sqrt{3}}{2} \\ \sin \theta_L \\ -\frac{1}{2} \sin \theta_L - \cos \theta_L \frac{\sqrt{3}}{2} \end{pmatrix} + \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad (2.188)$$

$$= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \begin{pmatrix} -\frac{1}{\sqrt{3}} \sin \theta_L + \cos \theta_L \\ \frac{2}{\sqrt{3}} \sin \theta_L \\ -\frac{1}{\sqrt{3}} \sin \theta_L - \cos \theta_L \end{pmatrix} + \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad (2.189)$$

Remark. The Lode angle is one of the Lode coordinates²⁹, or Haigh-Westergaard coordinates.

Remark. The Lode angle θ_L is essentially similar to the Lode parameter defined by $-\sqrt{3}\tan \theta$ [967].

Remark. There are 3 different Lode angles, as explained online³⁰:

$$\sin 3\theta_s = -\sin 3\bar{\theta}_s = \cos 3\theta_c = \frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{(\mathcal{I}_2(\boldsymbol{\tau}))^{3/2}}$$

and they are related by $\theta_s = \frac{\pi}{6} - \theta_c$ and $\theta_s = -\bar{\theta}_s$. The one used in this document is in fact the $\bar{\theta}_s$ above.

To recap:

²⁹https://en.wikipedia.org/wiki/Lode_coordinates

³⁰https://en.wikipedia.org/wiki/Lode_coordinates

$$\sigma_1 = \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta_L + \cos \theta_L \right) \quad (2.190)$$

$$\sigma_2 = \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\frac{2}{\sqrt{3}} \sin \theta_L \right) \quad (2.191)$$

$$\sigma_3 = \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta_L - \cos \theta_L \right) \quad (2.192)$$

We will later need $\sigma_1 - \sigma_3$ and $\sigma_1 + \sigma_3$ so we compute these quantities hereafter:

$$\begin{aligned} \sigma_1 - \sigma_3 &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta_L + \cos \theta_L + \frac{1}{\sqrt{3}} \sin \theta_L + \cos \theta_L \right) \\ &= 2 \cos \theta_L \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \end{aligned} \quad (2.193)$$

$$\begin{aligned} \sigma_1 + \sigma_3 &= \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta_L + \cos \theta_L \right) + \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta_L - \cos \theta_L \right) \\ &= \frac{2}{3} \mathcal{I}_1(\boldsymbol{\sigma}) - \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \frac{2}{\sqrt{3}} \sin \theta_L \end{aligned} \quad (2.194)$$

or,

$$\frac{\sigma_1 - \sigma_3}{2} = \cos \theta_L \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad (2.195)$$

$$\frac{\sigma_1 + \sigma_3}{2} = \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) - \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \frac{1}{\sqrt{3}} \sin \theta_L \quad (2.196)$$

Remark. The expression for the Lode angle is different in [1433, p101] than in [1423] or [1431, p62]. They all look suspiciously wrong too.

2.21.3 About the 2nd principal invariant of the deviatoric stress

$$\begin{aligned} \mathcal{K}_2(\boldsymbol{\tau}) &= \frac{1}{2} [\text{Tr}(\boldsymbol{\tau})^2 - \text{Tr}(\boldsymbol{\tau}^2)] \\ &= \frac{1}{2} [(\tau_{xx} + \tau_{yy})^2 - (\tau_{xx}^2 + 2\tau_{xy}^2 + \tau_{yy}^2)] \\ &= \frac{1}{2} [\tau_{xx}^2 + 2\tau_{xx}\tau_{yy} + \tau_{yy}^2 - \tau_{xx}^2 - 2\tau_{xy}^2 - \tau_{yy}^2] \\ &= \frac{1}{2} [2\tau_{xx}\tau_{yy} - 2\tau_{xy}^2] \\ &= \tau_{xx}\tau_{yy} - \tau_{xy}^2 \\ &= \left(\sigma_{xx} - \frac{\sigma_{xx} + \sigma_{yy}}{2} \right) \left(\sigma_{yy} - \frac{\sigma_{xx} + \sigma_{yy}}{2} \right) - \tau_{xy}^2 \\ &= \left(\sigma_{xx} - \frac{\sigma_{xx} + \sigma_{yy}}{2} \right) \left(\sigma_{yy} - \frac{\sigma_{xx} + \sigma_{yy}}{2} \right) - \tau_{xy}^2 \\ &= \left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right) \left(\frac{-\sigma_{xx} + \sigma_{yy}}{2} \right) - \tau_{xy}^2 \\ &= - \left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 - \tau_{xy}^2 \end{aligned}$$

Looking at Eq. (2.166), we can then write

$$\tau_{max} = \sqrt{-\mathcal{K}_2(\boldsymbol{\tau})} = \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2}$$

2.22 Tensor (moment) invariants

physics.tex

There are many different notations used in the literature for invariants and these can prove to be confusing³¹. Note that we only consider symmetric tensors in what follows. Given a tensor \mathbf{T} , one can compute its (moment) invariants as follows (see [1051, p.339], or Appendix A.2 of [1431])

$$\mathcal{I}_1(\mathbf{T}) = \text{tr}[\mathbf{T}] \quad (2.197)$$

$$= T_{xx} + T_{yy} + T_{zz} \quad (2.198)$$

$$\mathcal{I}_2(\mathbf{T}) = \frac{1}{2}\text{tr}[\mathbf{T} \cdot \mathbf{T}] \quad (2.199)$$

$$= \frac{1}{2} \sum_{ij} T_{ij} T_{ji} \quad (2.200)$$

$$= \frac{1}{2}(T_{xx}^2 + T_{yy}^2 + T_{zz}^2) + T_{xy}^2 + T_{xz}^2 + T_{yz}^2 \quad (2.201)$$

$$\mathcal{I}_3(\mathbf{T}) = \frac{1}{3}\text{tr}[\mathbf{T} \cdot \mathbf{T} \cdot \mathbf{T}] \quad (2.202)$$

$$= \frac{1}{3} \sum_i \sum_j \sum_k T_{ij} T_{jk} T_{ki} \quad (2.203)$$

³¹No kidding, true story.

| i | j | k | $T_{ij}T_{jk}T_{ki}$ | $symm$ |
|-----|-----|-----|----------------------|----------------------|
| x | x | x | $T_{xx}T_{xx}T_{xx}$ | T_{xx}^3 |
| y | x | x | $T_{yx}T_{xx}T_{xy}$ | $T_{xx}T_{xy}^2$ |
| z | x | x | $T_{zx}T_{xx}T_{xz}$ | $T_{xx}T_{xz}^2$ |
| x | y | x | $T_{xy}T_{yx}T_{xx}$ | $T_{xx}T_{xy}^2$ |
| y | y | x | $T_{yy}T_{yx}T_{xy}$ | $T_{yy}T_{xy}^2$ |
| z | y | x | $T_{zy}T_{yx}T_{xz}$ | $T_{xy}T_{xz}T_{yz}$ |
| x | z | x | $T_{xz}T_{zx}T_{xx}$ | $T_{xx}T_{xz}^2$ |
| y | z | x | $T_{yz}T_{zx}T_{xy}$ | $T_{xy}T_{xz}T_{yz}$ |
| z | z | x | $T_{zz}T_{zx}T_{xz}$ | $T_{zz}T_{xz}^2$ |
| x | x | y | $T_{xx}T_{xy}T_{yx}$ | $T_{xx}T_{xy}^2$ |
| y | x | y | $T_{yx}T_{xy}T_{yy}$ | $T_{yy}T_{xy}^2$ |
| z | x | y | $T_{zx}T_{xy}T_{yz}$ | $T_{xy}T_{xz}T_{yz}$ |
| x | y | y | $T_{xy}T_{yy}T_{yx}$ | $T_{yy}T_{xy}^2$ |
| y | y | y | $T_{yy}T_{yy}T_{yy}$ | T_{yy}^3 |
| z | y | y | $T_{zy}T_{yy}T_{yz}$ | $T_{yy}T_{yz}^2$ |
| x | z | y | $T_{xz}T_{zy}T_{yx}$ | $T_{xy}T_{xz}T_{yz}$ |
| y | z | y | $T_{yz}T_{zy}T_{yy}$ | $T_{yy}T_{yz}^2$ |
| z | z | y | $T_{zz}T_{zy}T_{yz}$ | $T_{zz}T_{yz}^2$ |
| x | x | z | $T_{xx}T_{xz}T_{zx}$ | $T_{xx}T_{xz}^2$ |
| y | x | z | $T_{yx}T_{xz}T_{zy}$ | $T_{xy}T_{xz}T_{yz}$ |
| z | x | z | $T_{zx}T_{xz}T_{zz}$ | $T_{zz}T_{xz}^2$ |
| x | y | z | $T_{xy}T_{yz}T_{zx}$ | $T_{xy}T_{yz}T_{yz}$ |
| y | y | z | $T_{yy}T_{yz}T_{zy}$ | $T_{yy}T_{yz}^2$ |
| z | y | z | $T_{zy}T_{yz}T_{zz}$ | $T_{zz}T_{yz}^2$ |
| x | z | z | $T_{xz}T_{zz}T_{zx}$ | $T_{zz}T_{xz}^2$ |
| y | z | z | $T_{yz}T_{zz}T_{zy}$ | $T_{zz}T_{yz}^2$ |
| z | z | z | $T_{zz}T_{zz}T_{zz}$ | T_{zz}^3 |

In the end

$$\sum_{i=x,y,z} \sum_{j=x,y,z} \sum_{k=x,y,z} T_{ij}T_{jk}T_{ki} = T_{xx}(T_{xx}^2 + 3T_{xy}^2 + 3T_{xz}^2) + T_{yy}(3T_{xy}^2 + T_{yy}^2 + 3T_{yz}^2) + T_{zz}(3T_{xz}^2 + 3T_{yz}^2 + T_{zz}^2) + 6T_{xy}T_{yz}T_{yz}$$

and then the third moment invariant of the symmetric tensor \mathbf{T} is given by:

$$\begin{aligned}
\mathcal{I}_3(\mathbf{T}) &= \frac{1}{3}T_{xx}(T_{xx}^2 + 3T_{xy}^2 + 3T_{xz}^2) \\
&+ \frac{1}{3}T_{yy}(3T_{xy}^2 + T_{yy}^2 + 3T_{yz}^2) \\
&+ \frac{1}{3}T_{zz}(3T_{xz}^2 + 3T_{yz}^2 + T_{zz}^2) \\
&+ 2T_{xy}T_{xz}T_{yz} \\
&= \frac{1}{3}(T_{xx}^3 + T_{yy}^3 + T_{zz}^3) + T_{xx}(T_{xy}^2 + T_{xz}^2) + T_{yy}(T_{xy}^2 + T_{yz}^2) + T_{zz}(T_{xz}^2 + T_{yz}^2) + 2T_{xy}T_{xz}T_{yz}
\end{aligned} \tag{2.204}$$

2.23 Stress & strain rate invariants

stress_sr_invariants.tex

The implementation of the plasticity criterions relies essentially on the invariants of the (deviatoric) stress $\boldsymbol{\tau}$ and the (deviatoric) strainrate tensors $\dot{\boldsymbol{\epsilon}}$:

$$\mathcal{I}_1(\boldsymbol{\sigma}) = \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \quad (2.206)$$

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2}(\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2) + \tau_{xy}^2 + \tau_{xz}^2 + \tau_{yz}^2 \quad (2.207)$$

$$\begin{aligned} \mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3}\tau_{xx}(\tau_{xx}^2 + 3\tau_{xy}^2 + 3\tau_{xz}^2) \\ &+ \frac{1}{3}\tau_{yy}(3\tau_{xy}^2 + \tau_{yy}^2 + 3\tau_{yz}^2) \\ &+ \frac{1}{3}\tau_{zz}(3\tau_{xz}^2 + 3\tau_{yz}^2 + \tau_{zz}^2) \\ &+ 2\tau_{xy}\tau_{xz}\tau_{yz} \end{aligned} \quad (2.208)$$

and also the second invariant of the deviatoric strain rate is:

$$\begin{aligned} \mathcal{I}_2(\dot{\boldsymbol{\epsilon}}^d) &= \frac{1}{2}[(\dot{\epsilon}_{xx}^d)^2 + (\dot{\epsilon}_{yy}^d)^2 + (\dot{\epsilon}_{zz}^d)^2] + (\dot{\epsilon}_{xy}^d)^2 + (\dot{\epsilon}_{xz}^d)^2 + (\dot{\epsilon}_{yz}^d)^2 \\ &= \frac{1}{6}[(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + (\dot{\epsilon}_{yy} - \dot{\epsilon}_{zz})^2 + (\dot{\epsilon}_{xx} - \dot{\epsilon}_{zz})^2] + \dot{\epsilon}_{xy}^2 + \dot{\epsilon}_{xz}^2 + \dot{\epsilon}_{yz}^2 \end{aligned} \quad (2.209)$$

Proofs of these relationships are given in Appendix T.

We have

$$\begin{aligned} \tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2 &= \left(\sigma_{xx} - \frac{1}{3}I_1\right)^2 + \left(\sigma_{yy} - \frac{1}{3}I_1\right)^2 + \left(\sigma_{zz} - \frac{1}{3}I_1\right)^2 \\ &= \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{2}{3}I_1(\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) + 3\frac{1}{9}I_1^2 \\ &= \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{2}{3}I_1^2 + \frac{1}{3}I_1^2 \\ &= \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{1}{3}I_1^2 \\ &= \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{1}{3}(\sigma_{xx} + \sigma_{yy} + \sigma_{zz})^2 \\ &= \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{1}{3}(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xx}\sigma_{yy} + 2\sigma_{xx}\sigma_{zz} + 2\sigma_{yy}\sigma_{zz}) \\ &= \frac{1}{3}(3\sigma_{xx}^2 + 3\sigma_{yy}^2 + 3\sigma_{zz}^2 - \sigma_{xx}^2 - \sigma_{yy}^2 - \sigma_{zz}^2 - 2\sigma_{xx}\sigma_{yy} - 2\sigma_{xx}\sigma_{zz} - 2\sigma_{yy}\sigma_{zz}) \\ &= \frac{1}{3}(2\sigma_{xx}^2 + 2\sigma_{yy}^2 + 2\sigma_{zz}^2 - 2\sigma_{xx}\sigma_{yy} - 2\sigma_{xx}\sigma_{zz} - 2\sigma_{yy}\sigma_{zz}) \end{aligned} \quad (2.210)$$

$$= \frac{1}{3}((\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{xx} - \sigma_{zz})^2 + (\sigma_{yy} - \sigma_{zz})^2) \quad (2.211)$$

so that

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{6}[(\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{yy} - \sigma_{zz})^2 + (\sigma_{xx} - \sigma_{zz})^2] + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2$$

Remark. $\mathcal{I}_2(\boldsymbol{\tau})$ is often called J_2 or J_2' so that one sometimes speaks of J_2 -plasticity.

These (second) invariants are almost always used under a square root so we define:

$$\tau_e = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad \dot{\epsilon}_e = \sqrt{\mathcal{I}_2(\dot{\boldsymbol{\epsilon}}^d)} \quad (2.212)$$

Note that these quantities have the same dimensions as their tensor counterparts, i.e. Pa for stresses and s⁻¹ for strain rates.

If the stress tensor is such that it is diagonal, i.e.

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{pmatrix} \quad \text{and} \quad \boldsymbol{\tau} = \begin{pmatrix} \tau_1 & 0 & 0 \\ 0 & \tau_2 & 0 \\ 0 & 0 & \tau_3 \end{pmatrix}$$

then the invariants are

$$\begin{aligned} \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_1 + \sigma_2 + \sigma_3 \\ \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{6} [(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_1 - \sigma_3)^2] \end{aligned} \quad (2.213)$$

$$\begin{aligned} \mathcal{I}_3(\boldsymbol{\tau}) &= \tau_1 \tau_2 \tau_3 \\ &= \frac{1}{3} \text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau} \cdot \boldsymbol{\tau}] \\ &= \frac{1}{3} \text{tr} \left[\begin{pmatrix} \tau_1 & 0 & 0 \\ 0 & \tau_2 & 0 \\ 0 & 0 & \tau_3 \end{pmatrix} \cdot \begin{pmatrix} \tau_1 & 0 & 0 \\ 0 & \tau_2 & 0 \\ 0 & 0 & \tau_3 \end{pmatrix} \cdot \begin{pmatrix} \tau_1 & 0 & 0 \\ 0 & \tau_2 & 0 \\ 0 & 0 & \tau_3 \end{pmatrix} \right] \\ &= \frac{1}{3} \text{tr} \begin{pmatrix} \tau_1^3 & 0 & 0 \\ 0 & \tau_2^3 & 0 \\ 0 & 0 & \tau_3^3 \end{pmatrix} \\ &= \frac{1}{3} (\tau_1^3 + \tau_2^3 + \tau_3^3) \\ &= \frac{1}{3} [(\sigma_1 - \mathcal{I}_1(\boldsymbol{\sigma})/3)^3 + (\sigma_2 - \mathcal{I}_1(\boldsymbol{\sigma})/3)^3 + (\sigma_3 - \mathcal{I}_1(\boldsymbol{\sigma})/3)^3] \\ &= \frac{1}{3 \cdot 27} [(3\sigma_1 - \mathcal{I}_1(\boldsymbol{\sigma}))^3 + (3\sigma_2 - \mathcal{I}_1(\boldsymbol{\sigma}))^3 + (3\sigma_3 - \mathcal{I}_1(\boldsymbol{\sigma}))^3] \\ &= \frac{1}{81} [(2\sigma_1 - \sigma_2 - \sigma_3)^3 + (2\sigma_2 - \sigma_1 - \sigma_3)^3 + (2\sigma_3 - \sigma_1 - \sigma_2)^3] \end{aligned} \quad (2.214)$$

The formulation of the third invariant of $\boldsymbol{\tau}$ in Eq. 2.214 is used in Wojciechowski [1367].

2.24 Two-dimensional plane strain calculations

plane_strain.tex

We start from the 3D strain rate tensor

$$\dot{\boldsymbol{\epsilon}}(\vec{v}) = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} & \dot{\epsilon}_{yz} \\ \dot{\epsilon}_{zx} & \dot{\epsilon}_{zy} & \dot{\epsilon}_{zz} \end{pmatrix}$$

The plane strain assumption is such that the problem at hand is assumed to be infinite in a given direction. In the case of computational geodynamics, most 2D modelling is a vertical section of the crust-lithosphere-mantle and the underlying implicit assumption is then that the orogen/rift/subduction/etc ... is infinite in the direction perpendicular to the screen/paper.

Let us assume that the deformation takes place in the x, y -plane, so that $w = 0$ (velocity in the z direction is zero) and $\partial_z \rightarrow 0$ (no change in the z direction). We then have $\dot{\epsilon}_{zz} = 0$ as well as $\dot{\epsilon}_{xz} = 0$ and $\dot{\epsilon}_{yz} = 0$, so that the strain rate tensor is

$$\dot{\boldsymbol{\epsilon}}(\vec{v}) = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & 0 \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Incompressible flow

If the flow is incompressible then the deviatoric stress tensor is given by

$$\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}^d(\vec{\nabla}) = 2\eta \left(\dot{\boldsymbol{\epsilon}}(\vec{\nabla}) - \frac{1}{3} \underbrace{\text{tr}[\dot{\boldsymbol{\epsilon}}]}_{=0} \mathbf{1} \right) = 2\eta\dot{\boldsymbol{\epsilon}}(\vec{\nabla}) = \begin{pmatrix} \tau_{xx} & \tau_{xy} & 0 \\ \tau_{yx} & \tau_{yy} & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

One then discards the unnecessary line and column in the tensor, leaving a 2×2 matrix. Finding the principal stress components is then trivial since we have done it in 2D already.

It is important to keep in mind that the invariants we need to implement the rheologies are $\mathcal{I}_1(\boldsymbol{\sigma})$, $\mathcal{I}_2(\boldsymbol{\tau})$ and $\mathcal{I}_3(\boldsymbol{\tau})$. By formulating our yield surfaces with pressure $p = -\mathcal{I}_1(\boldsymbol{\sigma})/3$ we can then avoid confusion, and since the other two invariants are functions of $\boldsymbol{\tau}$ the pressure term does not pose any problem: simply set τ_{xz} , τ_{yz} and τ_{zz} to zero in the equations of Section 2.23 and we obtain:

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2}(\tau_{xx}^2 + \tau_{yy}^2) + \tau_{xy}^2 \quad (2.215)$$

$$\begin{aligned} \mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3}\tau_{xx}(\tau_{xx}^2 + 3\tau_{xy}^2) + \frac{1}{3}\tau_{yy}(3\tau_{xy}^2 + \tau_{yy}^2) \\ &= \frac{1}{3}(\tau_{xx}^3 + 3\tau_{xx}\tau_{xy}^2 + 3\tau_{yy}\tau_{xy}^2 + \tau_{yy}^3) \\ &= \frac{1}{3}(\tau_{xx}^3 + 3(\tau_{xx} + \tau_{yy})\tau_{xy}^2 + \tau_{yy}^3) \\ &= \frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^3) \quad \text{since } \tau_{ii} = 0 \end{aligned} \quad (2.216)$$

The principal stresses of the deviatoric stress tensor $\boldsymbol{\tau}$ are given by

$$\begin{aligned} \tau_1 &= \frac{\tau_{xx} + \tau_{yy}}{2} + \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \\ \tau_2 &= \frac{\tau_{xx} + \tau_{yy}}{2} - \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \end{aligned} \quad (2.217)$$

The full stress tensor is then

$$\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau} = \begin{pmatrix} -p + \tau_{xx} & \tau_{xy} & 0 \\ \tau_{yx} & -p + \tau_{yy} & 0 \\ 0 & 0 & -p \end{pmatrix}$$

so it remains a 3×3 tensor!

However, looking at the conservation of momentum,

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} = \vec{0}$$

Given the conditions for plane-strain then \vec{g} is likely to be in the xy -plane so that the z component of the equation becomes:

$$-\partial_z p = 0$$

and since we have $\partial_z \rightarrow 0$ anyways this equation is automatically fulfilled. Then, we might as well proceed by considering that the stress tensor is in fact 2D as the third row/column has no incidence. In that case the pressure is given by $p = -\mathcal{I}_1(\boldsymbol{\sigma})/2$. In the plasticity yield criterion or plastic potential we will need the full stress $\boldsymbol{\sigma}$ only via its first invariant (i.e. the pressure). The other two invariants are those of the deviatoric stress.

Let us start from the deviatoric stress tensor:

$$\boldsymbol{\tau} = \boldsymbol{\sigma} - \frac{1}{2}\mathcal{I}_1(\boldsymbol{\sigma}) = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy} \end{pmatrix} - \frac{\sigma_{xx} + \sigma_{yy}}{2} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} (\sigma_{xx} - \sigma_{yy})/2 & \sigma_{xy} \\ \sigma_{xy} & -(\sigma_{xx} - \sigma_{yy})/2 \end{pmatrix}$$

The second invariant of the deviatoric stress tensor is then

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2}\boldsymbol{\tau} : \boldsymbol{\tau} = \frac{1}{2} (2(\sigma_{xx} - \sigma_{yy})^2/4 + 2\sigma_{xy}^2)$$

or

$$\mathcal{I}_2(\boldsymbol{\tau}) = \left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \sigma_{xy}^2$$

and the effective deviatoric stress

$$\tau_e = \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \sigma_{xy}^2}$$

Remark: Using the form of $\mathcal{I}_2(\boldsymbol{\tau})$ above one arrives at

$$\begin{aligned} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xx}} &= 2 \frac{1}{2} \frac{\sigma_{xx} - \sigma_{yy}}{2} = \tau_{xx} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yy}} &= -2 \frac{1}{2} \frac{\sigma_{xx} - \sigma_{yy}}{2} = \tau_{yy} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xy}} &= 2\sigma_{xy} = 2\tau_{xy} \end{aligned}$$

which is ...wrong! One should first write the second invariant for the generic case of the deviatoric stress tensor (without assuming it is symmetric):

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2}\boldsymbol{\tau} : \boldsymbol{\tau} = \frac{1}{2} (2(\sigma_{xx} - \sigma_{yy})^2/4 + \sigma_{xy}^2 + \sigma_{yx}^2) = \left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \frac{1}{2}\sigma_{xy}^2 + \frac{1}{2}\sigma_{yx}^2$$

Then

$$\begin{aligned} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xx}} &= 2 \frac{1}{2} \frac{\sigma_{xx} - \sigma_{yy}}{2} = \tau_{xx} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yy}} &= -2 \frac{1}{2} \frac{\sigma_{xx} - \sigma_{yy}}{2} = \tau_{yy} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xy}} &= \sigma_{xy} = \tau_{xy} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yx}} &= \sigma_{yx} = \tau_{yx} \end{aligned}$$

which can be simply written as

$$\frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} = \boldsymbol{\tau}$$

Compressible flow

If the flow is not incompressible, then the deviatoric strain rate tensor is

$$\dot{\epsilon}^d(\vec{v}) = \dot{\epsilon}(\vec{v}) - \frac{1}{3}\text{tr}[\dot{\epsilon}]\mathbf{1} = \dot{\epsilon}(\vec{v}) - \frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy})\mathbf{1} = \begin{pmatrix} \frac{2}{3}\dot{\epsilon}_{xx} - \frac{1}{3}\dot{\epsilon}_{yy} & \dot{\epsilon}_{xy} & 0 \\ \dot{\epsilon}_{yx} & -\frac{1}{3}\dot{\epsilon}_{xx} + \frac{2}{3}\dot{\epsilon}_{yy} & 0 \\ 0 & 0 & -\frac{1}{3}\dot{\epsilon}_{xx} - \frac{1}{3}\dot{\epsilon}_{yy} \end{pmatrix}$$

The deviatoric stress tensor now has the form

$$\boldsymbol{\tau} = \begin{pmatrix} \tau_{xx} & \tau_{xy} & 0 \\ \tau_{yx} & \tau_{yy} & 0 \\ 0 & 0 & \tau_{zz} \end{pmatrix}$$

We are interested in the principal components of the deviatoric stress tensor $\boldsymbol{\tau}$ so that we now have the following determinant to compute:

$$\begin{vmatrix} \tau_{xx} - \lambda & \tau_{xy} & 0 \\ \tau_{xy} & \tau_{yy} - \lambda & 0 \\ 0 & 0 & \tau_{zz} - \lambda \end{vmatrix} = 0$$

which yields the following characteristic equation:

$$(\tau_{zz} - \lambda)(\lambda - \tau_1)(\lambda - \tau_2) = 0$$

where $\tau_{1,2}$ have previously been obtained in the 2D case:

$$\begin{aligned} \tau_1 &= \frac{\tau_{xx} + \tau_{yy}}{2} + \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \\ \tau_2 &= \frac{\tau_{xx} + \tau_{yy}}{2} - \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \end{aligned} \quad (2.218)$$

We have

$$\begin{aligned} \tau_{xx} + \tau_{yy} &= 2\eta\frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy}) \\ \tau_{xx} - \tau_{yy} &= 2\eta(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy}) \end{aligned} \quad (2.219)$$

Then

$$\begin{aligned} \tau_1 &= \frac{\tau_{xx} + \tau_{yy}}{2} + \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \\ &= \eta\frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy}) + \eta\sqrt{(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + 4\dot{\epsilon}_{xy}^2} \\ \tau_2 &= \frac{\tau_{xx} + \tau_{yy}}{2} - \sqrt{\left(\frac{\tau_{xx} - \tau_{yy}}{2}\right)^2 + \tau_{xy}^2} \\ &= \eta\frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy}) - \eta\sqrt{(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + 4\dot{\epsilon}_{xy}^2} \end{aligned} \quad (2.220)$$

It does not look like it is going to simplify down the road ... Also, the third eigenvalue/principal stress remains and it is not clear whether it is larger or smaller than the other two. The 3D framework is then probably the most appropriate.

Let us now turn to the second invariant of the deviatoric strain rate (see Eq. (2.209)):

$$\mathcal{I}_2(\dot{\epsilon}^d) = \frac{1}{2} \dot{\epsilon}^d : \dot{\epsilon}^d \quad (2.221)$$

$$= \frac{1}{2} [(\dot{\epsilon}_{xx}^d)^2 + (\dot{\epsilon}_{yy}^d)^2 + (\dot{\epsilon}_{zz}^d)^2] + (\dot{\epsilon}_{xy}^d)^2 + (\dot{\epsilon}_{xz}^d)^2 + (\dot{\epsilon}_{yz}^d)^2 \quad (2.222)$$

But there is also an expression for $\mathcal{I}_2(\dot{\epsilon}^d)$ directly as a function of the $\dot{\epsilon}_{ij}$ components (see Eq. (2.209)):

$$\mathcal{I}_2(\dot{\epsilon}^d) = \frac{1}{6} [(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + (\dot{\epsilon}_{yy} - \dot{\epsilon}_{zz})^2 + (\dot{\epsilon}_{xx} - \dot{\epsilon}_{zz})^2] + \dot{\epsilon}_{xy}^2 + \dot{\epsilon}_{xz}^2 + \dot{\epsilon}_{yz}^2 \quad (2.223)$$

$$= \frac{1}{6} [(\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy})^2 + (\dot{\epsilon}_{yy})^2 + (\dot{\epsilon}_{xx})^2] + \dot{\epsilon}_{xy}^2 \quad (2.224)$$

$$= \frac{1}{6} [\dot{\epsilon}_{xx}^2 - 2\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + \dot{\epsilon}_{yy}^2 + \dot{\epsilon}_{yy}^2 + \dot{\epsilon}_{xx}^2] + \dot{\epsilon}_{xy}^2 \quad (2.225)$$

$$= \frac{1}{6} [2\dot{\epsilon}_{xx}^2 - 2\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy}^2] + \dot{\epsilon}_{xy}^2 \quad (2.226)$$

$$= \frac{1}{3} [\dot{\epsilon}_{xx}^2 - \dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + \dot{\epsilon}_{yy}^2] + \dot{\epsilon}_{xy}^2 \quad (2.227)$$

If we now do things the old/wrong(?) way, one would start directly from the 2D strain rate tensor

$$\dot{\epsilon} = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} \end{pmatrix}$$

The deviatoric strain rate tensor is then logically defined as

$$\dot{\epsilon}^d = \dot{\epsilon} - \frac{1}{2} Tr[\dot{\epsilon}] \mathbf{1} = \dot{\epsilon} - \frac{1}{2} (\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy}) \mathbf{1}$$

or,

$$\dot{\epsilon}^d = \begin{pmatrix} \frac{1}{2}\dot{\epsilon}_{xx} - \frac{1}{2}\dot{\epsilon}_{yy} & \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{yx} & -\frac{1}{2}\dot{\epsilon}_{xx} + \frac{1}{2}\dot{\epsilon}_{yy} \end{pmatrix}$$

Let us now turn to the second invariant of the deviatoric strain rate (see Section 3.21 in fieldstone)

$$\begin{aligned} \mathcal{I}_2(\dot{\epsilon}^d) &= \frac{1}{2} \dot{\epsilon}^d : \dot{\epsilon}^d \\ &= \frac{1}{2} [(\frac{1}{2}\dot{\epsilon}_{xx} - \frac{1}{2}\dot{\epsilon}_{yy})^2 + (-\frac{1}{2}\dot{\epsilon}_{xx} + \frac{1}{2}\dot{\epsilon}_{yy})^2] + \dot{\epsilon}_{xy}^2 \\ &= \frac{1}{2} [\frac{1}{4}(2\dot{\epsilon}_{xx}^2 - 4\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy}^2)] + \dot{\epsilon}_{xy}^2 \\ &= \frac{1}{4} [\dot{\epsilon}_{xx}^2 - 2\dot{\epsilon}_{xx}\dot{\epsilon}_{yy} + \dot{\epsilon}_{yy}^2] + \dot{\epsilon}_{xy}^2 \end{aligned} \quad (2.228)$$

which is not the same as the previous expression!

2.25 Alternative principal stresses notations

physics.tex

The principal stresses of the stress tensor σ are σ_1 , σ_2 and σ_3 with $\sigma_1 \geq \sigma_2 \geq \sigma_3$. Following Wojciechowski [1367], we start by stating that the intermediate principal stress can always be represented as a linear combination of two other stresses:

$$\sigma_2 = (1 - b)\sigma_1 + b\sigma_3 \quad \text{where} \quad b = \frac{\sigma_1 - \sigma_2}{\sigma_1 - \sigma_3} \in [0, 1] \quad (2.229)$$

The quantity b is called the principal stress ratio. Let us now introduce the maximum shear plane stresses σ_m and τ_m such that ³²

$$\boxed{\sigma_m = \frac{\sigma_1 + \sigma_3}{2}} \quad \boxed{\tau_m = \frac{\sigma_1 - \sigma_3}{2}} \quad (2.230)$$

so that we have

$$\sigma_1 = \sigma_m + \tau_m \quad (2.231)$$

$$\sigma_2 = \sigma_m - a\tau_m \quad (2.232)$$

$$\sigma_3 = \sigma_m - \tau_m \quad (2.233)$$

The quantity $a \in [-1, 1]$ is an equivalent measure of the principal stress ratio and is defined as

$$a = 2b - 1 = 2 \frac{\sigma_1 - \sigma_2}{\sigma_1 - \sigma_3} - 1 = \frac{\sigma_1 - 2\sigma_2 + \sigma_3}{\sigma_1 - \sigma_3} \quad (2.234)$$

We can introduce a , σ_m and τ_m in the invariants above:

$$\begin{aligned} \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_1 + \sigma_2 + \sigma_3 \\ &= (\sigma_m + \tau_m) + (\sigma_m - a\tau_m) + (\sigma_m - \tau_m) \\ &= 3\sigma_m - a\tau_m \end{aligned} \quad (2.235)$$

$$\begin{aligned} \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{6} [(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_1 - \sigma_3)^2] \\ &= \frac{1}{6} [(\sigma_m + \tau_m - \sigma_m + a\tau_m)^2 + (\sigma_m - a\tau_m - \sigma_m + \tau_m)^2 + (\sigma_m + \tau_m - \sigma_m + \tau_m)^2] \\ &= \frac{1}{6} [(\tau_m + a\tau_m)^2 + (-a\tau_m + \tau_m)^2 + (\tau_m + \tau_m)^2] \\ &= \frac{\tau_m^2}{6} [(1+a)^2 + (-a+1)^2 + 4] \\ &= \frac{\tau_m^2}{6} [1 + 2a + a^2 + 1 - 2a + a^2 + 4] \\ &= \frac{\tau_m^2}{3} (a^2 + 3) \end{aligned} \quad (2.236)$$

Using the definition of the third invariant of Eq. (2.214):

$$\begin{aligned} \mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{81} [(2\sigma_1 - \sigma_2 - \sigma_3)^3 + (2\sigma_2 - \sigma_1 - \sigma_3)^3 + (2\sigma_3 - \sigma_1 - \sigma_2)^3] \\ &= \frac{1}{81} [(2\sigma_m + 2\tau_m - \sigma_m + a\tau_m - \sigma_m + \tau_m)^3 + (2\sigma_m - 2a\tau_m - \sigma_m - \tau_m - \sigma_m + \tau_m)^3 + (2\sigma_m - 2\tau_m - \sigma_m + \tau_m)^3] \\ &= \frac{1}{81} [(2\tau_m + a\tau_m + \tau_m)^3 + (-2a\tau_m - \tau_m + \tau_m)^3 + (-2\tau_m - \tau_m + a\tau_m)^3] \\ &= \frac{\tau_m^3}{81} [(3+a)^3 + (-2a)^3 + (-3+a)^3] \\ &= \frac{\tau_m^3}{81} [27 + 9a + 3a^2 + a^3 - 8a^3 - 27 + 9a - 3a^2 + a^3] \\ &= \frac{\tau_m^3}{81} (18a - 6a^3) \\ &= \frac{2a\tau_m^3}{27} (3 - a^2) \end{aligned}$$

³²Although most of this section is inspired by Wojciechowski [1367], I have decided not to use his notations which are very confusing since he denotes σ_m by p

which is different than Eq. (14) of Wojciechowski [1367]!!

To recap:

$$\boxed{\mathcal{I}_1(\boldsymbol{\sigma}) = 3\sigma_m - a\tau_m} \quad \boxed{\mathcal{I}_2(\boldsymbol{\tau}) = \frac{\tau_m^2}{3} (a^2 + 3)} \quad \boxed{\mathcal{I}_3(\boldsymbol{\tau}) = \frac{2a\tau_m^3}{27} (3 - a^2)} \quad (2.238)$$

Remark. *Wojciechowski [1367] defines the Lode angle as being the opposite of my definition in Eq. 2.180.*

Finally, we can show using Eqs. (2.190,2.191,2.192) that

$$\begin{aligned} a &= \frac{\sigma_1 - 2\sigma_2 + \sigma_3}{\sigma_1 - \sigma_3} \\ &= \frac{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta + \cos \theta \right) - 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\frac{2}{\sqrt{3}} \sin \theta \right) + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta - \cos \theta \right)}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta + \cos \theta \right) - \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\frac{1}{\sqrt{3}} \sin \theta - \cos \theta \right)} \\ &= \frac{\left(-\frac{1}{\sqrt{3}} \sin \theta + \cos \theta \right) - 2 \left(\frac{2}{\sqrt{3}} \sin \theta \right) + \left(-\frac{1}{\sqrt{3}} \sin \theta - \cos \theta \right)}{\left(-\frac{1}{\sqrt{3}} \sin \theta + \cos \theta \right) - \left(-\frac{1}{\sqrt{3}} \sin \theta - \cos \theta \right)} \\ &= \frac{-\frac{6}{\sqrt{3}} \sin \theta}{2 \cos \theta} \\ &= -\frac{3}{\sqrt{3}} \frac{\sin \theta}{\cos \theta} \\ &= -\sqrt{3} \tan \theta \end{aligned} \quad (2.239)$$

Here again we arrive at the opposite of Eq. (16) of Wojciechowski [1367].

2.26 Recap of notations and definitions of stress invariants

recap_invariants.tex

When it comes to stress invariants, I urge the reader to be extremely careful when considering their source(s). As we have seen these come in two main flavours (principal and moment invariants) and they are often written for the full stress or deviatoric tensor. On top of it all, typos are common like in any source and the occasional minus sign or factor 2 or 3 can be missing. This is the reason why I have spent substantial time re-deriving those in the past pages with a consistent set of notations:

| | |
|--------------------------------------------------------------------------|----------------------------------------------------|
| $\boldsymbol{\sigma}$ | (full) stress tensor |
| $\sigma_1, \sigma_2, \sigma_3$ | principal stresses |
| $\boldsymbol{\tau}$ | deviatoric stress tensor |
| τ_1, τ_2, τ_3 | principal deviatoric stresses |
| $\mathcal{I}_1(\boldsymbol{T})$ | first moment invariant of tensor \boldsymbol{T} |
| $\mathcal{I}_2(\boldsymbol{T})$ | second moment invariant of tensor \boldsymbol{T} |
| $\mathcal{I}_3(\boldsymbol{T})$ | third moment invariant of tensor \boldsymbol{T} |
| $\tau_e = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}$ | effective deviatoric stress |
| $\dot{\epsilon}_e = \sqrt{\mathcal{I}_2(\dot{\boldsymbol{\epsilon}}^d)}$ | effective deviatoric strain rate |

Proofs of all the following relationships are given in Appendix T.

- First invariant

$$\begin{aligned}
 \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\
 \mathcal{I}_1(\boldsymbol{\tau}) &= 0 \\
 \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} &= \mathbf{1}
 \end{aligned}$$

- Second invariant

$$\begin{aligned}
 \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \boldsymbol{\tau} : \boldsymbol{\tau} \\
 &= \frac{1}{2} \text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau}] \\
 &= \frac{1}{2} \sum_{ij} \tau_{ij} \tau_{ji} \\
 &= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2 + 2\tau_{xy}^2 + 2\tau_{xz}^2 + 2\tau_{yz}^2) \\
 &= \frac{1}{6} [(\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{yy} - \sigma_{zz})^2 + (\sigma_{xx} - \sigma_{zz})^2] + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
 &= -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma}) \\
 \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} &= \boldsymbol{\tau}
 \end{aligned}$$

- Third invariant

$$\begin{aligned}
\mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3} \sum_{ijk} \tau_{ij} \tau_{jk} \tau_{ki} \\
&= \det(\boldsymbol{\tau}) \\
&= \frac{1}{3} \text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau} \cdot \boldsymbol{\tau}] \\
&= \frac{2}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma}) \\
\frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} &=
\end{aligned} \tag{2.240}$$

$$\theta_L = \frac{1}{3} \sin^{-1} \left(-\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \tag{2.241}$$

$$\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \mathbf{1} \tag{2.242}$$

$$\frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \boldsymbol{\sigma} \tag{2.243}$$

$$\frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \boldsymbol{\sigma} \cdot \boldsymbol{\sigma} \tag{2.244}$$

2.27 Rheology in geodynamics

rheology.tex

The reader is referred to Barnes [45] for a discussion and review of non-linear viscous rheologies and to Coussot [281] for a review of experimental data for yield stress fluid flows. See also Tanner & Tanner [1236] for a summary of Heinrich Hencky's scientific work on rheology.

The Cauchy stress tensor is given by $\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau}$ so that $\mathcal{I}_1(\boldsymbol{\sigma}) = -p\mathcal{I}_1(\mathbf{1}) + \mathcal{I}_1(\boldsymbol{\tau})$. Since $\boldsymbol{\tau}$ is deviatoric, its first invariant is zero. We then have $\mathcal{I}_1(\boldsymbol{\sigma}) = -p n_D$ where n_D is the number of dimensions.

Books:

- Plasticity and Geomechanics, Davis and Selvadurai. [319]
- Elasticity and Geomechanics, Davis and Selvadurai. [318]
- Rheology of the Earth, Ranalli. [1041]
- Deformation of Earth materials, Karato. [674]
- Fundamentals of the Theory of Plasticity, Kachanov. [661]
- Computational methods for plasticity, de Souza Neto et al. [1183]
- Computer simulation of dynamic phenomena, M. wilkins [1357]
- Continuum theory of plasticity, Khan and Huang [698]
- Theory of plasticity, Chakrabarty [217]
- Zienkiewicz Taylor [1431]
- Rheology Principles, Macosko [821]
- Computational Inelasticity, Simo and Hughes [622]
- Lectures on Visco-Plastic Fluid Mechanics, G. Ovarlez & S. Hormozi [966]
- Complex fluids, P. Saramito [1112]

2.27.1 Linear viscous aka Newtonian

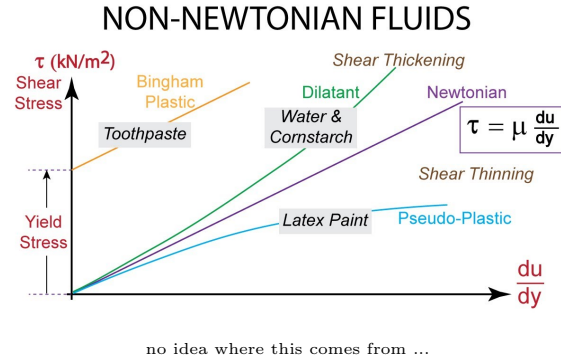
Simply put, a Newtonian fluid is a fluid in which the viscous stresses at every point are linearly proportional to the local strain rate. Mathematically speaking, this means that the fourth-order tensor \mathbf{C} relating the viscous stress tensor to the strain rate tensor does not depend on the stress state and velocity of the flow.

$$\boldsymbol{\tau} = \mathbf{C} : \dot{\boldsymbol{\epsilon}} \quad (2.245)$$

One very often makes the assumption that the fluid is isotropic, i.e. its mechanical properties are the same along any direction. As a consequence the fourth order viscosity tensor \mathbf{C} is symmetric and will have only two independent real parameters: a bulk viscosity coefficient, that defines the resistance of the medium to gradual uniform compression; and a dynamic viscosity coefficient η that expresses its resistance to gradual shearing³³.

³³We here neglect the so-called rotational viscosity coefficient which results from a coupling between the fluid flow and the rotation of the individual particles

Rather logically we denote by non-Newtonian fluids which are not Newtonian, i.e. their viscosity (tensor) depends on stress. Such fluids are part of our daily life, e.g. honey, toothpaste, paint, blood, or shampoo. They are also sometimes denoted as Generalized Newtonian Fluid .



2.27.2 Power-law model

One of the simplest non-Newtonian viscosity model is the power law model, for which the viscosity depends on the (effective) deviatoric strain rate as follows:

$$\eta(\dot{\epsilon}_e) = K\dot{\epsilon}_e^{n-1} \quad \text{or} \quad \sigma = 2K\dot{\epsilon}_e^n \quad (2.246)$$

where n and K are parameters. n is called the power law index. $\dot{\epsilon}_e$ is defined in (2.212) and in the table here above. Note that a Newtonian viscosity is recovered when $n = 1$. Also n and K may depend on temperature (see Reddy [1051, p339]).

A so-called 'generalised' power law rheology is proposed in Iaffaldano & bunge (2009) [618]:

$$\eta = K(\dot{\epsilon}_e + \dot{\epsilon}_0)^{n-1} \quad (2.247)$$

so that in the rigid areas where $\dot{\epsilon}_e \rightarrow 0$ the rheology uses instead a minimum strain rate value $\dot{\epsilon}_0$.

📖 Relevant Literature: England & Molnar (1997) [377]

2.27.3 Carreau model

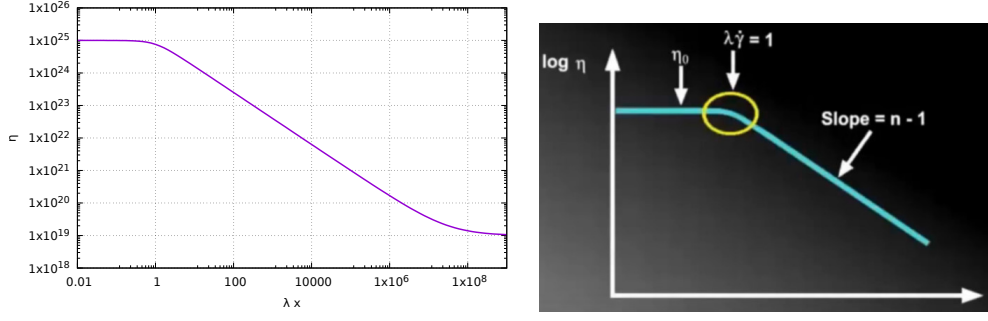
Note that this model is sometimes called Bird-Carreau in the literature. As explained in Reddy [1051], the power-law model poses no restriction on how small or large the viscosity may become, which may prove problematic once implemented as it can lead to runaway effects (strain rate becomes large \rightarrow viscosity becomes smaller \rightarrow strain rate becomes larger, etc ...). This problem is alleviated in the so-called Carreau ³⁴ model [210] (see for example Zinani & Frey (2007) [1439]). The viscosity is then given by

$$\eta(\dot{\epsilon}_e) = \eta_\infty + (\eta_0 - \eta_\infty) \left(1 + (\lambda\dot{\epsilon}_e)^2\right)^{(n-1)/2} \quad (2.248)$$

where η_0 , η_∞ , λ and $n \in [0, 1]$ are material parameters. λ is called the relaxation time: it is the inverse of the shear rate at which the fluid changes from Newtonian to power-law behavior.

At low strain rate a Carreau fluid behaves as a Newtonian fluid with viscosity η_0 . At intermediate strain rates $\dot{\epsilon}_e\lambda \sim 1$ a Carreau fluid behaves as a Power-law fluid. At high strain rate, a Carreau fluid behaves as a Newtonian fluid again with viscosity η_∞ .

³⁴https://en.wikipedia.org/wiki/Carreau_fluid



Left: Carreau model effective viscosity as a function of the product $\lambda \dot{\epsilon}_e$. Right: taken from video at <https://youtu.be/qErs5zZV4BQ>.

Note that the (Bird)-Carreau-Yasuda model [1382, 963] is very similar to the standard (Bird)-Carreau:

$$\eta = \eta_\infty + (\eta_0 - \eta_\infty) (1 + (\lambda \dot{\epsilon}_e)^a)^{(n-1)/a} \quad (2.249)$$

It is for instance used in van de Vosse *et al.* (2003) [1303] to model blood.

Flows in a Lid-Driven Cavity with this rheology are presented in [1439, 1153].

 **Relevant Literature:** Bercovici (1993) [76], Bercovici (1995) [77], Marcotte (2000) [835], Huerta & Liu (1988) [603].

2.27.4 Bingham model

Bingham [88] fluids can sustain an applied stress without any motion occurring. Only when the applied stress exceeds a yield stress τ_0 then the fluid flows. This translates as follows [1051]:

$$\boldsymbol{\tau} = \left(\frac{\tau_0}{\dot{\epsilon}} + 2\eta_0 \right) \dot{\epsilon}^d \quad \text{if } \tau_e > \tau_0 \quad (2.250)$$

$$\boldsymbol{\tau} = \mathbf{0} \quad \text{if } \tau_e \leq \tau_0 \quad (2.251)$$

When flow occurs, the effective viscosity is then given by:

$$\eta(\dot{\epsilon}_e) = \frac{\tau_0}{\dot{\epsilon}_e} + 2\eta_0 \quad (2.252)$$

and when the strain rate is large we recover a Newtonian behaviour. Typical Bingham fluids are mud, slurry, toothpaste.


When using a velocity-based FEM code, the implementation of this rheological behaviour is complicated by the no-flow condition under a given stress. However, our codes require a relationship between stress and strain rate in the form of an effective viscosity which cannot be zero. This difficulty can be circumvented by implementing Bingham fluids as follows [1051]:

$$\boldsymbol{\tau} = \left(\frac{\tau_0(1 - \eta/\eta_r)}{\dot{\epsilon}_e} + 2\eta_0 \right) \dot{\epsilon} \quad \text{if } \tau_e > \tau_0 \quad (2.253)$$

$$\boldsymbol{\tau} = 2\eta_r \dot{\epsilon} \quad \text{if } \tau_e \leq \tau_0 \quad (2.254)$$

where η_r is a pre-yield viscosity and $\eta/\eta_r \ll 1$ (typically 1% or less). This is a form of regularisation, and we will see a similar one in the next section.

Note the interesting paper by Barnes and Walter (1985) [46] who argue that "the yield stress concept is an idealization, and that, given accurate measurements, no yield stress exists. The simple Cross model is shown to be a useful empiricism for many non-Newtonian fluids, including those which have hitherto been thought to possess a yield stress." The Cross model is presented in Section 2.27.8.

 **Relevant Literature:** Papanastasiou (1987) [973], Blackery & Mitsoulis (1997) [94], Mitsoulis & Zisis (2001) [884], Mahmood *et al.* (2017) [824], Syrakos *et al.* (2014) [1223], Bingham [88], Balmforth

& Rust (2009) [41], Grinevich & Olshanskii (2009) [497], Sverdrup *et al.* (2018) [1221] FE method for incompressible non-Newtonian flow (Bercovier & Engelman (1980) [80]); Flow around a rigid sphere (Liu *et al.* (2002) [794]), Conduit flow of an incompressible, yield-stress fluid, Taylor and Wilson [1241] (1997).

2.27.5 Herschel-Bulkley visco-plastic model

The Herschel-Bulkley model is effectively a combination of the power-law model and a simple plastic model:

$$\boldsymbol{\tau} = 2 \left(K \dot{\epsilon}_e^{n-1} + \frac{\tau_0}{\dot{\epsilon}} \right) \dot{\epsilon} \quad \text{if } \tau_e > \tau_0 \quad (2.255)$$

$$\dot{\epsilon} = \mathbf{0} \quad \text{if } \tau_e \leq \tau_0 \quad (2.256)$$

in which τ_0 is the yield stress, K the consistency, and n is the flow index [85]. The flow index measures the degree to which the fluid is shear-thinning ($n < 1$) or shear-thickening ($n > 1$). If $n = 1$ and $\tau_0 = 0$ the model reduces to the Newtonian model.

The term between parenthesis above is the nonlinear effective viscosity. Concretely, the implementation goes as follows³⁵:

$$\eta(\dot{\epsilon}) = \begin{cases} \eta_0 & \dot{\epsilon}_e \leq \dot{\epsilon}_0 \\ K \dot{\epsilon}_e^{n-1} + \frac{\tau_0}{\dot{\epsilon}_e} & \dot{\epsilon}_e \geq \dot{\epsilon}_0 \end{cases} \quad (2.257)$$

The limiting viscosity η_0 is chosen such that $\eta_0 = K \dot{\epsilon}_0^{n-1} + \frac{\tau_0}{\dot{\epsilon}_0}$

A large limiting viscosity means that the fluid will only flow in response to a large applied force. This feature captures the Bingham-type behaviour of the fluid. Note that when strain rates are large, the power-law behavior dominates.

As we have seen for Bingham fluids, the equations above are not easily amenable to implementation so that one usually resorts to regularisation, which is a modification of the equations by introducing a new material parameter which controls the exponential growth of stress. This way the equation is valid for both yielded and unyielded areas (Blackery & Mitsoulis (1997) [94], Papanastasiou (1987) [973], Zinani & Frey (2007) [1439], Sverdrup *et al.* (2018) [1221]):

$$\eta(\dot{\epsilon}_e) = K \dot{\epsilon}_e^{n-1} + \frac{\tau_0}{\dot{\epsilon}_e} [1 - \exp(-m \dot{\epsilon}_e)] \quad (2.258)$$

When the strain rate becomes (very) small a Taylor expansion of the regularisation term yields $1 - \exp(-m \dot{\epsilon}) \sim m \dot{\epsilon}$ so that $\eta_{eff} \rightarrow m \tau_0$. However, it seems more physically meaningful to replace m by a reference strain rate value $\dot{\epsilon}_0$ so that

$$\eta_{eff}(\dot{\epsilon}) = K \dot{\epsilon}_e^{n-1} + \frac{\tau_0}{\dot{\epsilon}_e} \left[1 - \exp \left(-\frac{\dot{\epsilon}_e}{\dot{\epsilon}_0} \right) \right] \quad (2.259)$$

In this case, when strain rate becomes (very) small a Taylor expansion of the regularisation term yields

$$\frac{\tau_0}{\dot{\epsilon}_e} \left[1 - \exp \left(-\frac{\dot{\epsilon}_e}{\dot{\epsilon}_0} \right) \right] \simeq \frac{\tau_0}{\dot{\epsilon}_e} \frac{\dot{\epsilon}_e}{\dot{\epsilon}_0} = \frac{\tau_0}{\dot{\epsilon}_0} \quad (2.260)$$

This has the dimensions of a viscosity and this is effectively the definition of a maximum viscosity η_{max} .

 **Relevant Literature:**

- Viscous flow with large free surface motion (Huerta & Liu (1988) [603]);

³⁵https://en.wikipedia.org/wiki/Herschel-Bulkley_fluid

- Numerical simulation of thermal plumes (Massmeyer *et al.* [840]);
- Flows Through a Sudden Axisymmetric Expansion (Machado *et al.* [818], Jay *et al.* (2001) [640]);
- Dam break problem (Ancey & Cochard (2009) [17], Cochard & Ancey (2009) [262], Balmforth *et al.* [40];
- Weakly compressible Poiseuille flow (Taliadorou (2009) [1231]);
- Flow past cylinders in tubes (Mitsoulis & Galazoulas (2009) [883]);
- Determination of yield surfaces (Burgos & Alexandrou (1999) [180]);
- Carbopol hydrogel rheology for experimental tectonics and geodynamics Di Giuseppe *et al.* [333] (2015);
- Flow past a sphere(disc) Besses, Magnin, and Jay [85] (2004); Gavrilov *et al.* (2017) [442]).
- Progress in numerical simulation of yield stress fluid flows, Saramito and Wachs [1113] (2017);
- elastoviscoplastic model based on the Herschel–Bulkley viscoplastic model Saramito [1111] (2019).

2.27.6 The Casson model

It is described in Barnes (1999) [45]:

$$\sqrt{\sigma} = \sqrt{\sigma_y} + \sqrt{\eta_p \dot{\epsilon}_e} \quad (2.261)$$

or, when squaring it:

$$\sigma = \sigma_y + \eta_p \dot{\epsilon}_e + 2\sqrt{\sigma_p \eta_p \dot{\epsilon}_e} \quad (2.262)$$

This model has been found to accurately describe the behaviour of synthetic based muds [2]. See also Section 2.5.1 of Macosko [821].

2.27.7 The Ellis model

An Ellis equation would be of the form [1077]

$$\frac{\eta - \eta_\infty}{\eta_0 - \eta_\infty} = \frac{1}{1 + (\sigma/\sigma_c)^m} \quad (2.263)$$

where σ is the shear stress, σ_c is a critical shear stress and m is a large number. See also Section 2.4.3 of Macosko [821].

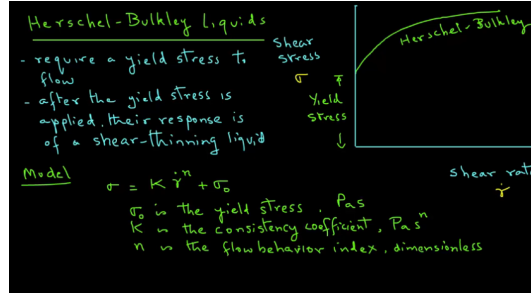
2.27.8 One model to rule them all?

Let us consider the base equation

$$\boxed{\frac{\eta - \eta_\infty}{\eta_0 - \eta_\infty} = [1 + (K \dot{\epsilon}_e)^a]^{-(1-n)/a}} \quad (2.264)$$

This equation is purposefully generic and specific parameter combination choices allow to recover any of the above models (and more) [963]. See also an early paper by Cross (1965) [288] for a somewhat similar equation. See also Section 2.4.2 of Macosko [821].

Similar conclusions are reached in the following video:



<https://youtu.be/dVCb11dZR7Y>

| Model | 1D viscosity | 3D generalization |
|----------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Newtonian | $\mu_0 = \text{constant}$ | $\sigma' = 2\mu_0 D'$ |
| Power law | $\mu = m\dot{\gamma}^{n-1}$ | $\sigma' = 2m(\sqrt{2 \text{tr}(D'^2)})^{n-1} D'$ |
| Truncated Power Law | $\mu = \mu_0 (\dot{\gamma}/\dot{\gamma}_0)^{n-1} \quad \dot{\gamma} \leq \dot{\gamma}_0$ $\mu = \mu_0 (\dot{\gamma}/\dot{\gamma}_0)^{n-1} \quad \dot{\gamma} \geq \dot{\gamma}_0$ | $\sigma' = 2\mu_0 D' \frac{\sqrt{2 \text{tr}(D'^2)}}{(\dot{\gamma}_0)^{n-1}} \quad D' \leq \dot{\gamma}_0$ $\sigma' = 2\mu_0 (\sqrt{2 \text{tr}(D'^2)})^{n-1} D' \quad D' \geq \dot{\gamma}_0$ |
| Carreau | $\frac{\mu - \mu_\infty}{\mu_0 - \mu_\infty} = [1 + (\lambda\dot{\gamma})^2]^{(n-1)/2}$ | $\sigma' = 2(\mu_0 - \mu_\infty) \left[(1 + 2\lambda^2 \text{tr}(D'^2))^{(n-1)/2} D' + 2\mu_\infty D' \right]$ |
| Carreau-A | $\mu_\infty = 0$ | $\sigma' = 2\mu_0 [1 + 2\lambda^2 \text{tr}(D'^2)]^{(n-1)/2} D'$ |
| Bingham | $\mu = \infty \quad \tau \leq \tau_0$ $\mu = \mu_0 + \tau_0/\dot{\gamma} \quad \tau \geq \tau_0$ | $D' = 0 \quad \frac{1}{2} \text{tr}(\sigma'^2) \leq \tau_0^2$ $\sigma' = 2\mu_0 [1 + \tau_0/\sqrt{2 \text{tr}(D'^2)}] D' \quad D' \geq \tau_0^2$ |
| Herschel and Bulkley | $\mu = \infty \quad \tau \leq \tau_0$ $\mu = m\dot{\gamma}^{n-1} + \tau_0/\dot{\gamma} \quad \tau \geq \tau_0$ | $D' = 0 \quad \frac{1}{2} \text{tr}(\sigma'^2) \leq \tau_0^2$ $\sigma' = 2m(\sqrt{2 \text{tr}(D'^2)})^{n-1} D' + 2\tau_0/\sqrt{2 \text{tr}(D'^2)} D' \quad \frac{1}{2} \text{tr}(\sigma'^2) \geq \tau_0^2$ |

where σ' and D' are the deviatoric part of the stress and stretch [i.e. $\frac{1}{2}(\nabla\mathbf{v} + \nabla\mathbf{v}^T)$] tensors, respectively

Taken from Huerta and Liu [603] (1988)

2.27.9 Dislocation and Diffusion creep

insert here background and links to relevant textbooks

The standard dislocation creep effective viscosity is given by:

$$\eta^{ds}(p, T, \dot{\epsilon}) = \eta^{ds}(p, T, \dot{\epsilon}_e) = \frac{1}{2} f A^{-1/n} \dot{\epsilon}_e^{(1-n)/n} \exp\left(\frac{Q + pV}{nRT}\right)$$

where A is the pre-exponential scaling factor, f is a scaling factor representing viscous weakening or strengthening, Q is the activation energy, V is the activation volume, T is the absolute temperature, n is the power-law exponent, R is the universal gas constant.

The coefficients A , n , Q , V are material parameters and are obtained in the laboratory by means of high pressure/temperature experiments (see for instance Karato & Wu (1993) [673]). Unfortunately these experiments cannot be run at Earth-like strain rate values ($\sim 10^{-15}\text{s}^{-1}$) so that extrapolations must be carried out over several orders of magnitude to arrive at values we can use in our numerical models. The 1/2 factor arises from the relationship between deviatoric stress and strain rate which involves a factor 2.

The factor f is in fact a tuning parameter used to explore end members (e.g. 'weak crust' vs 'strong crust'), see discussion in the supplementary material in Huismans & Beaumont (2011) [612]. This approach has been extensively used by the SOPALE users community, see for instance Warren *et al.* (2008) [1343, 1344, 1345] or Gray & Pysklywec (2012) [482].

insert here equation for diffusion creep


Furthermore, we know that several other factors will strongly affect the rheology:

- water content, or as often mentioned: 'dry' vs 'wet'. Following [673], dry means water-free and wet means water-saturated conditions.

| Mechanism | Dry | Wet |
|----------------------------|----------------------|----------------------|
| Dislocation creep | | |
| A (s^{-1}) | 3.5×10^{22} | 2.0×10^{18} |
| n | 3.5 | 3.0 |
| m | 0 | 0 |
| E^* ($kJ\ mol^{-1}$) | 540 | 430 |
| V^* ($cm^3\ mol^{-1}$) | 15 to 25† | 10 to 20† |
| Diffusion creep | | |
| A (s^{-1}) | 8.7×10^{15} | 5.3×10^{15} |
| n | 1.0 | 1.0 |
| m | 2.5‡ | 2.5‡ |
| E^* ($kJ\ mol^{-1}$) | 300 | 240 |
| V^* ($cm^3\ mol^{-1}$) | 6§ | 5§ |

†The activation volume for dislocation creep is not well constrained. Values from $13\ cm^3\ mol^{-1}$ at wet conditions to $27\ cm^3\ mol^{-1}$ at nearly dry conditions have been reported (33). Considering this uncertainty, we used a range of activation volumes, 15 to $25\ cm^3\ mol^{-1}$ for dry olivine and 10 to $20\ cm^3\ mol^{-1}$ for wet olivine. ‡The grain-size exponent is reported as 2 for dry olivine and 3 for wet olivine (16). However, during extrapolation to a large grain size, the grain-size exponent may change from 2 to 3 or from 3 to 2 (16). Taking this uncertainty into account, we chose a grain-size exponent of 2.5, and the preexponential factors were modified accordingly. §The activation volume for diffusion is determined as $6\ cm^3\ mol^{-1}$ at dry conditions (26). No data are available for the activation volume for diffusion at wet conditions; we assumed that this value is ~80% of the activation volume at dry conditions.

Taken from Karato and Wu [673].

 **Relevant Literature:** Quinquis & Buiter (2014) [1027] and refs therein for the effects of water migration on models of subduction dynamics.

- composition: while one typically assigns olivine properties to the mantle in models, the mineral olivine³⁶ is actually a magnesium iron silicate with the formula $(Mg^{2+}, Fe^{2+})_2SiO_4$. and the ratio of magnesium to iron varies between the two endmembers of the solid solution series: forsterite (Mg-endmember: Mg_2SiO_4) and fayalite (Fe-endmember: Fe_2SiO_4).
- grain size: this only affects diffusion creep mechanisms [673]. Grain size varies over several orders of magnitude and also evolves over time and its evolution is affected by the ambient deformation and the deformation history. Dannberg *et al.* [301] then used a diffusion creep effective viscosity given by:

$$\eta^{df} = \frac{1}{2} A_{df}^{-1} d^m \exp \left(\frac{Q_{df} + pV_{df}}{RT} \right)$$

where d is the (variable) grain size and m the grain size exponent. Grain growth/evolution is usually approximated using semi-empirical expressions [301, section 2.2]. Smaller grains facilitating faster creep.

Relevant literature on this topic is in Section ??.

- anisotropy, LPO: see relevant literature in Section ??.
- phase changes

Remark. *It is not uncommon to find in the literature effective viscosity formulations written as a function of B with $B = A^{-1/n}$ [1343, 1344, 1345]. Also, this B coefficient often contains the conversion factor of the next remark.*

Remark. *Material parameters obtained in the lab are often measured on a uniaxial machine. An additional coefficient is added to the effective viscosity formula (see [482, 483], or Table 1a of Warren *et al.* (2008) [1343]): $3^{-(1+n)/2n} 2^{(1-n)/n}$. See page 77 of Ranalli [1041] for an explanation.*

Remark. *In Tullis, Horowitz, and Tullis [1287] it is explained how to formulate a flow law for a polyphase aggregate from end-member flow laws.*

³⁶<https://en.wikipedia.org/wiki/Olivine>

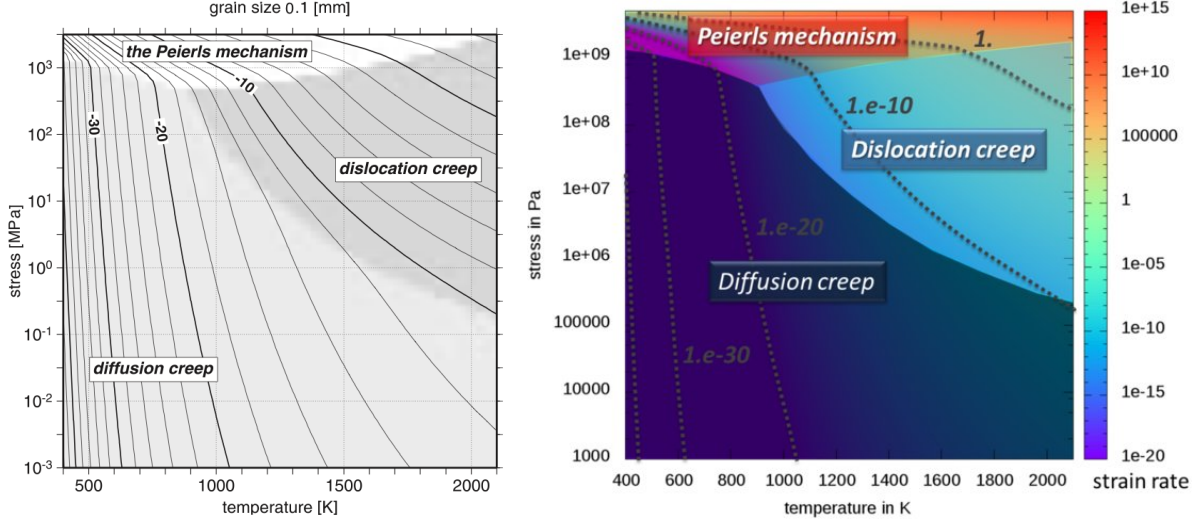


Fig-B: Left: Taken from Kameyama *et al.* (1999) [666]. Deformation mechanism map calculated for grain size $a = 0.1\text{mm}$. The lightly shaded area indicates that deformation mainly occurs by diffusion creep. The densely shaded area indicates that deformation mainly occurs by power-law creep. The white region indicates that deformation mainly occurs by the Peierls mechanism. The solid curves are lines of constant strain rate. The numbers attached to each contour indicate the logarithm of the strain rate in the unit of s^{-1} . Right: Taken from Elbeshhausen & Melosh (1998) [366]. Strain rate as a function of stress and temperature. The parameter space dominated by each of the terms is denoted by different hues, depending on strain rate. Note that the original equations implicitly include a pressure dependence in the enthalpy term, which is not strong and therefore neglected here. The diffusion regime is highly temperature dependent and important for small strain rates only. Dislocation creep occurs for intermediate stresses and higher temperatures, while the Peierls mechanism dominates at higher stresses ($\geq 500\text{MPa}$) and shows a strong stress dependence for low temperatures. The dashed contours are strain rate in units of s^{-1} .

A closer look at the diffusion creep of Karato & Wu (1993) In the article, the following equation is used:

$$\dot{\epsilon} = A \left(\frac{\tau}{\mu} \right) \left(\frac{b}{d} \right)^m \exp \left(-\frac{Q + pV}{RT} \right)$$

where μ is the shear modulus ($\sim 80\text{GPa}$), b is the length of the Burgers vector ($\sim 0.5\text{nm}$) and d is the grain size. One can express the above equation in terms of second invariants (see Section 2.22):

$$\dot{\epsilon}_e = A \left(\frac{\tau_e}{\mu} \right) \left(\frac{b}{d} \right)^m \exp \left(-\frac{Q + pV}{RT} \right)$$

and assuming a Newtonian linearisation/relation between deviatoric stress and strain rate $\tau_e = 2\eta^{df}\dot{\epsilon}_e$, one arrive at

$$\eta^{df} = \frac{1}{2} \left(\frac{A}{\mu} \right)^{-1} \left(\frac{b}{d} \right)^{-m} \exp \left(\frac{Q + pV}{RT} \right)$$

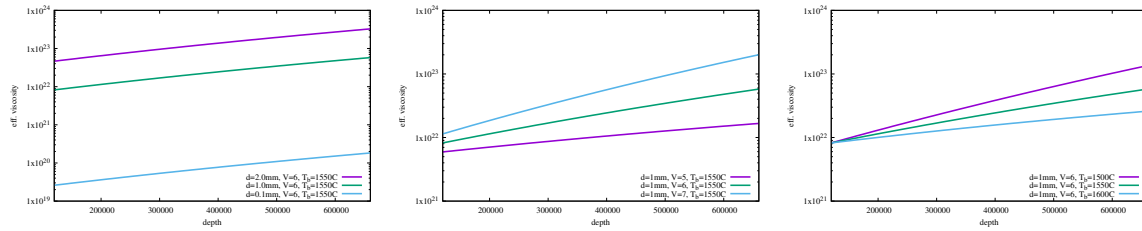
or,

$$\eta^{df} = \frac{1}{2} \left[\frac{A}{\mu} \left(\frac{b}{d} \right)^m \right]^{-1} \exp \left(\frac{Q + pV}{RT} \right)$$

The effective diffusion creep viscosity is independent of strain-rate so that one could substitute the total pressure for lithostatic pressure in the equation, assume a geotherm and easily compute the predicted viscosity as a function of grain size d .

Let us assume that the 1D profile starts from the base of the lithosphere (say 120km depth) and ends at the 660 boundary. Assume the temperature to increase linearly from 1300°C to T_{bottom} (to be specified). At the bottom of the lithosphere, the lithostatic pressure is of the order of $\rho \cdot g \cdot L \simeq 3000 \cdot 10 \cdot 120e3 \simeq 4\text{GPa}$. At the bottom of the domain, the pressure has increased by $3300 \cdot 10 \cdot 630e3 \simeq 21\text{GPa}$.


The viscosity profile is plotted hereunder for three different grain sizes, bottom temperature and activation volumes ($4, 5, 6 \text{ cm}^3/\text{mol}$).



Effective diffusion creep viscosity for various grain size, activation volume and basal temperature values.

images/rheology/kawudiff/

Although this exercise only provides us with first-order results, we can conclude that one can essentially change the diffusion creep effective viscosity by up to 2 orders of magnitude simply by choosing key parameters within acceptable ranges.

 **Relevant Literature:** Dixon and Durham [335] “Measurement of activation volume for creep of dry olivine at upper-mantle conditions”

2.27.10 The von Mises failure criterion

vMcriterion.tex

The von Mises yield criterion suggests that the yielding of materials begins when the second deviatoric stress invariant $\mathcal{I}_2(\boldsymbol{\tau})$ reaches a critical value. For this reason, it is sometimes called the J_2 -plasticity or J_2 flow theory³⁷. It is part of a plasticity theory that applies best to ductile materials, such as metals.

In material science and engineering the von Mises yield criterion can be also formulated in terms of the von Mises stress or equivalent tensile stress, σ_v , a scalar stress value that can be computed from the stress tensor. In this case, a material is said to start yielding when its von Mises stress reaches a critical value known as the yield strength, σ_Y . The von Mises stress is used to predict yielding of materials under any loading condition from results of simple uniaxial tensile tests. The von Mises stress satisfies the property that two stress states with equal distortion energy have equal von Mises stress.

Because the von Mises yield criterion is independent of the first stress invariant, $\mathcal{I}_1(\boldsymbol{\sigma})$, it is applicable for the analysis of plastic deformation for ductile materials such as metals, as the onset of yield for these materials does not depend on the hydrostatic component of the stress tensor.

Although formulated by Maxwell in 1865, it is generally attributed to von Mises [876]. Huber (1904), in a paper in Polish, anticipated to some extent this criterion. Heinrich Hencky formulated the same criterion as von Mises independently in 1924 [562, 1236]. This criterion is also referred to as the Maxwell-Huber-Hencky-von Mises theory.

The von Mises yield criterion (also known as Prandtl-Reuss yield criterion) is expressed in the principal stresses as

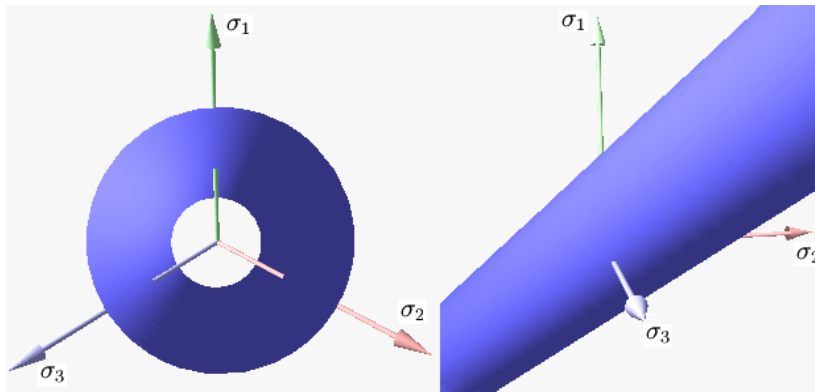
$$\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} = c \quad \text{or,} \quad \frac{1}{6}[(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_3 - \sigma_1)^2] = c^2$$

where c is the yield stress in uniaxial tension. The von Mises yield criterion writes:

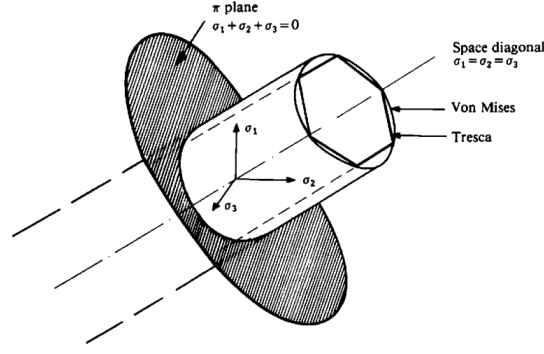
$$F^{\text{VM}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - c \quad (2.265)$$

which is the Drucker-Prager criterion with $\phi = 0$ (see Section 2.27.13).

The following figure shows the von Mises yield surface in the three-dimensional space of principal stresses.




³⁷ J_2 is the common notation for $\mathcal{I}_2(\boldsymbol{\tau})$



Left: Taken from https://en.wikipedia.org/wiki/Yield_surface. Right: Taken from Owen and Hinton [967].

It is a circular cylinder of infinite length with its axis inclined at equal angles to the three principal stresses.

 **Relevant Literature:** Tasos C Papanastasiou. “Flows of materials with yield”. In: *Journal of Rheology* 31.5 (1987), pp. 385–404, F. Tin-Loi and N.S. Ngo. “Performance of the p-version finite element method for limit analysis”. In: *International Journal of Mechanical Sciences* 45 (2003), pp. 1149–1166. DOI: 10.1016/j.ijmecsci.2003.08.004

The yield surface Let us try to draw the yield function in the space $\sigma_1, \sigma_2, \sigma_3$. It is given by

$$\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} = c \quad (2.266)$$

$$\Rightarrow \mathcal{I}_2(\boldsymbol{\tau}) = c^2 \quad (2.267)$$

$$\Rightarrow \frac{1}{6} [(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_1 - \sigma_3)^2] = c^2 \quad (2.268)$$

$$\Rightarrow (\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_1 - \sigma_3)^2 = 6c^2 \quad (2.269)$$

or, temporarily setting $x = \sigma_1$, $y = \sigma_2$ and $z = \sigma_3$:

$$(x - y)^2 + (y - z)^2 + (x - z)^2 = 6c^2 \quad (2.270)$$

$$(x - y)^2 + y^2 - 2yz + z^2 + x^2 - 2xz + z^2 = 6c^2 \quad (2.271)$$

$$2z^2 - 2(x + y)z + (x - y)^2 + x^2 + y^2 - 6c^2 = 0 \quad (2.272)$$

This is a second order polynomial in z . Its discriminant Δ is

$$\begin{aligned} \Delta &= 4(x + y)^2 - 4 \cdot 2 \cdot [(x - y)^2 + x^2 + y^2 - 6c^2] \\ &= 4x^2 + 8xy + 4y^2 - 8[x^2 - 2xy + y^2 + x^2 + y^2 - 6c^2] \\ &= 4x^2 + 8xy + 4y^2 - 8[2x^2 - 2xy + 2y^2 - 6c^2] \\ &= 4x^2 + 8xy + 4y^2 - 16x^2 + 16xy - 16y^2 + 48c^2 \\ &= -12x^2 + 24xy - 12y^2 + 48c^2 \\ &= -12(x^2 - 2xy + y^2) + 48c^2 \\ &= -12(x - y)^2 + 48c^2 \end{aligned}$$

Since I am looking for $z(x, y) \in \mathbb{R}$ then $\Delta > 0$ and this imposes a restriction on admissible x, y pairs:

$$-12(x - y)^2 + 48c^2 > 0$$

$$(x - y)^2 < 4c^2$$

$$x - y < 2c \quad \text{or,} \quad y - x < 2c$$

$$y > x - 2c \quad \text{or,} \quad y < x + 2c$$

So the discriminant is positive in the band given by $y > x - 2c$ and $y < x + 2c$ in the x, y -plane, which is a band centered around the line $y = x$. When $\Delta > 0$ we have then

$$z = \frac{2(x + y) \pm \sqrt{\Delta}}{4}$$

which means that for each pair x, y there are 2 z values. The middle of this surface is given by the line $z = (x + y)/2$. The plane normal to this line is given by $z = -2(x + y)$.

This approach is reasonably simple for the von Mises criterion but quickly becomes intractable for other criteria.

We now look into the derivatives of the von Mises plastic potential $Q^{\text{vM}}(\boldsymbol{\sigma})$. We have

$$Q^{\text{vM}}(\boldsymbol{\sigma}) = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - c \quad (2.273)$$

Then

$$\frac{\partial Q^{\text{vM}}}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} = 0 \quad (2.274)$$

$$\frac{\partial Q^{\text{vM}}}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} = 1 \quad (2.275)$$

$$\frac{\partial Q^{\text{vM}}}{\partial \theta_L(\boldsymbol{\tau})} = 0 \quad (2.276)$$

so

$$C_1^{\text{vM}} = 0 \quad (2.277)$$

$$C_2^{\text{vM}} = \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}(1 - 0) = \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \quad (2.278)$$

$$C_3^{\text{vM}} = 0 \quad (2.279)$$

2.27.11 The Tresca failure criterion

trcriterion.tex

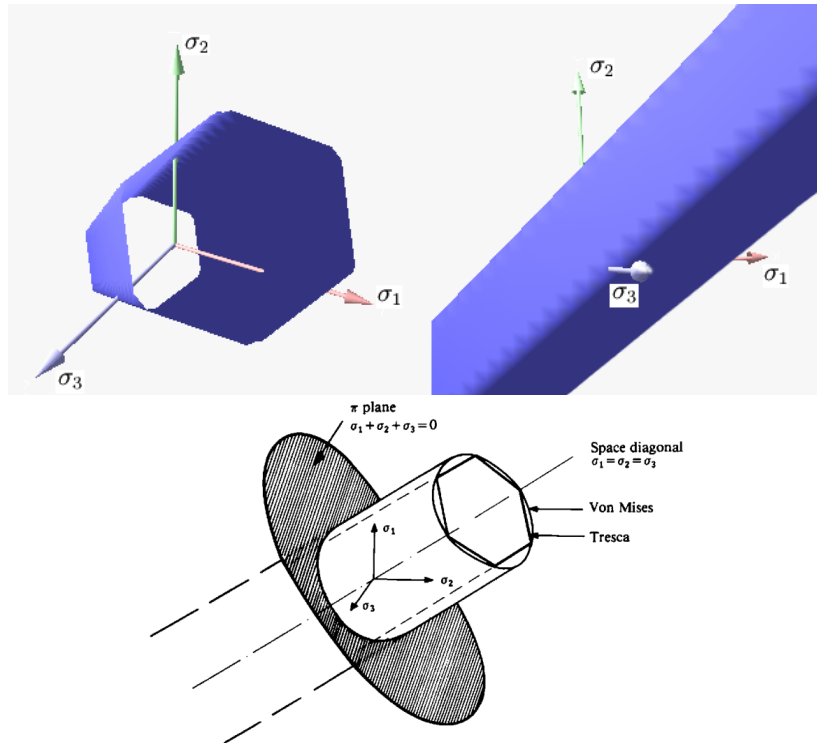
The Tresca or maximum shear stress yield criterion is taken to be the work of Henri Tresca. It is also referred as the Tresca-Guest (TG) criterion. The functional form of this yield criterion is

$$f(\sigma_1, \sigma_2, \sigma_3) = 0$$

In terms of the principal stresses the Tresca criterion is expressed as

$$\max(|\sigma_1 - \sigma_2|, |\sigma_2 - \sigma_3|, |\sigma_3 - \sigma_1|) = c$$

The following figure shows the Tresca-Guest yield surface in the three-dimensional space of principal stresses.



Left: Taken from https://en.wikipedia.org/wiki/Yield_surface. Right: Taken from Owen and Hinton [967].

It is a prism of six sides and having infinite length. This means that the material remains viscous when all three principal stresses are roughly equivalent (a hydrostatic pressure), no matter how much it is compressed or stretched. However, when one of principal stresses becomes smaller (or larger) than the others the material is subject to shearing. In such situations, if the shear stress reaches the yield limit then the material enters the plastic domain.


Remark. *The yield function presents sharp corners, making its numerical implementation more difficult (directional derivatives are needed)*

We have already established in Eq. (2.195):

$$\sigma_1 - \sigma_3 = 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \cos \theta$$

with $\sigma_1 > \sigma_2 > \sigma_3$, so that the failure criterion is given by

$$F^{\text{Tr}} = 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \cos \theta - c$$

 **Relevant Literature:** F. Tin-Loi and N.S. Ngo. “Performance of the p-version finite element method for limit analysis”. In: *International Journal of Mechanical Sciences* 45 (2003), pp. 1149–1166. DOI: 10.1016/j.ijmecsci.2003.08.004

We now look into the derivatives of the plastic potential $Q^{\text{Tr}}(\boldsymbol{\sigma})$. We have

$$Q^{\text{Tr}} = 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \cos \theta_L(\boldsymbol{\tau}) - c$$

$$\frac{\partial Q^{\text{Tr}}}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} = 0 \quad (2.280)$$

$$\begin{aligned} \frac{\partial Q^{\text{Tr}}}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} &= 2 \cos \theta_L - 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L \frac{\partial \theta_L}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \\ &= 2 \cos \theta_L - 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L \frac{\partial \theta_L}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \\ &= 2 \cos \theta_L - 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L \left(-\frac{1}{2} \tan 3\theta_L \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \right) 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \\ &= 2 \cos \theta_L + 2 \sin \theta_L \tan 3\theta_L \\ &= 2 \cos \theta_L (1 + \tan \theta_L \tan 3\theta_L) \end{aligned} \quad (2.281)$$

$$\frac{\partial Q^{\text{Tr}}}{\partial \theta_L(\boldsymbol{\tau})} = -2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L \quad (2.282)$$

so

$$C_1^{\text{Tr}} = 0 \quad (2.283)$$

$$\begin{aligned} C_2^{\text{Tr}} &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} - \frac{\tan 3\theta_L}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \right) \\ &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(2 \cos \theta_L (1 + \tan \theta_L \tan 3\theta_L) + \frac{\tan 3\theta_L}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L \right) \\ &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} (2 \cos \theta_L (1 + \tan \theta_L \tan 3\theta_L) + 2 \tan 3\theta_L \sin \theta_L) \\ &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} (2 \cos \theta_L (1 + \tan \theta_L \tan 3\theta_L) + 2 \cos \theta_L \tan 3\theta_L \tan \theta_L) \\ &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} (2 \cos \theta_L (1 + 2 \tan \theta_L \tan 3\theta_L)) \end{aligned} \quad (2.284)$$

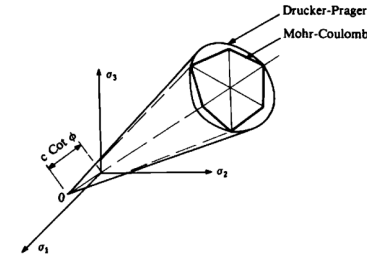
$$\begin{aligned} C_3^{\text{Tr}} &= -\frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \\ &= -\frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} (-2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \sin \theta_L) \\ &= \frac{\sqrt{3}}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\sin \theta_L}{\cos 3\theta_L} \end{aligned} \quad (2.285)$$

2.27.12 The Mohr-Coulomb failure criterion

mccriterion.tex

Mohr-Coulomb theory is a model describing the response of a material such as rubble piles or concrete to shear stress as well as normal stress. Most of the classical engineering materials somehow follow this rule in at least a portion of their shear failure envelope. In geology it is used to define shear strength of soils at different effective stresses [525].

In structural engineering it is used to determine failure load as well as the angle of fracture of a displacement fracture in concrete and similar materials. Coulomb's friction hypothesis is used to determine the combination of shear and normal stress that will cause a fracture of the material. Mohr's circle is used to determine the principal stresses that will produce this combination of shear and normal stress, and the angle of the plane in which this will occur. According to the principle of normality, the stress introduced at failure will be perpendicular to the line describing the fracture condition.



(a) Geometrical representation of the Mohr-Coulomb and Drucker-Prager yield surfaces in principal stress space.

Taken from Owen and Hinton [967].

The Mohr-Coulomb failure criterion represents the linear envelope that is obtained from a plot of the shear strength of a material versus the applied normal stress. This relation is expressed as (Owen & Hinton book [967, p219])

$$\tau_m = -\sigma_m \sin \phi + c \cos \phi \quad (2.286)$$

where τ_m is the magnitude of the shear stress, σ_m is the normal stress, c is the intercept of the failure envelope with the τ axis, and ϕ is the slope of the failure envelope. The minus sign in the above equation is for the case where compression is assumed to be negative³⁸. The quantity c is often called the cohesion and the angle ϕ is called the angle of internal friction.

We have

$$\tau_m = \frac{\sigma_1 - \sigma_3}{2} \quad \sigma_m = \frac{\sigma_1 + \sigma_3}{2}$$

with σ_1 is the maximum principal stress and σ_3 is the minimum principal stress, or

$$\frac{\sigma_1 - \sigma_3}{2} = -\frac{\sigma_1 + \sigma_3}{2} \sin \phi + c \cos \phi \quad (2.287)$$

Using Eqs. (2.195) and (2.196) for $(\sigma_1 - \sigma_3)/2$ and $(\sigma_1 + \sigma_3)/2$ we get³⁹:

$$\begin{aligned} \frac{\sigma_1 - \sigma_3}{2} &= -\frac{\sigma_1 + \sigma_3}{2} \sin \phi + c \cos \phi \\ \Rightarrow \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \cos \theta &= -\left(\frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) - \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \frac{1}{\sqrt{3}} \sin \theta \right) \sin \phi + c \cos \phi \\ \Rightarrow \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta - \frac{1}{\sqrt{3}} \sin \theta \sin \phi \right) - c \cos \phi &= 0 \end{aligned}$$

³⁸https://en.wikipedia.org/wiki/Mohr-Coulomb_theory

³⁹This is Eq. (7.16) of Owen and Hinton [967]

$$F^{\text{MC}} = \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta - \frac{1}{\sqrt{3}} \sin \theta \sin \phi \right) - c \cos \phi \quad (2.288)$$

This formula (without the cohesion) is used in Willett [1359]. Since $p = -\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})$, we also have:

$$F^{\text{MC}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta - \frac{1}{\sqrt{3}} \sin \theta \sin \phi \right) - (p \sin \phi + c \cos \phi) \quad (2.289)$$

Remark. The expression for F in the Mohr-Coulomb case in Zienkiewicz & Corneau (1974) [1423] contains errors which are later corrected in Zienkiewicz, Taylor, and Fox [1433, p102].

 **Relevant Literature:** this criterion is also used in computer graphics animation [1420]

when
 $\phi = 0$
we
should
recover
Tresca
but fac-
tor 2 is
wrong ?

We now look into the derivatives of the Drucker-Prager plastic potential $Q^{\text{MC}}(\boldsymbol{\sigma})$. We have

$$Q^{\text{MC}} = \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta_L(\boldsymbol{\tau}) - \frac{1}{\sqrt{3}} \sin \theta_L(\boldsymbol{\tau}) \sin \phi \right) - c \cos \phi$$

Then

$$\frac{\partial Q^{\text{MC}}}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} = \frac{1}{3} \sin \phi \quad (2.290)$$

$$\begin{aligned} \frac{\partial Q^{\text{MC}}}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} &= \cos \theta_L - \frac{1}{\sqrt{3}} \sin \theta_L \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\sin \theta_L - \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi \right) \frac{\partial \theta_L}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \\ &= \cos \theta_L - \frac{1}{\sqrt{3}} \sin \theta_L \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\sin \theta_L - \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi \right) \frac{\partial \theta_L}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \\ &= \cos \theta_L - \frac{1}{\sqrt{3}} \sin \theta_L \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\sin \theta_L - \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi \right) \left(-\frac{1}{2} \tan 3\theta_L \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \right) 2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \\ &= \cos \theta_L - \frac{1}{\sqrt{3}} \sin \theta_L \sin \phi + \left(\sin \theta_L + \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi \right) \tan 3\theta_L \\ &= \cos \theta_L \left[1 - \frac{1}{\sqrt{3}} \tan \theta_L \sin \phi + \left(\tan \theta_L + \frac{1}{\sqrt{3}} \sin \phi \right) \tan 3\theta_L \right] \\ &= \cos \theta_L \left[(1 + \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (\tan 3\theta_L - \tan \theta_L) \right] \end{aligned} \quad (2.291)$$

$$\begin{aligned} \frac{\partial Q^{\text{MC}}}{\partial \theta_L(\boldsymbol{\tau})} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(-\sin \theta_L - \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi \right) \\ &= -\frac{1}{\sqrt{3}} \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} (\sqrt{3} \sin \theta_L + \cos \theta_L \sin \phi) \end{aligned} \quad (2.292)$$

so

$$\begin{aligned}
C_1^{\text{MC}} &= \frac{1}{3} \sin \phi \\
C_2^{\text{MC}} &= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} - \frac{\tan 3\theta_L}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \right) \\
&= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\cos \theta_L \left[(1 + \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (\tan 3\theta_L - \tan \theta_L) \right] + \frac{\tan 3\theta_L}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{1}{\sqrt{3}} \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} (\sqrt{3} \sin \theta_L + \cos \theta_L \sin \phi) \right) \\
&= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\cos \theta_L \left[(1 + \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (\tan 3\theta_L - \tan \theta_L) \right] + \tan 3\theta_L (\sin \theta_L + \frac{1}{\sqrt{3}} \cos \theta_L \sin \phi) \right) \\
&= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\cos \theta_L \left[(1 + \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (\tan 3\theta_L - \tan \theta_L) + \tan 3\theta_L (\tan \theta_L + \frac{1}{\sqrt{3}} \sin \phi) \right] \right) \\
&= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \cos \theta_L \left[(1 + 2 \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (2 \tan 3\theta_L - \tan \theta_L) \right] \\
C_3^{\text{MC}} &= -\frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \\
&= -\frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \left[-\frac{1}{\sqrt{3}} \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} (\sqrt{3} \sin \theta_L + \cos \theta_L \sin \phi) \right] \\
&= \frac{\sqrt{3} \sin \theta_L + \sin \phi \cos \theta_L}{2\mathcal{I}_2(\boldsymbol{\tau}) \cos 3\theta_L}
\end{aligned}$$

2.27.13 The Drucker-Prager failure criterion

dpcriterion.tex

The von Mises yield criterion is not suitable for modelling the yielding of frictional material as it does not include the effect of mean stress as observed in experiments. To overcome this limitation, Drucker and Prager (1952) [347] proposed a revised function for frictional materials.

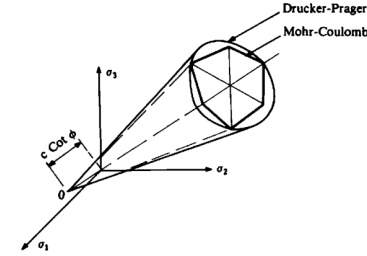
The Drucker-Prager yield criterion has the function form

$$F^{\text{DP}}(\boldsymbol{\sigma}) = F(\mathcal{I}_1(\boldsymbol{\sigma}), \mathcal{I}_2(\boldsymbol{\tau})) = 0 \quad (2.295)$$

This criterion is most often used for concrete where both normal and shear stresses can determine failure. The Drucker-Prager yield criterion may be expressed as

$$F^{\text{DP}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \alpha \mathcal{I}_1(\boldsymbol{\sigma}) + k = 0 \quad (2.296)$$

should it not be $-k$?



(a) Geometrical representation of the Mohr-Coulomb and Drucker-Prager yield surfaces in principal stress space.

Taken from Owen and Hinton [967].

Using the parameters σ_m , τ_m , $a = -\sqrt{3} \tan \theta$, $\mathcal{I}_1(\boldsymbol{\sigma})$ and $\mathcal{I}_2(\boldsymbol{\tau})$ of Section 2.25 we have

$$\begin{aligned} F^{\text{DP}} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \alpha \mathcal{I}_1(\boldsymbol{\sigma}) + k \\ &= \sqrt{\frac{\tau_m^2}{3}(a^2 + 3)} + \alpha(3\sigma_m - a\tau_m) + k \\ &= \tau_m \sqrt{(a^2/3 + 1)} + \alpha(3\sigma_m + \tau_m \sqrt{3} \tan \theta) + k \quad (\text{since } \tau_m > 0) \\ &= \tau_m \sqrt{\tan^2 \theta + 1} + \alpha(3\sigma_m + \tau_m \sqrt{3} \tan \theta) + k \\ &= \tau_m \sqrt{\frac{1}{\cos^2 \theta}} + \alpha(3\sigma_m + \tau_m \sqrt{3} \tan \theta) + k \\ &= \tau_m \frac{1}{\cos \theta} + \alpha(3\sigma_m + \tau_m \sqrt{3} \tan \theta) + k \quad (\text{since } \cos \theta > 0) \end{aligned}$$

$F = 0$ then leads to write

$$\begin{aligned} \tau_m + (3\alpha\sigma_m + k) \cos \theta + \tau_m \alpha \sqrt{3} \sin \theta &= 0 \\ \Rightarrow \tau_m(1 + \alpha \sqrt{3} \sin \theta) + (3\alpha\sigma_m + k) \cos \theta &= 0 \end{aligned}$$

and finally

$$\tau_m = -\frac{(3\alpha\sigma_m + k) \cos \theta}{1 + \alpha \sqrt{3} \sin \theta} = -\frac{3\alpha \cos \theta}{1 + \alpha \sqrt{3} \sin \theta} \sigma_m - \frac{k \cos \theta}{1 + \alpha \sqrt{3} \sin \theta}$$

Remark. This is the same equation as Eq. 19 of Wojciechowski [1367] but with $\theta \rightarrow -\theta$.

The Mohr-Coulomb yield criterion writes (see Eq. (2.286))

$$\tau_m = -\sigma_m \sin \phi + c \cos \phi$$

so that equating both expressions of τ_m for the Drucker-Prager and Mohr-Coulomb criteria leads to:

$$-\frac{3\alpha \cos \theta}{1 + \alpha\sqrt{3} \sin \theta} = -\sin \phi \quad (2.297)$$

$$-\frac{k \cos \theta}{1 + \alpha\sqrt{3} \sin \theta} = c \cos \phi \quad (2.298)$$

Eq. (2.297) yields

$$\begin{aligned} 3\alpha \cos \theta &= \sin \phi (1 + \alpha\sqrt{3} \sin \theta) \\ \Rightarrow 3\alpha \cos \theta - \alpha\sqrt{3} \sin \theta \sin \phi &= \sin \phi \end{aligned}$$

and finally

$$\boxed{\alpha(\phi) = \frac{\sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi}}$$

Inserting this into Eq. (2.298):

$$\begin{aligned} -k \cos \theta &= c \cos \phi (1 + \alpha\sqrt{3} \sin \theta) \\ &= c \cos \phi \left(1 + \frac{\sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \sqrt{3} \sin \theta \right) \\ &= c \cos \phi \left(1 + \frac{\sqrt{3} \sin \phi \sin \theta}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \right) \\ &= c \cos \phi \left(\frac{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} + \frac{\sqrt{3} \sin \phi \sin \theta}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \right) \\ &= c \cos \phi \left(\frac{3 \cos \theta}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \right) \end{aligned}$$

so that

$$\boxed{k(c, \phi) = -\frac{3 c \cos \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi}}$$

The Drucker-Prager yield criterion which for a given θ is equal to the Mohr-Coulomb yield is then:

$$\begin{aligned} F^{\text{DP}} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \alpha(\phi) \mathcal{I}_1(\boldsymbol{\sigma}) + k(c, \phi) \\ &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{\sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \mathcal{I}_1(\boldsymbol{\sigma}) - \frac{3 c \cos \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \\ &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{3 \sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{3 c \cos \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \right] \quad (2.299) \end{aligned}$$

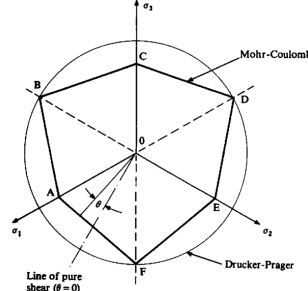
$$\begin{aligned} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[\frac{3 p \sin \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} + \frac{3 c \cos \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \right] \\ &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \frac{3 p \sin \phi + 3 c \cos \phi}{3 \cos \theta - \sqrt{3} \sin \theta \sin \phi} \\ &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \frac{p \sin \phi + c \cos \phi}{\cos \theta - \frac{1}{\sqrt{3}} \sin \theta \sin \phi} \quad (2.300) \end{aligned}$$

which, when multiplied by $\cos \theta - \frac{1}{\sqrt{3}} \sin \theta \sin \phi$, gives the Mohr-Coulomb criterion of Eq. (2.288).

For $\theta = \pi/6$, the DP yield surface **circumscribes** the MC yield surface and Eq. (2.299) writes:

$$\begin{aligned}
 F^{\text{DP}} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{3 \sin \phi}{3\sqrt{3}/2 - \sqrt{3}/2 \sin \phi} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{3 c \cos \phi}{3\sqrt{3}/2 - \sqrt{3}/2 \sin \phi} \right] \\
 &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{6 \sin \phi}{\sqrt{3}(3 - \sin \phi)} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{6 c \cos \phi}{\sqrt{3}(3 - \sin \phi)} \right] \\
 &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \frac{6p \sin \phi + 6 c \cos \phi}{\sqrt{3}(3 - \sin \phi)}
 \end{aligned} \tag{2.301}$$

i.e.



(b) Two-dimensional, π plane, representation of the Mohr-Coulomb and Drucker-Prager yield criteria.

Taken from Owen and Hinton [967].

$$F^{\text{DP}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{6 \sin \phi}{\sqrt{3}(3 - \sin \phi)} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} - \frac{6 c \cos \phi}{\sqrt{3}(3 - \sin \phi)} \tag{2.302}$$

which is the formula used in Glerum *et al.* (2018) [467]. This is also Eq. (14a) in Zienkiewicz & Corneau (1974) [1423], Eq. (7.18) in Owen and Hinton [967], and Eq. (13.10a) in Zienkiewicz (1975) [1422] provided it is divided altogether by $\sqrt{3}$.

For $\theta = -\pi/6$, the DP yield surface **middle circumscribes** the MC yield surface and Eq. (2.299) writes:

$$\begin{aligned}
 F^{\text{DP}} &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{3 \sin \phi}{3\sqrt{3}/2 + \sqrt{3}/2 \sin \phi} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{3 c \cos \phi}{3\sqrt{3}/2 + \sqrt{3}/2 \sin \phi} \right] \\
 &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{6 \sin \phi}{\sqrt{3}(3 + \sin \phi)} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{6 c \cos \phi}{\sqrt{3}(3 + \sin \phi)} \right] \\
 &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \frac{6p \sin \phi + 6c \cos \phi}{\sqrt{3}(3 + \sin \phi)}
 \end{aligned} \tag{2.303}$$

This is Eq. (7.19) of Owen and Hinton [967].

Another DP formulation which **inscribes** the MC yield surface is found on the wikipedia page of the Drucker-Prager yield criterion ⁴⁰ (but I have no idea how it is arrived at):

$$F^{\text{DP}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - \left[-\frac{3 \sin \phi}{\sqrt{9 + 3 \sin^2 \phi}} \frac{\mathcal{I}_1(\boldsymbol{\sigma})}{3} + \frac{3 c \cos \phi}{\sqrt{9 + 3 \sin^2 \phi}} \right] \tag{2.304}$$

The yield surfaces of these three Drucker-Prager formulations are plotted against the Mohr-Coulomb yield surface in Section 2.27.16.

⁴⁰https://en.wikipedia.org/wiki/Drucker-Prager_yield_criterion

Remark. Leroy & Ortiz [772] use the Drucker-Prager plasticity model also and match it to the Mohr-Coulomb model in the triaxial test and formulate it as follows (Their definition of the second invariant of stress contains a $3/2$ term):

$$\begin{aligned}
F &= \tau_e \sqrt{3} + \frac{6 \sin \phi}{3 - \sin \phi} \left(-p - \frac{c}{\tan \phi} \right) \\
&= \tau_e \sqrt{3} - \left(\frac{6 \sin \phi}{3 - \sin \phi} p + c \frac{6 \cos \phi}{3 - \sin \phi} \right) \\
&= \sqrt{3} \left[\tau_e - \left(\frac{6 \sin \phi}{\sqrt{3}(3 - \sin \phi)} p + c \frac{6 \cos \phi}{\sqrt{3}(3 - \sin \phi)} \right) \right]
\end{aligned} \tag{2.305}$$


Except for the $\sqrt{3}$ this is identical to Eq. (2.301).

Remark. Bui et al. (2008) [160] use yet again another formulation:

$$\begin{aligned}
F &= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{\tan \phi}{\sqrt{9 + 12 \tan^2 \phi}} \mathcal{I}_1(\boldsymbol{\sigma}) - \frac{3c}{\sqrt{9 + 12 \tan^2 \phi}} \\
&= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{\sin \phi}{\sqrt{9 \cos^2 \phi + 12 \sin^2 \phi}} \mathcal{I}_1(\boldsymbol{\sigma}) - \frac{3c \cos \phi}{\sqrt{9 \cos^2 \phi + 12 \sin^2 \phi}} \\
&= \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{\sin \phi}{\sqrt{9 + 3 \sin^2 \phi}} \mathcal{I}_1(\boldsymbol{\sigma}) - \frac{3c \cos \phi}{\sqrt{9 + 3 \sin^2 \phi}}
\end{aligned}$$

which is identical to (2.304).

Remark. Cacace & Jacquety (2017) [197] replace $\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}$ by $\sqrt{\mathcal{I}_2(\boldsymbol{\tau}) + \epsilon_0^2}$ where ϵ_0 is a small non-hardening parameters here introduced to relax the singularity at the cone's tip of the Drucker-Prager yield envelope.

 **Relevant Literature:** Gilda Currenti and Charles A Williams. “Numerical modeling of deformation and stress fields around a magma chamber: Constraints on failure conditions and rheology”. In: *Physics of the Earth and Planetary Interiors* 226 (2014), pp. 14–27. DOI: 10.1016/j.pepi.2013.11.003

Dissecting the original paper by Drucker and Prager (1952) The authors state that a yield function which is a proper generalisation of the M-C hypothesis is:

$$F = \alpha \mathcal{I}_1(\boldsymbol{\sigma}) + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - k$$

where α and k are positive constants at each point of the material.

According to the concept of plastic potential, the stress-strain relation corresponding to this yield function is

$$\dot{\epsilon}_{ij}^p = \lambda \frac{\partial F}{\partial \sigma_{ij}}$$

where $\dot{\epsilon}_{ij}^p$ is the plastic strain rate and λ is a positive factor of proportionality which may assume different values in space. Using the above expression for F :

$$\dot{\epsilon}_{ij}^p = \lambda \left(\alpha \delta_{ij} + \frac{\tau_{ij}}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \right) \tag{2.306}$$

A very important feature of this equation is that the plastic rate of cubical dilation is

$$\text{tr}[\dot{\boldsymbol{\epsilon}}^p] = \dot{\epsilon}_{ii}^p = 3\alpha\lambda \neq 0 \tag{2.307}$$

This equation shows that plastic deformation must be accompanied by an increase in volume if $\alpha \neq 0$. This property is known as dilatancy.

Plane strain We need to establish three expressions. First, from Eq. (2.306) we can write

$$\dot{\varepsilon}_{zz}^p = \lambda \left(\alpha + \frac{\tau_{zz}}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \right)$$

but since $\dot{\varepsilon}_{zz} = 0$ in plane strain then we find

$$\tau_{zz} = -2\alpha\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad (2.308)$$

which is Eq. (6) of the paper.

Second, we start from the definition of the first invariant and use the equation above:

$$\begin{aligned} \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\ \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} + \tau_{zz} + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \\ \mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} - 2\alpha\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \\ \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} - 2\alpha\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \\ \mathcal{I}_1(\boldsymbol{\sigma}) &= \frac{3}{2}(\sigma_{xx} + \sigma_{yy}) - 3\alpha\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \end{aligned} \quad (2.309)$$

which is Eq. (7) of the paper.

Finally, we start from (and we use the fact that $\sigma_{xx} - \sigma_{yy} = \tau_{xx} - \tau_{yy}$)

$$\begin{aligned} \left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 &= \frac{1}{4}(\sigma_{xx} - \sigma_{yy})^2 + \tau_{xy}^2 \\ &= \frac{1}{4}(\sigma_{xx} - \sigma_{yy})^2 + \underbrace{\tau_{xy}^2 + \frac{1}{2}(\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2)}_{\mathcal{I}_2(\boldsymbol{\tau})} - \frac{1}{2}(\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2) \\ &= \frac{1}{4}(\tau_{xx} - \tau_{yy})^2 + \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{2}\tau_{xx}^2 - \frac{1}{2}\tau_{yy}^2 - \frac{1}{2}\tau_{zz}^2 \\ &= \mathcal{I}_2(\boldsymbol{\tau}) + \frac{1}{4}\tau_{xx}^2 - \frac{1}{2}\tau_{xx}\tau_{yy} + \frac{1}{4}\tau_{yy}^2 - \frac{1}{2}\tau_{xx}^2 - \frac{1}{2}\tau_{yy}^2 - \frac{1}{2}4\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{4}\tau_{xx}^2 - \frac{1}{2}\tau_{xx}\tau_{yy} - \frac{1}{4}\tau_{yy}^2 - 2\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{4}(\tau_{xx}^2 + 2\tau_{xx}\tau_{yy} + \tau_{yy}^2) - 2\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{4}\underbrace{(\tau_{xx} + \tau_{yy})^2}_{-\tau_{zz}} - 2\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{4}\tau_{zz}^2 - 2\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{4}4\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) - 2\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau}) - 3\alpha^2\mathcal{I}_2(\boldsymbol{\tau}) \\ &= \mathcal{I}_2(\boldsymbol{\tau})(1 - 3\alpha^2) \end{aligned} \quad (2.310)$$

so that

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{1 - 3\alpha^2} \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right] \quad (2.311)$$

which is Eq. (8) of the paper.

In the paper the authors propose the yield function

$$F = \alpha \mathcal{I}_1(\boldsymbol{\sigma}) + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - k$$

We first replace the first (plane strain) invariant (see Eq. (2.309)):

$$\begin{aligned} F &= \alpha \left[\frac{3}{2}(\sigma_{xx} + \sigma_{yy}) - 3\alpha \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \right] + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - k \\ &= \alpha \frac{3}{2}(\sigma_{xx} + \sigma_{yy}) - 3\alpha^2 \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - k \\ &= \alpha \frac{3}{2}(\sigma_{xx} + \sigma_{yy}) + (1 - 3\alpha^2) \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - k \end{aligned}$$

and we now introduce the second invariant of Eq. (2.311):

$$\begin{aligned} F &= \alpha \frac{3}{2}(\sigma_{xx} + \sigma_{yy}) + (1 - 3\alpha^2) \frac{1}{(1 - 3\alpha^2)^{1/2}} \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - k \\ &= \alpha \frac{3}{2}(\sigma_{xx} + \sigma_{yy}) + (1 - 3\alpha^2)^{1/2} \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - k \\ &= \frac{3\alpha}{(1 - 3\alpha^2)^{1/2}} \frac{1}{2}(\sigma_{xx} + \sigma_{yy}) + \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - \frac{k}{(1 - 3\alpha^2)^{1/2}} \\ &= \underbrace{\frac{3\alpha}{(1 - 3\alpha^2)^{1/2}} \frac{1}{2}(\sigma_{xx} + \sigma_{yy})}_{\sin \phi} + \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - \underbrace{\frac{k}{(1 - 12\alpha^2)^{1/2}}}_{c} \underbrace{\frac{(1 - 12\alpha^2)^{1/2}}{(1 - 3\alpha^2)^{1/2}}}_{\cos \phi} \quad (2.312) \end{aligned}$$

Note that if we define a triangle with sides $(1 - 12\alpha^2)^{1/2}$ and 3α with hypotenuse $(1 - 3\alpha^2)^{1/2}$ then the angle ϕ makes sense and we recover $\cos^2 \phi + \sin^2 \phi = 1$.

In the end:

$$F = \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - \left(-\frac{1}{2}(\sigma_{xx} + \sigma_{yy}) \sin \phi + c \cos \phi \right)$$

which is the Mohr-Coulomb yield criterion of Eq. (1) in the paper.

Note that when $\alpha = 0$ (yield criterion independent of the mean stress - incompressible flow see Eq. (2.307)) then $c = k$, $\cos \phi = 1$ and $\sin \phi = 0$ and we find the Tresca yield criterion

$$F^{TR} = \left[\left(\frac{\sigma_{xx} - \sigma_{yy}}{2} \right)^2 + \tau_{xy}^2 \right]^{1/2} - k$$

Also, setting $\alpha = 0$ in Eq. (2.311) yields a criterion that writes

$$F^{vM} = \mathcal{I}_2(\boldsymbol{\tau}) - k$$

which is the von Mises criterion!

We now look into the derivatives of the Drucker-Prager plastic potential $Q^{\text{DP}}(\boldsymbol{\sigma})$. We have

$$Q^{\text{DP}}(\boldsymbol{\sigma}) = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + \alpha \mathcal{I}_1(\boldsymbol{\sigma}) + k$$

Then

$$\frac{\partial Q^{\text{DP}}}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} = \alpha \quad (2.313)$$

$$\frac{\partial Q^{\text{DP}}}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} = 1 \quad (2.314)$$

$$\frac{\partial Q^{\text{DP}}}{\partial \theta_{\text{L}}(\boldsymbol{\tau})} = 0 \quad (2.315)$$

The parameters α and k can be expressed as a function of the angle of friction and cohesion so as to match the Mohr-Coulomb criterion in some sense (see above). Then

$$C_1^{\text{DP}} = \alpha \quad (2.316)$$

$$C_2^{\text{DP}} = \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \quad (2.317)$$

$$C_3^{\text{DP}} = 0 \quad (2.318)$$

ToDo: check Alejano and Bobet [6] (2012) and compare with my notes above.

2.27.14 The Griffith-Murrell failure criterion

The Griffith-Murrell yield criterion [140, 143, 54] is not often used. Extending the work of Griffith (1921) to three dimensional stress distributions, Murrell (1963) suggested the following criterion for rock failure expressed in terms of the principal stresses:

$$(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_3 - \sigma_1)^2 + 24T_0(\sigma_1 + \sigma_2 + \sigma_3) = 0$$

where T_0 is a material property called the tensile strength. In principal stress space, this criterion is represented by a paraboloid of revolution around the pressure (or hydrostatic) axis.

Using the definition of $\mathcal{I}_2(\boldsymbol{\tau})$ and $\mathcal{I}_1(\boldsymbol{\sigma})$, it also writes:

$$\mathcal{I}_2(\boldsymbol{\tau}) - 12T_0p = 0$$

which is the formulation used in Hansen *et al.* (2000) [532], although the authors use the lithostatic pressure instead of the full pressure. They also use a tensile strength parameter T_0^e and a compressive strength parameter T_0^c , both around a few tens of MPas.

2.27.15 The Cam-clay failure criterion

camclay.tex

The Original Cam-Clay model is based on the assumption that the soil is isotropic, elasto-plastic, deforms as a continuum, and it is not affected by creep.

 Relevant Literature: [991]

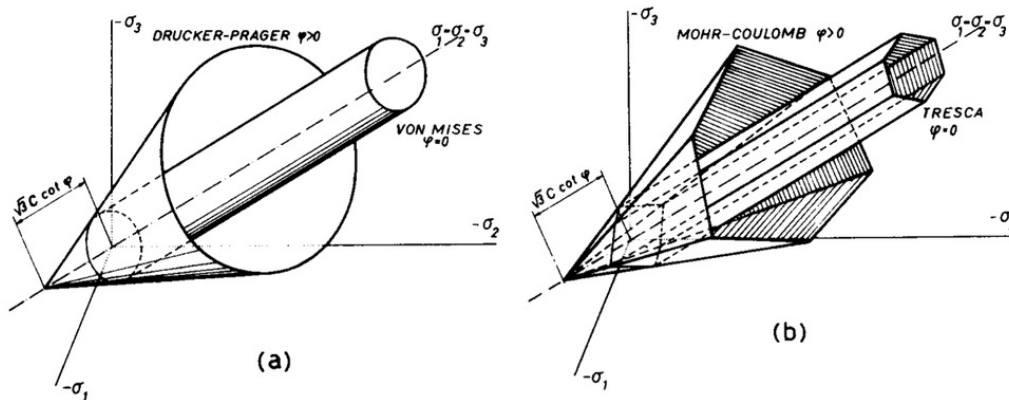
ask Chris Spiers, Pijenburg *et al.*, JGR 2019

2.27.16 The failure envelope, or yield surface

 Relevant Literature: Schöpfer *et al.* (2013) [1136].

A yield surface is a five-dimensional surface in the six-dimensional space of stresses. The state of stress of inside the yield surface is elastic. When the stress state lies on the surface the material is said to have reached its yield point and the material is said to have become plastic. Further deformation of the material causes the stress state to remain on the yield surface, even though the surface itself may change shape and size as the plastic deformation evolves, this is because stress states that lie outside the yield surface are non-permissible.

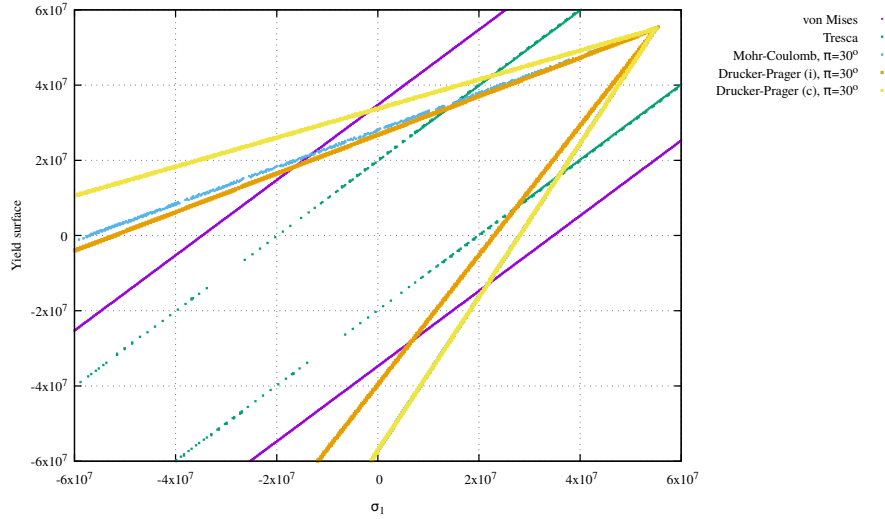
The yield surface is usually expressed in terms of (and visualized in) a three-dimensional principal stress space $(\sigma_1, \sigma_2, \sigma_3)$, a two- or three-dimensional space spanned by stress invariants or a version of the three-dimensional Haigh-Westergaard space.



Yield surfaces in stress space [1423]. Note that the axes are $-\sigma_1, -\sigma_2, -\sigma_3$

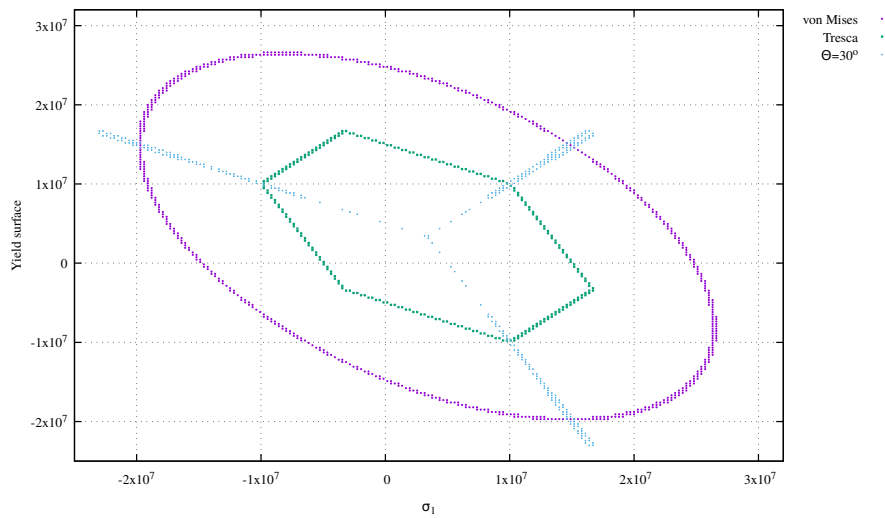
Having obtained the equations for the yield functions in the previous sections, we can easily test them as follows: in the $(\sigma_1, \sigma_2, \sigma_3)$ space we can look for stress states that fulfil the yield equations. I set $c = 20\text{MPa}$ and $\phi = 20^\circ$ and restrain the search to the space $[-100\text{MPa}; 100\text{MPa}]^3$. The python code and the gnuplot script used to generate the plots hereafter are in `images/rheology/surfaces`. The implemented algorithm is somewhat naive and quite inefficient: discretise the space in N^3 points and for each point check whether any of the von Mises, Tresca, Mohr-Coulomb and (the three variants of) Drucker-Prager criteria is satisfied and when the point is in the space $\sigma_1 + \sigma_2 + \sigma_3 = 10\text{MPa}$ (perpendicular to the $x = y = z$ line) write it to the corresponding file.

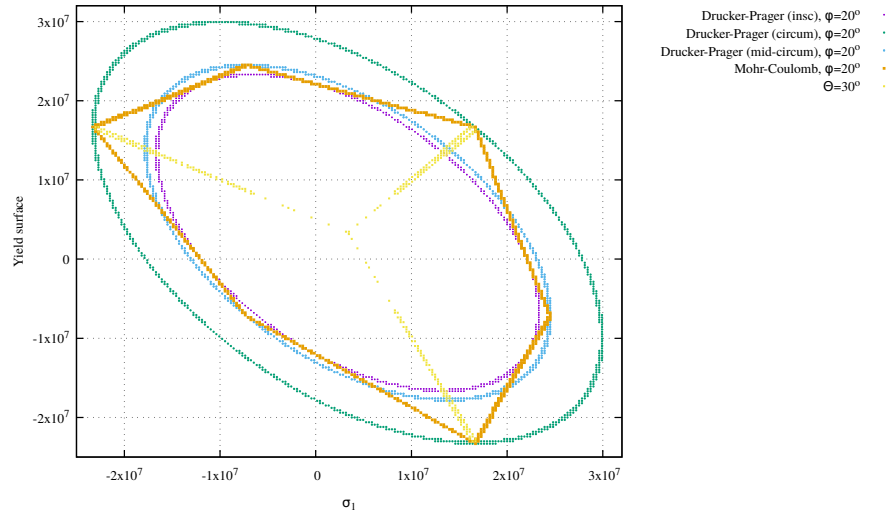
The recovered surfaces are similar to those of the figure above but their plot in a 3D space is difficult. I have therefore isolated two sub-plots. The first one is for $\sigma_1 = \sigma_2$:



We see that the von Mises and Tresca envelopes are parallel to the line $\sigma_1 = \sigma_2 = \sigma_3$ (which is expected since they do not depend on pressure).

The second plot is in the plane $\sigma_1 + \sigma_2 + \sigma_3 = 0$ which is perpendicular to the middle line $\sigma_1 = \sigma_2 = \sigma_3 = 0$. To facilitate plotting the envelopes are plotted as a function of σ_1 only (so that even though they are circles in the chosen plane they appear here as ellipses):





We see that we indeed recover that the three Drucker-Prager formulations inscribe (purple), middle-circumscribe (blue) and circumscribe (green) the Mohr-Coulomb one.

2.27.17 Peierls creep

peierls.tex

Looking at the literature, there seem to be many formulations for the Peierls creep deformation mechanism, but it appears that a standard formulation for the Peierls creep writes:

$$\dot{\varepsilon} = A\sigma^n \exp \left[-\frac{Q + pV}{RT} \left(1 - \left(\frac{\sigma}{\sigma_P} \right)^k \right)^q \right]$$

and it seems common to take $k = 1$, and $n = 2$ [455, 675]

$$\dot{\varepsilon} = A\sigma^2 \exp \left[-\frac{Q + pV}{RT} \left(1 - \frac{\sigma}{\sigma_P} \right)^q \right]$$

Elbeshhausen & Melosh (2018) [366] use

$$\dot{\varepsilon} = A \exp \left[-\frac{Q}{RT} \left(1 - \frac{\sigma}{\sigma_P} \right)^q \right]$$

In Chenin *et al.* (2019) [229] the authors state that their Peierls creep implementation relies on parameters from Evans and Goetze (1979) [383] using the approach of Kameyama *et al.* (1999) [666]:

$$\eta^{pe} = \frac{2}{3} \frac{(1-s)/s}{(1+s)/2s} A (\varepsilon_e^{ds})^{\frac{1}{n}-1}$$


with A for this formulation:

$$A = \left[A_p \exp \left(-\frac{Q(1-\gamma)^2}{RT} \right) \right]^{-1/s} \gamma \sigma_p$$

where s is an effective stress exponent that depends on the temperature:

$$s = 2\gamma \frac{Q}{RT} (1 - \gamma)$$

where γ is a fitting parameter.

 **Relevant Literature** Babeyko, Sobolev, Vietor, Oncken, and Trumbull [36], Burov [187], Faul, Gerald, Farlai, Ahlefeldt, and Jackson [388], Garel, Goes, Davies, Davies, Kramer, and Wilson [435], Gerya [455], Goetze and Evans [470], Katayama and Karato [675], Kawazoe, Karato, Otsuka, Jing, and Mookherjee [685], Karato, Riedel, and Yuen [671], Mei, Suzuki, Kohlstedt, Dixon, and Durham [860], Zhong and Watts [1413], Chenin, Schmalholz, Manatschal, and Karner [228], Shi, Wei, Li, Liu, and Liu [1160], Review article from 1966: Guyot & Dorn [518]

2.27.18 Stress limiting rheology

Taken from van Hunen *et al.* (2002) [616]:

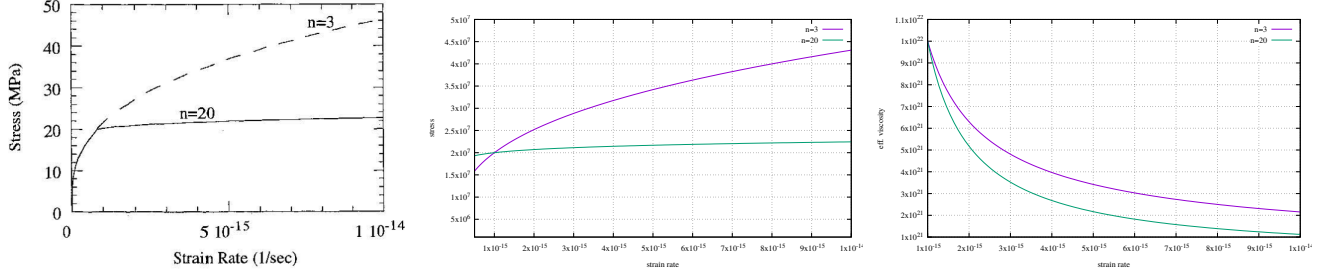
$$\eta_y = \tau_y \dot{\varepsilon}_y^{-1/n_y} \dot{\varepsilon}_e^{(1/n_y)-1}$$

where the yield stress τ_y , the yield strain rate $\dot{\varepsilon}_y$ and the yield exponent n_y are prescribed parameters. In this article, $n_y = 10$, $\dot{\varepsilon}_y = 10^{-15} \text{s}^{-1}$. When $n_y = 1$ the viscosity is constant and given by $\eta_{eff} = \tau_y / \dot{\varepsilon}_y$.

This rheology has also been coined pseudo-plastic in Zhong *et al.* (1998) [1411]. Their equation is simply

$$\eta_{eff} = A^{1/n} \dot{\varepsilon}_e^{-1+1/n}$$

where A is the preexponent which depends on temperature, pressure, and composition.



Left figure is taken from [1411]. Authors report $A = 7.9 \cdot 10^{-8} \text{Pa}^3 \text{s}$ for the $n = 3$ case, which makes no sense. See gnuplot script for actual values of A .

2.27.19 Arrhenius law

A purely temperature-dependent dimensional Arrhenius law that emulates the temperature dependence of viscosity in silicate rock is often employed for mantle rocks [5, 1406, 557, 126, 923, 1196, 98, 504]:

$$\eta(T) = \eta_0 \exp\left(\frac{Q}{R}\left(\frac{1}{T} - \frac{1}{T_0}\right)\right) \quad \text{or} \quad \eta(T) = \eta_0 \exp\left(\frac{Q}{RT}\right) \quad (2.319)$$

where η_0 is a reference viscosity and T_0 its corresponding reference temperature.

It can also account for pressure effects as in [812] where the diffusion creep viscosity (under the assumption of homogeneous grain size) is temperature- and pressure-dependent:

$$\eta(T) = \eta_0 \exp\left(\frac{1}{R}\left(\frac{Q - pV}{T} - \frac{Q}{T_0}\right)\right)$$

(I find the minus sign rather suspicious)

2.27.20 Simple parametrisation of the mantle

Many CITCOMS-based publications [166, 165] have used the following (dimensionless) viscosity for the mantle:

$$\eta(T, z) = \eta_r(r) \exp(A(0.5 - T))$$

where η_r is a depth-dependent viscosity profile (usually defined as discontinuous linear profiles for various shells)

The non-dimensional activation coefficient is chosen to be $A = 9.2103$ in [165] which leads to a temperature-induced viscosity contrast of 10^4 (for $T \in [0, 1]$).

This is also called the Frank-Kamenetskii flow rule, as used in [1196, 756]:

$$\eta' = \eta_0 \exp(-\theta T)$$

where the parameters η_0 , θ account for the local chemical composition of the rock. Note that the Frank-Kamenetskii approximation takes many forms in the literature [943].

Another temperature-dependent common expression is as follows [397]:

$$\eta(T) = \eta_\infty \exp\left(\frac{Q}{R}\left(\frac{1}{T} - \frac{1}{T_\infty}\right)\right)$$

Also, following [397]: For studying transient convection in a non-Newtonian rheological fluid, it is expedient from a computational point of view to employ a law which behaves linearly for low stresses initially and becomes gradually non-Newtonian only after a certain threshold stress level has been surpassed [243, 245]:


$$\eta(T, p, \tau_2) = \eta(T, p) \frac{1}{A_2 + A_3 \tau_2^2}$$

where A_2 is a parameter describing the linear creep at low stress levels and A_3 governs the transition stress between Newtonian and non-Newtonian rheologies.

Coltice and Sheppard (2018) [274] use a depth- and temperature-dependent viscosity formulation:

$$\eta(z, T) = \eta_0(z) \exp \frac{Q}{RT}$$

Note that this expression is supplemented with a pseudo-plastic formulation [1082].

 Relevant Literature: [705]

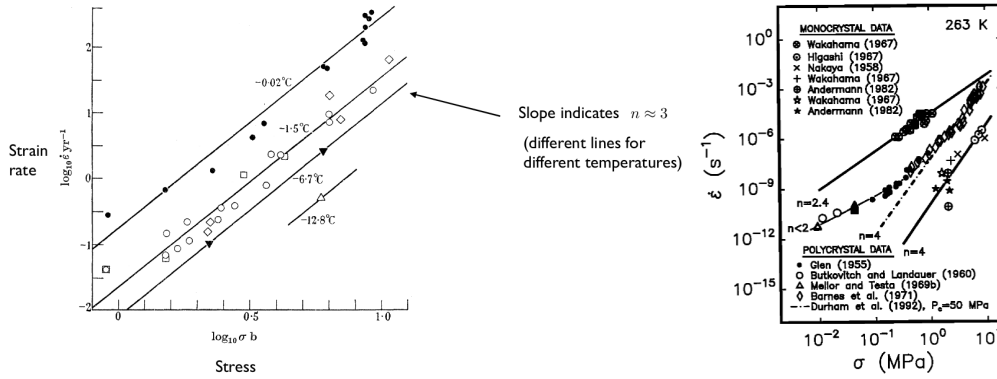
2.27.21 Glen's law for ice

Ice and rocks share similarities in terms of (viscous) rheology. Glen's law is the most commonly used flow law for ice in glaciers and ice sheets [466] and it is actually a power-law type rheology:

$$\dot{\epsilon} = A \tau^n$$

with $n \sim 3$ and $A \sim 2.4 \cdot 10^{-24} \text{Pa}^{-3} \cdot \text{s}^{-1}$ at 0°C . The effective viscosity is then given by

$$\eta = \frac{1}{2A\tau_e^{n-1}}$$



Left: Taken from Glen [466]; Right: taken from [471].

Most of these studies suggest values of the power-law exponent $n \sim 2 - 4$, and there seems to be a general indication that the exponent is lower at lower stresses.

The A coefficient above has been found to depend on temperature and is reasonably described with an Arrhenius law:

$$A(T) = A_0 \exp \left(-\frac{Q}{RT} \right)$$


A standard formulation is the Paterson-Budd law with a fixed Glen exponent $n = 3$ and a split Arrhenius term [982]:

$$A = 3.615 \cdot 10^{-13} \text{Pa}^{-3} \cdot \text{s}^{-1}, \quad Q = 60 \text{ kJ/mol}, \quad \text{if } T < 263 \text{ K}$$

$$A = 1.733 \cdot 10^3 \text{Pa}^{-3} \cdot \text{s}^{-1}, \quad Q = 139 \text{ kJ/mol}, \quad \text{if } T > 263 \text{ K}$$

Be careful that in these two equations the temperature T is the pressure-adjusted temperature [982]. Note that A is also affected by the water content and the presence of impurities.

Finally, Glen's law is the standard rheology used for ice-sheet modelling but it does not account for the complex evolution of fabric and resulting anisotropy. Indeed the grain size evolution (growth & reduction) plays a large role in the rheology Behn, Goldsby, and Hirth [68] (2021). See STONE 59.

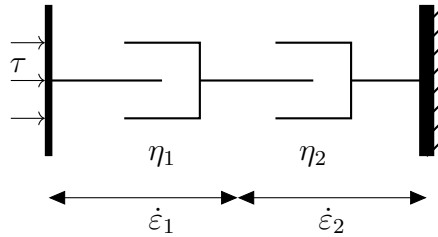
 Relevant Literature Alley [8] (1992), Greve [492] (1997), Greve and Blatter [493] (2009), Isaac, Stadler, and Ghattas [623] (2015), Krabbendam [728] (2016), Jiménez, Duddu, and Bassis [648] (2017), Helanow and Ahlkrone [561] (2018), Very mathematics heavy papers: [658, 222].

2.27.22 Strain rate partitioning across deformation mechanisms

When multiple viscous deformation mechanisms are present, one needs more dashpots, and more complicated element diagrams than the ones above occur (also when adding plastic deformation). Two important rules are to be remembered: 1) for parallel components, stresses are additive, strain rates are equal in each; 2) for components in series, stresses are equal in each and strain rates are additive.

Let us then look at various assemblies of dashpots and plastic elements:

- two viscous dampers in series:



each is subjected to the same stress τ but deforms with its own strain rate $\dot{\epsilon}_1$ and $\dot{\epsilon}_2$ and we have

$$\dot{\epsilon}_T = \dot{\epsilon}_1 + \dot{\epsilon}_2 = \frac{\tau}{2\eta_1} + \frac{\tau}{2\eta_2} \quad (2.320)$$

The effective viscosity of this combination is denoted η_{eff} and is such that $\eta_{eff} = \tau/2\dot{\epsilon}_T$, which means that

$$\frac{\tau}{2\eta_{eff}} = \frac{\tau}{2\eta_1} + \frac{\tau}{2\eta_2}$$

or,

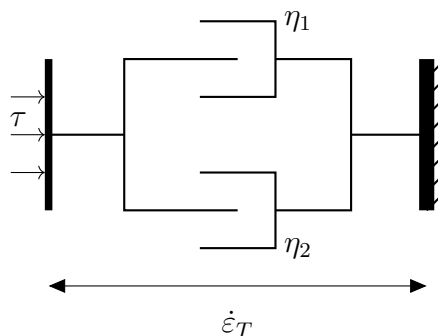
$$\eta_{eff} = \left(\frac{1}{\eta_1} + \frac{1}{\eta_2} \right)^{-1}$$

i.e. it follows that the effective viscosity of two or more viscous dampers in series is the harmonic average of the individual viscosities of the dampers.

In general, for n dampers in series:

$$\eta_{eff} = \left(\sum_{i=1}^n \frac{1}{\eta_i} \right)^{-1}$$

- two viscous dampers in parallel:



each is deformed with the same strain rate $\dot{\epsilon}_T$ and their stresses add up:

$$\tau = \tau_1 + \tau_2 = 2\eta_1\dot{\epsilon}_T + 2\eta_2\dot{\epsilon}_T$$

and since we define the effective viscosity as $\tau = 2\eta_{eff}\dot{\epsilon}_T$ then it follows:

$$2\eta_{eff}\dot{\epsilon}_T = 2\eta_1\dot{\epsilon}_T + 2\eta_2\dot{\epsilon}_T$$

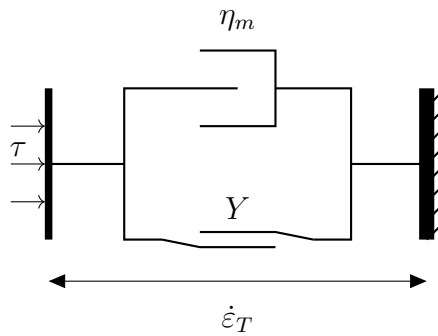
or,

$$\eta_{eff} = \eta_1 + \eta_2$$

i.e., the effective viscosity of two or more viscous dampers is the sum of their viscosities (*but not their arithmetic mean!*).

- one viscous damper and a plastic element in parallel:

(tikz_vp.tex)



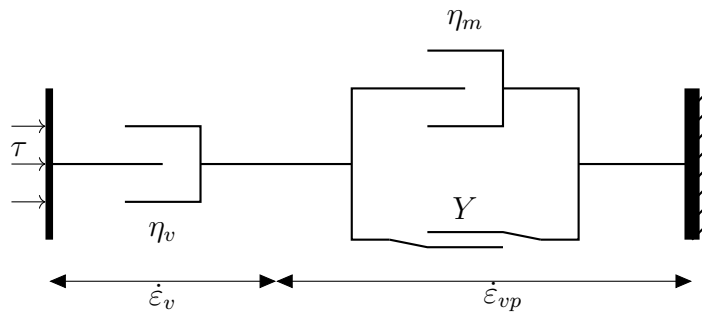
The effective 'plastic' viscosity of the plastic element is $\eta_p = \frac{Y}{2\dot{\epsilon}_T}$ so the effective viscosity of this setup is then

$$\eta_{eff} = \frac{Y}{2\dot{\epsilon}_T} + \eta_m$$

which is the viscosity of a Bingham fluid (see Section 2.27.4).

- two viscous dampers and a plastic element arranged as follows:

(tikz_vvp.tex)



This rheology would be called visco-viscoplastic. The algorithm goes then as follows:

1. Assume we know η_v and $\dot{\epsilon}_T$ (from previous iteration), as well as the plasticity parameters Y and η_m .
2. if $2\eta_v\dot{\epsilon}_T < Y$ the stress is below the yield stress value and plasticity is not active. Use η_v in the material model and $\dot{\epsilon}_v = \dot{\epsilon}_T$.

3. if $2\eta_v\dot{\epsilon}_T > Y$ the stress is above the yield value, which is not allowed. In this case the plastic element is 'switched on'. In that case the viscous damper is in series with the (visco)plastic element. The former deforms with a strain rate $\dot{\epsilon}_v$ while the latter with $\dot{\epsilon}_{vp}$ (both under the same stress τ) and we have $\dot{\epsilon}_T = \dot{\epsilon}_v + \dot{\epsilon}_{vp}$.

$$\begin{aligned}
\dot{\epsilon}_T &= \dot{\epsilon}_v + \dot{\epsilon}_{vp} \\
&= \dot{\epsilon}_v + \frac{\tau}{2\eta_{vp}} \\
&= \dot{\epsilon}_v + \frac{\tau}{2\left(\frac{Y}{2\dot{\epsilon}_{vp}} + \eta_m\right)} \\
&= \dot{\epsilon}_v + \frac{\tau}{2\left(\frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v)} + \eta_m\right)} \\
\dot{\epsilon}_T - \dot{\epsilon}_v &= \frac{\tau}{2\left(\frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v)} + \eta_m\right)} \\
2(\dot{\epsilon}_T - \dot{\epsilon}_v) \left(\frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v)} + \eta_m \right) &= \tau \\
Y + 2(\dot{\epsilon}_T - \dot{\epsilon}_v)\eta_m &= \tau \\
Y + 2\left(\dot{\epsilon}_T - \frac{\tau}{2\eta_v}\right)\eta_m &= \tau \\
Y + (2\eta_v\dot{\epsilon}_T - \tau)\frac{\eta_m}{\eta_v} &= \tau \\
Y + 2\eta_m\dot{\epsilon}_T &= \tau\left(1 + \frac{\eta_m}{\eta_v}\right)
\end{aligned}$$

and finally

$$\tau = \frac{Y + 2\eta_m\dot{\epsilon}_T}{1 + \frac{\eta_m}{\eta_v}} \quad (2.321)$$

Note that this solution exists even when $\eta_m = 0$, and then rather logically $\tau = Y$.

4. Once we have τ , we can easily compute $\dot{\epsilon}_v = \frac{\tau}{2\eta_v}$
5. We then compute $\dot{\epsilon}_{vp} = \dot{\epsilon}_T - \dot{\epsilon}_v$ which we use to compute η_{vp} :

$$\begin{aligned}
\eta_{vp} &= \frac{Y}{2\dot{\epsilon}_{vp}} + \eta_m \\
&= \frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v)} + \eta_m \\
&= \frac{Y}{2(\dot{\epsilon}_T - \frac{\tau}{2\eta_v})} + \eta_m \\
&= \frac{Y}{2(\dot{\epsilon}_T - \frac{Y+2\eta_m\dot{\epsilon}_T}{1+\frac{\eta_m}{\eta_v}} \frac{1}{2\eta_v})} + \eta_m \\
&= \frac{Y}{2\dot{\epsilon}_T - \frac{Y+2\eta_m\dot{\epsilon}_T}{\eta_v + \eta_m}} + \eta_m \tag{2.322}
\end{aligned}$$

$$= \frac{Y(\eta_v + \eta_m)}{2(\eta_v + \eta_m)\dot{\epsilon}_T - (Y + 2\eta_m\dot{\epsilon}_T)} + \eta_m \tag{2.323}$$

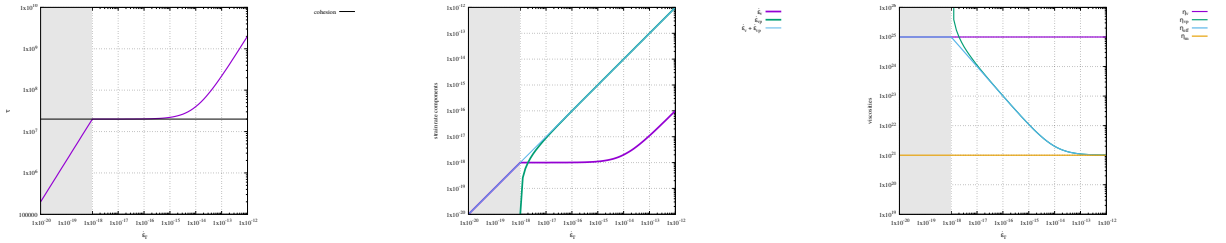
$$= \frac{Y(\eta_v + \eta_m)}{2\eta_v\dot{\epsilon}_T - Y} + \eta_m \tag{2.324}$$

$$= \frac{Y(\eta_v + \eta_m)/2\eta_v}{\dot{\epsilon}_T - Y/2\eta_v} + \eta_m \tag{2.325}$$

6. Having obtained η_{vp} we can compute the final effective viscosity

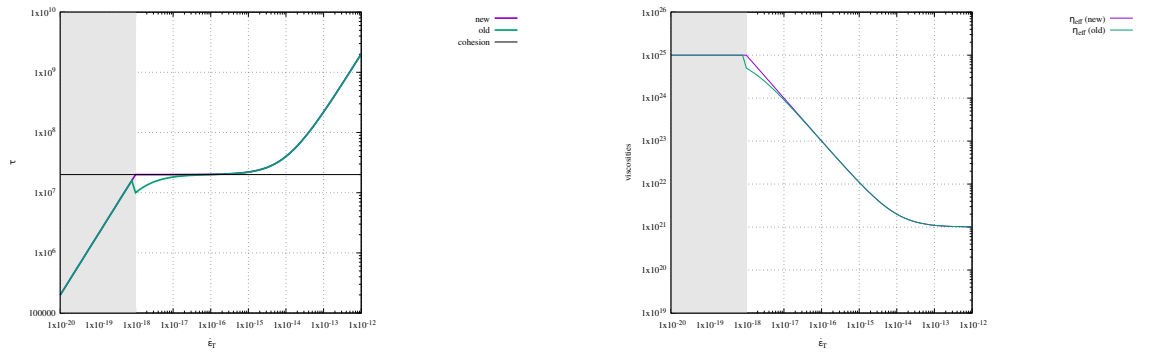
$$\eta_{eff} = \left(\frac{1}{\eta_v} + \frac{1}{\eta_{vp}} \right)^{-1}$$

On the following plots are shown τ , $\dot{\epsilon}_{vp}$, $\dot{\epsilon}_v$, η_{vp} , and η_{eff} as a function of $\dot{\epsilon}_T$:



Obtained for $\eta_m = 10^{21}$, $Y = 20\text{MPa}$ and $\eta_v = 10^{25}$. Python code in `images/rheology/vvp/`

In the following plots the resulting stress τ and effective viscosities η_{eff} are compared between the above approach ('new') and the simpler (and naive) approach where $\dot{\epsilon}_T$ is used in η_{vp} instead of $\dot{\epsilon}$ ('old'). In this particular case we see that it makes a difference at low strain rates close to the brittle-ductile transition.



Obtained for $\eta_m = 10^{21}$, $Y = 20\text{MPa}$ and $\eta_v = 10^{25}$. Python code in `images/rheology/vvp/`

Remark. The introduction of the damper η_m in parallel with the plastic element has an unavoidable effect: the stress τ becomes larger than Y at high strain rate values! Since the vp block is akin to a bingham fluid, this is no surprise.

Remark. The viscous dashpot η_v also acts as a maximum viscosity cutoff: if η_{vp} becomes (very) large, i.e. $\eta_{vp} \gg \eta_v$, then $\eta_{eff} \rightarrow \eta_v$. Conversely, if $\eta_p = Y/2\dot{\epsilon}_{vp}$ becomes (very) small, i.e. $\eta_p \ll \eta_m$ then η_m acts as a minimum viscosity limiter, i.e. $\eta_{vp} \rightarrow \eta_m$. Since $\eta_m \ll \eta_v$ then $\eta_{eff} \rightarrow \eta_m$.

A simple regularisation This idea originates in Massmeyer *et al.* (2013) [840]. We postulate

$$\tilde{\eta}_{eff} = \left(1 - \exp\left(-\frac{\dot{\epsilon}_T}{\dot{\epsilon}_T^c}\right)\right) \left(\frac{Y}{2\dot{\epsilon}_T} + \eta_m\right)$$

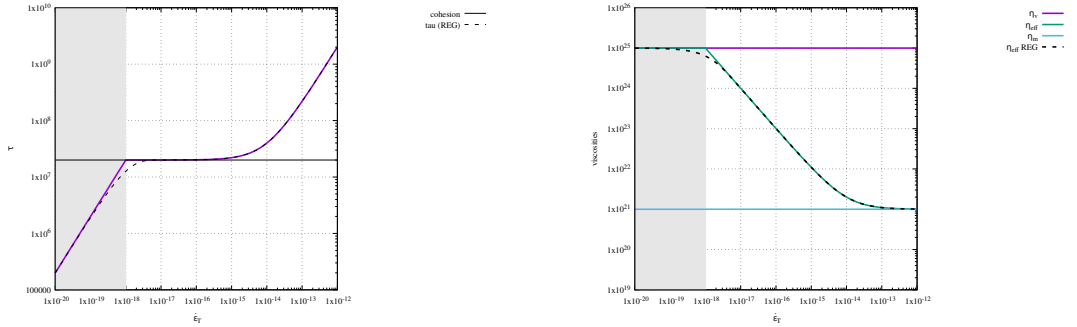
where $\dot{\epsilon}_T^c$ is the critical strain rate at which the transition viscous to viscous-viscoplastic occurs given by $\dot{\epsilon}_T^c = Y/2\eta_v$. When $\dot{\epsilon}_T \ll \dot{\epsilon}_T^c$ then the exponential term tends to zero and

$$\tilde{\eta}_{eff} \rightarrow \frac{Y}{2\dot{\epsilon}_T} + \eta_m$$

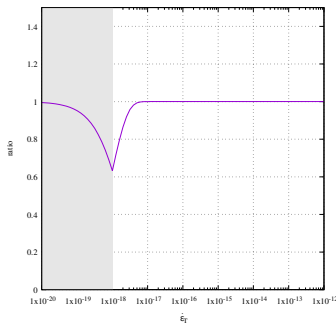
and if $\dot{\epsilon}_T \rightarrow \infty$ then $\tilde{\eta}_{eff} \rightarrow \eta_m$. Conversely if $\dot{\epsilon}_T \rightarrow 0$ then we can carry out a Taylor expansion of the exponential term ($\exp x \sim 1 + x$ when x is small).

$$\tilde{\eta}_{eff} \sim \left(\frac{\dot{\epsilon}_T}{\dot{\epsilon}_T^c}\right) \left(\frac{Y}{2\dot{\epsilon}_T} + \eta_m\right) \rightarrow \frac{\dot{\epsilon}_T}{\dot{\epsilon}_T^c} \frac{Y}{2\dot{\epsilon}_T} = \eta_v$$

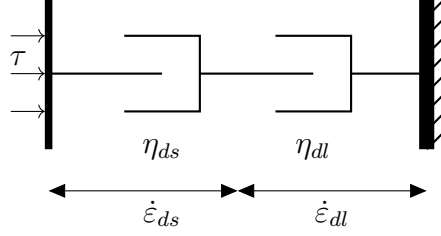
At low strain rates the viscosity does not 'explode' but actually converges to the background viscosity η_v . The stress τ corresponding to this viscosity is simply $\tilde{\tau} = 2\tilde{\eta}_{eff}$. Both $\tilde{\tau}$ and $\tilde{\eta}_{eff}$ are plotted hereunder:



Obtained for $\eta_m = 10^{21}$, $Y = 20\text{MPa}$ and $\eta_v = 10^{25}$. Python code in `images/rheology/vvp/`



- two nonlinear viscous dampers in series:



There are two dashpots in series, one accounts for dislocation creep, the other for diffusion creep. The algorithm goes then as follows:

1. Assume we know $\dot{\epsilon}_T$ (from previous iteration).
2. The dashpots are in series so

$$\dot{\epsilon}_T = \dot{\epsilon}_{ds} + \dot{\epsilon}_{df}$$

with

$$\dot{\epsilon}_{ds} = A_{ds} \tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) \quad (2.326)$$

$$\dot{\epsilon}_{df} = A_{df} \tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right) \quad (2.327)$$

such that we are in fact looking for the stress value τ so that

$$\dot{\epsilon}_T = A_{ds} \tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) + A_{df} \tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right)$$

or, we must find the zero of the function $\mathcal{F}(\tau)$:

$$\mathcal{F}(\tau) = \dot{\epsilon}_T - A_{ds} \tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) - A_{df} \tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right)$$

This equation can be solved with a Newton-Raphson algorithm and the iterations will be of the form:

$$\tau_{n+1} = \tau_n - \frac{\mathcal{F}(\tau_n)}{\mathcal{F}'(\tau_n)}$$

where the derivative of the function \mathcal{F} with respect to τ reads:

$$\mathcal{F}'(\tau) = \frac{\partial \mathcal{F}}{\partial \tau} = -A_{ds} n \tau^{n-1} \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) - A_{df} \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right)$$

Once the value of τ is found, the strain rate values of Eqs. (2.326) and (2.327) can be computed and so can the respective effective viscosities:

$$\eta_{ds} = \frac{1}{2} A_{ds}^{1/n} \dot{\epsilon}_{ds}^{\frac{1}{n}-1} \exp\left(\frac{Q_{ds} + pV_{ds}}{nRT}\right) \quad (2.328)$$

$$\eta_{df} = \frac{1}{2} A_{df}^{1/n} \exp\left(\frac{Q_{df} + pV_{df}}{RT}\right) \quad (2.329)$$

Their average effective viscosity $\tilde{\eta}_{eff}$ is given by

$$\tilde{\eta}_{eff} = \left(\frac{1}{\eta_{ds}} + \frac{1}{\eta_{df}} \right)^{-1}$$

Rather importantly, as we will see hereafter, the following variant is implemented in some codes (e.g. DOUAR , FANTOM , SOPALE , and probably many others) so as to bypass these costly Newton iterations:

1. compute η_{ds} and η_{df} with the *same* strainrate $\dot{\epsilon}_T$, pressure and temperature values
2. average them by means of an harmonic average

In this case, we have

$$\dot{\epsilon}_T = A_{df} \tau_{df} \exp \left(-\frac{Q_{df} + pV_{df}}{RT} \right) \quad \dot{\epsilon}_T = A_{ds} \tau_{ds}^n \exp \left(-\frac{Q_{ds} + pV_{ds}}{RT} \right)$$

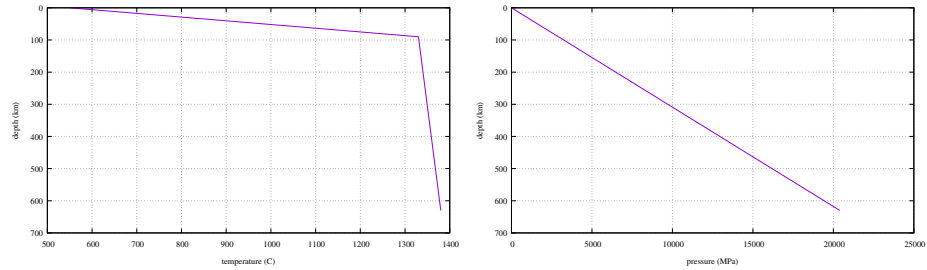
or,

$$\eta_{ds} = \frac{1}{2} A_{ds}^{1/n} \dot{\epsilon}_T^{\frac{1}{n}-1} \exp \left(\frac{Q_{ds} + pV_{ds}}{nRT} \right) \quad (2.330)$$

$$\eta_{df} = \frac{1}{2} A_{df}^{1/n} \exp \left(\frac{Q_{df} + pV_{df}}{RT} \right) \quad (2.331)$$

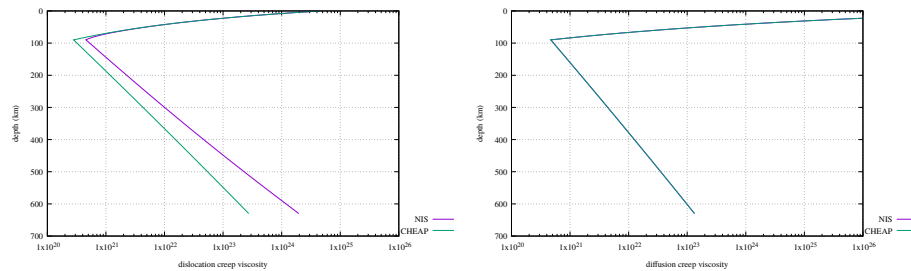
We see that this simplification has consequences on the dislocation creep viscosity only.

A concrete example Let us consider a vertical section of upper mantle, from 660km depth to 30km depth. The lithosphere is assumed to be 90km thick. The temperature at the moho (the top of the domain) is set to 550C, 1330C at the LMB and 1380C at the bottom. A constant strainrate $\dot{\epsilon}_T = 10^{-15} \text{s}^{-1}$ is assumed. We assume that the pressure is lithostatic (for simplicity the density is taken to be constant at 3300kg/m³). The temperature and pressure fields are shown hereunder:

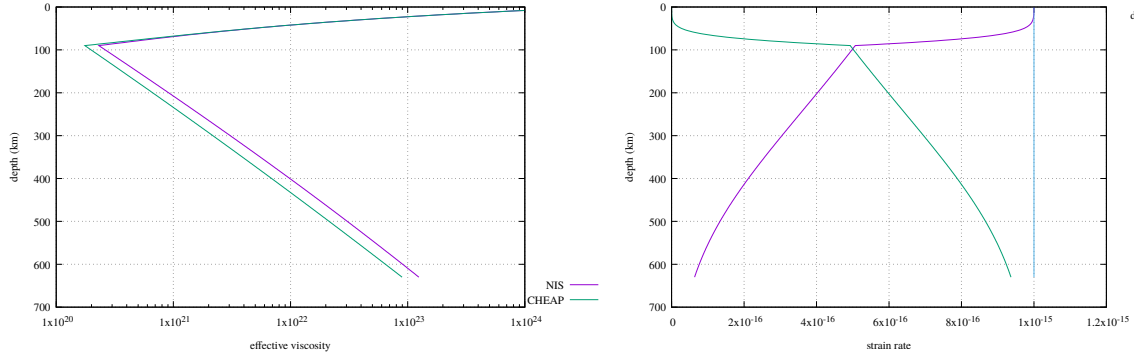


Material properties are taken from Karato & Wu (1993) [673]. The (fortran) code is available in `images/rheology/effvisc/`.

In what follows, the values obtained with Newton iterations are coined 'NR' and those obtained without are coined 'CHEAP'. The diffusion and dislocation creep viscosities can be computed for both algorithms and are shown hereunder (As mentioned earlier the diffusion creep viscosity is independent of strain rate so is the same for both):

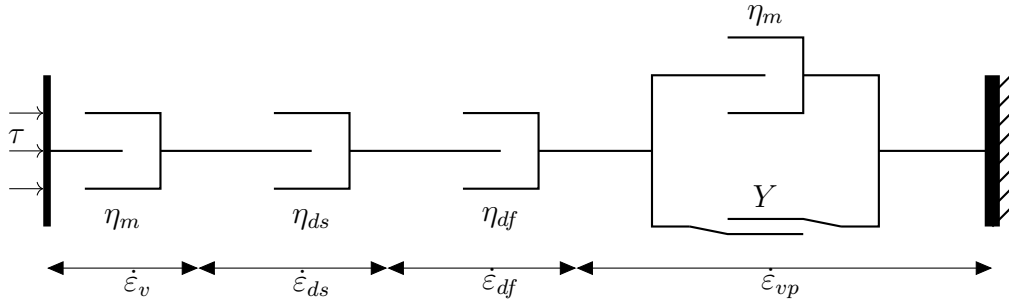


We can also plot the resulting effective viscosity η_{eff} for both approaches and we see that the differences are larger than 20%. This is shown here under on the left, alongside with the partitioning of the strain rate as a function of depth:



- multiple viscous dampers and a plastic element arranged as follows:

(tikz_vvp2.tex)



The algorithm goes then as follows:

1. Assume we know $\dot{\epsilon}_T$ (from previous iteration), as well as the plasticity parameters Y (a constant in the case of von Mises, or a pressure-dependent quantity otherwise) and η_m .
2. We start by assuming that the plasticity 'block' is not active ($\dot{\epsilon}_{vp} = 0$): we have then three dampers in series. We need their associated strain rates $\dot{\epsilon}_{df}$ and $\dot{\epsilon}_{ds}$ which are such that

$$\dot{\epsilon}_T = \dot{\epsilon}_v + \dot{\epsilon}_{ds} + \dot{\epsilon}_{df}$$

with

$$\dot{\epsilon}_v = \frac{\tau}{2\eta_v} \quad (2.332)$$

$$\dot{\epsilon}_{ds} = A_{ds} \tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) \quad (2.333)$$

$$\dot{\epsilon}_{df} = A_{df} \tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right) \quad (2.334)$$

such that we are in fact looking for the stress value τ so that

$$\dot{\epsilon}_T = A_{ds} \tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) + A_{df} \tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right) + \frac{\tau}{2\eta_v}$$

or, we must find the zero of the function \mathcal{F} :

$$\mathcal{F}(\tau) = \dot{\epsilon}_T - A_{ds}\tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) - A_{df}\tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right) - \frac{\tau}{2\eta_v}$$

This equation can be solved with a Newton-Raphson algorithm and the iterations will be of the form:

$$\tau_{n+1} = \tau_n - \frac{\mathcal{F}(\tau_n)}{\mathcal{F}'(\tau_n)}$$

where the derivative of the function \mathcal{F} with respect to τ reads:

$$\mathcal{F}'(\tau) = \frac{\partial \mathcal{F}}{\partial \tau} = -A_{df} \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right) - A_{ds}n\tau^{n-1} \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) - \frac{1}{2\eta_v}$$

Once the value of τ is found, the strain rate values of Eqs. (2.333), (2.334) and (2.332) can be computed and so can the respective effective viscosities:

$$\eta_{ds} = \frac{1}{2}A_{ds}^{1/n} \dot{\epsilon}_{ds}^{\frac{1}{n}-1} \exp\left(\frac{Q_{ds} + pV_{ds}}{nRT}\right) \quad (2.335)$$

$$\eta_{df} = \frac{1}{2}A_{df}^{1/n} \exp\left(\frac{Q_{df} + pV_{df}}{RT}\right) \quad (2.336)$$

Their average effective viscosity $\tilde{\eta}_{eff}$ is given by

$$\tilde{\eta}_{eff} = \left(\frac{1}{\eta_{ds}} + \frac{1}{\eta_{df}} + \frac{1}{\eta_v}\right)^{-1}$$

3. if $\tau = 2\tilde{\eta}_{eff}\dot{\epsilon}_T < Y$ the stress is below the yield stress value and the plasticity element is indeed not active. Use $\tilde{\eta}_{eff}$ in the material model.
4. if $\tau = 2\tilde{\eta}_{eff}\dot{\epsilon}_T > Y$ the stress is above the yield value, which is not allowed. In this case the plastic element must be present and active and the viscous dampers are then in series with the (visco)plastic element. The formers deform with a strain rate $\dot{\epsilon}_v$, $\dot{\epsilon}_{ds}$ and $\dot{\epsilon}_{df}$ while the latter with $\dot{\epsilon}_{vp}$ (all under the same stress τ) and we have $\dot{\epsilon}_T = \dot{\epsilon}_v + \dot{\epsilon}_{ds} + \dot{\epsilon}_{df} + \dot{\epsilon}_{vp}$ so:

$$\begin{aligned} \dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) - \dot{\epsilon}_{df}(\tau) &= \dot{\epsilon}_{vp} \\ &= \frac{\tau}{2\left(\frac{Y}{2\dot{\epsilon}_{vp}} + \eta_m\right)} \\ \dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) - \dot{\epsilon}_{df}(\tau) &= \frac{\tau}{2\left(\frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) + \dot{\epsilon}_{df}(\tau))} + \eta_m\right)} \\ 2[\dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) - \dot{\epsilon}_{df}(\tau)] \left(\frac{Y}{2(\dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) + \dot{\epsilon}_{df}(\tau))} + \eta_m\right) &= \tau \\ Y + 2(\dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) - \dot{\epsilon}_{df}(\tau))\eta_m &= \tau \end{aligned}$$

As before, we must find the zero of the function \mathcal{F} :

$$\begin{aligned} \mathcal{F}(\tau) &= Y + 2[\dot{\epsilon}_T - \dot{\epsilon}_v(\tau) - \dot{\epsilon}_{ds}(\tau) - \dot{\epsilon}_{df}(\tau)]\eta_m - \tau \\ &= Y + 2\left[\dot{\epsilon}_T - \frac{\tau}{2\eta_v} - A_{ds}\tau^n \exp\left(-\frac{Q_{ds} + pV_{ds}}{RT}\right) - A_{df}\tau \exp\left(-\frac{Q_{df} + pV_{df}}{RT}\right)\right]\eta_m - \tau \end{aligned}$$

Because dislocation creep involves the n -th power of the stress we will here also need to find the zero by means of a Newton-Raphson algorithm.

We have:

$$\frac{\partial \mathcal{F}}{\partial \tau} = \left[-\frac{1}{\eta_v} - 2\frac{\partial \dot{\epsilon}_{ds}(\tau)}{\partial \tau} - 2\frac{\partial \dot{\epsilon}_{df}(\tau)}{\partial \tau} \right] \eta_m - 1 \quad (2.337)$$

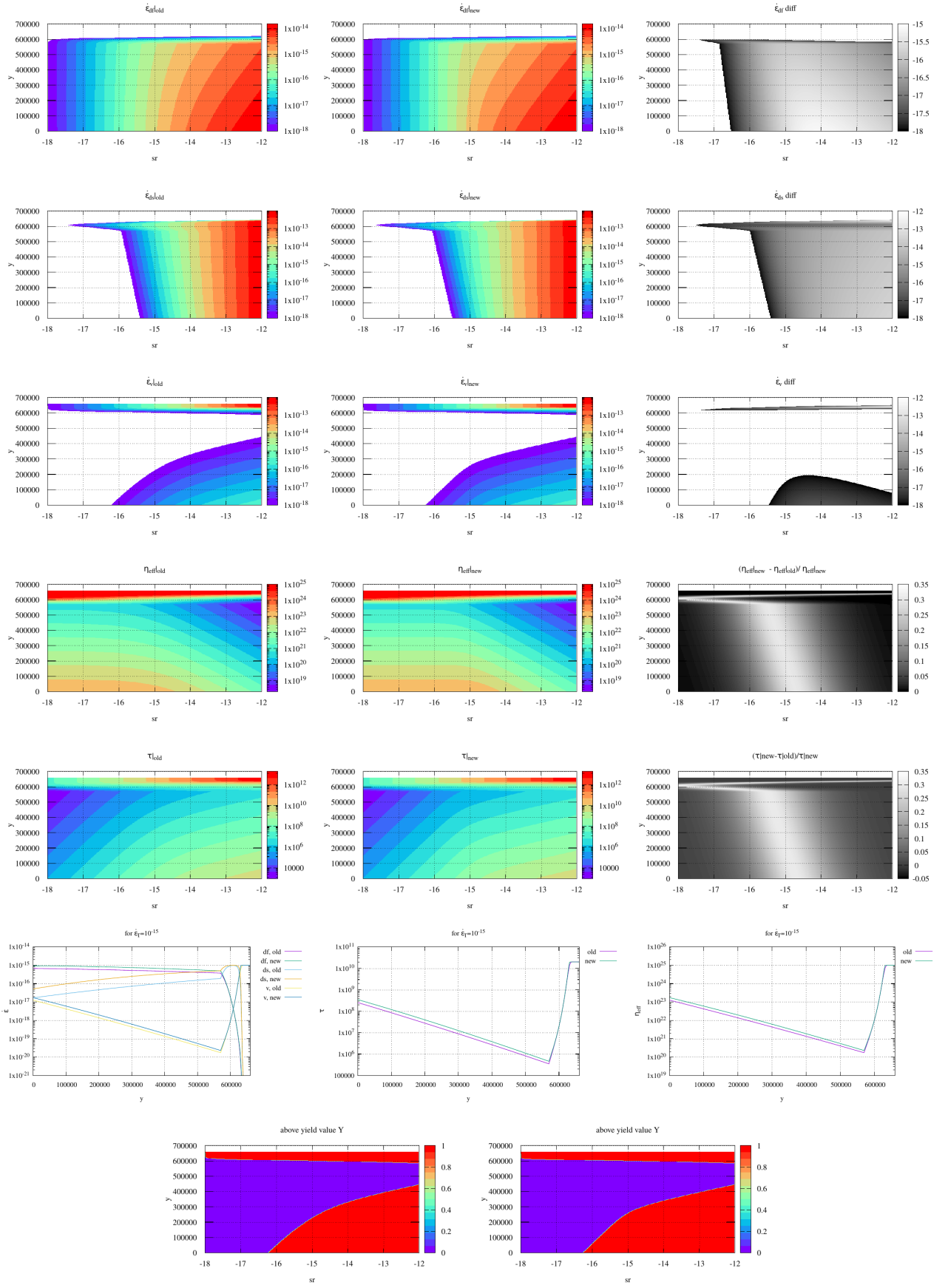
$$\begin{aligned} \mathcal{F}(\tau)/2\eta_m &= \frac{Y}{2\eta_m} + \dot{\epsilon}_T - \frac{\tau}{2\eta_v} - A_{ds}(p, T)\tau^n - A_{df}(p, T)\tau - \frac{\tau}{2\eta_m} \\ &= -A_{ds}(p, T)\tau^n - \left(A_{df}(p, T) + \frac{1}{2\eta_v} - \frac{1}{2\eta_m} \right) \tau + \left(\frac{Y}{2\eta_m} + \dot{\epsilon}_T \right) \end{aligned} \quad (2.338)$$

Note that when $\eta_m = 0$ we logically recover $\tau = Y$ as the stress cannot exceed the yield strength Y .

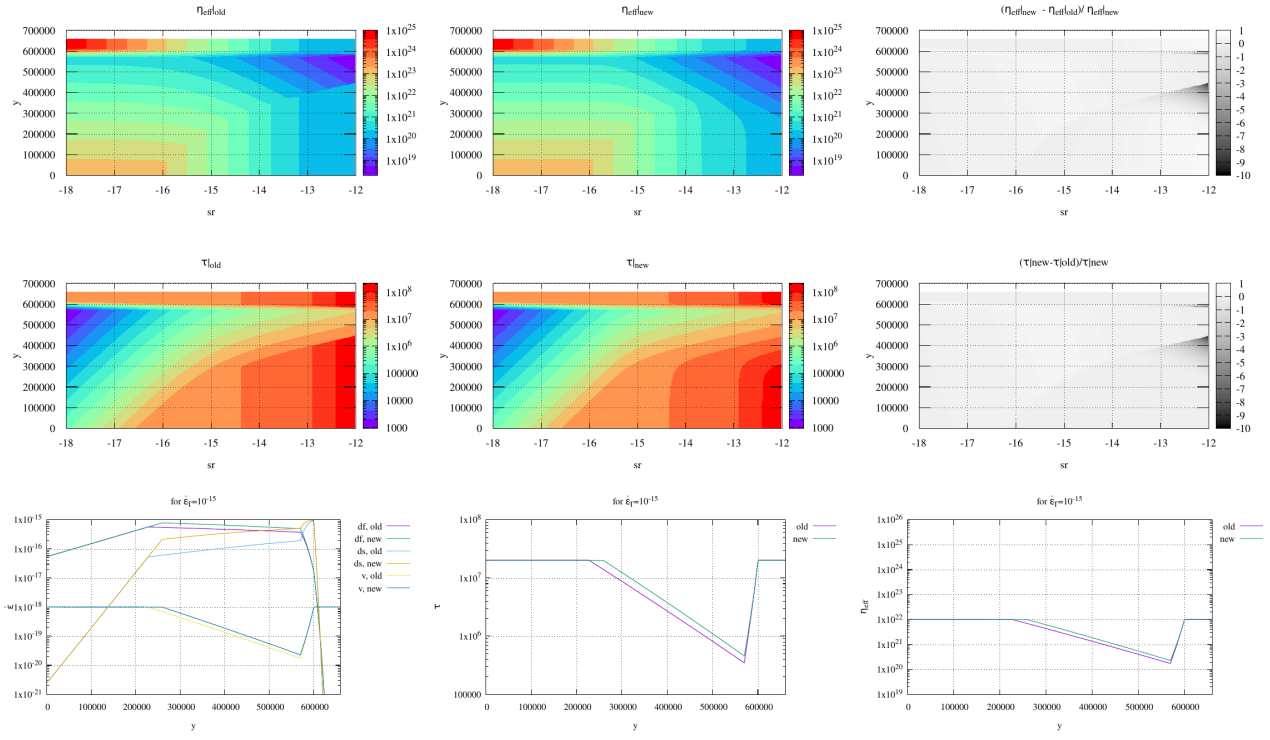
Although this approach is probably the most consistent in terms of physics, the presence of the Newton-Raphson iterations makes it very expensive since this procedure is to be repeated for every quadrature point or every particle.

Let us consider a concrete example: we set $Y = 20\text{MPa}$, $\eta_v = 10^{25}\text{pascal}$, $\eta_m = 10^{20}\text{pascal}$. The domain is one-dimensional of depth 660km. The density is assumed to be constant at 3300kg m^{-3} . Dislocation and diffusion creep parameters are taken from Karato & Wu (1993) [673]. The temperature is linear is 20°C at the surface, 550°C at 30km depth, 1330°C at 90km depth and 1380°C at the bottom. Pressure is assumed to be lithostatic. The python program and the gnuplot script are in *images/rheology/example*.

In the code I consider two cases: 'old' and 'new'. The latter is described above. 'old' goes as follows: loop over total strain rate values. Compute dislocation and diffusion creep viscosities with it. Compute harmonic average of these with linear viscosity. Compute deviatoric stress value. use it in dislocation and diffusion formulae to arrive at respective strainrates.



Viscous branch: $(ds+df+v)$ 'old' stands for the old approach when $\dot{\epsilon}_T$ was used for all mechanisms. 'new' stands for the new approach and the right strain rate decomposition.




Visco-viscoplastic rheology: (ds+df+v+vp)

Remark. *Chenin et al. (2019) [229], base their rheological model on the additive decomposition of the following deviatoric strain rate tensor ϵ^d :*

$$\epsilon^d = \epsilon^{el} + \epsilon^{pl} + \epsilon^{ds} + \epsilon^{df} + \epsilon^{pe}$$

where the five strain rate terms correspond respectively to the elastic, plastic, and viscous creep (dislocation, diffusion, peierls) contributions. This implies that all these elements are in series and the associated viscosities are then averaged with an harmonic mean. Rather interestingly, it is then stated that "this strain rate equation is nonlinear and solved locally on cell centroids and vertices in order to define the current effective viscosity and stress [1011]."

 Relevant Literature: [580, 808, 579, 1117, 805, 25, 367]

2.27.23 Anisotropic viscosity

Following the paper by Lev and Hager (2008) [776], the anisotropic viscosity enters the equation of momentum through a 'correction' term added to the isotropic part of the constitutive equation relating stress and strain rate [914]:

$$\sigma_{ij} = -p\delta_{ij} + 2\eta_N\dot{\epsilon}_{ij} - 2(\eta_N - \eta_S)\Lambda_{ijkl}\dot{\epsilon}_{kl}$$

where η_N is the normal viscosity and η_S is the shear viscosity. The fourth order tensor Λ reflects the orientation of the directors in space, denoted by \vec{n} :

$$\Lambda_{ijkl} = \frac{1}{2}(n_i n_k \delta_{lj} + n_j n_k \delta_{il} + n_i n_l \delta_{kj} + n_j n_l \delta_{ik}) - 2n_i n_j n_k n_l$$

Following [898, 914], the 'directors' are advected through the model and are analogous to particles. The directors are vector-particles pointing normal to the easy-glide plane or layer, thus defining the directions associated with η_N and η_S . In each time step of the calculation, the directors are advected and rotated by the flow, and in return determine the viscosity structure for the next time step [912].

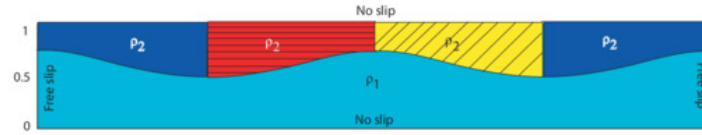


Figure 4. A schematic description of the model geometry and initial conditions. The colours denote the densities and rheologies: blue—isotropic, $\rho = 1$, $\eta_{\text{iso}} = 1$, red—anisotropic with horizontal fabric, $\rho = 1$, $\delta = 0.1$, yellow—anisotropic with dipping fabric, $\rho = 1$, $\delta = 10$, cyan—isotropic, $\rho = 0$, $\eta_{\text{iso}} = 1$. There is no slip on the top and bottom boundaries, and free slip is allowed along the side walls. The thickness of the top layer and the amplitude of the interface perturbation were exaggerated for clarity.

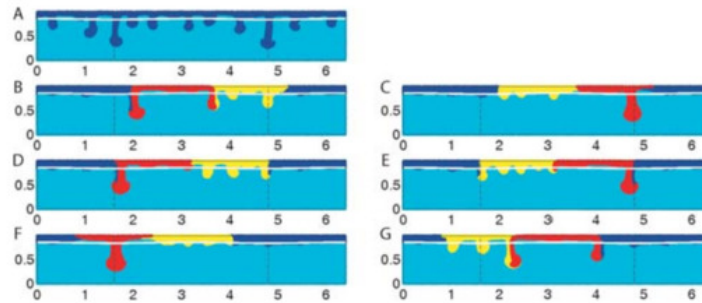


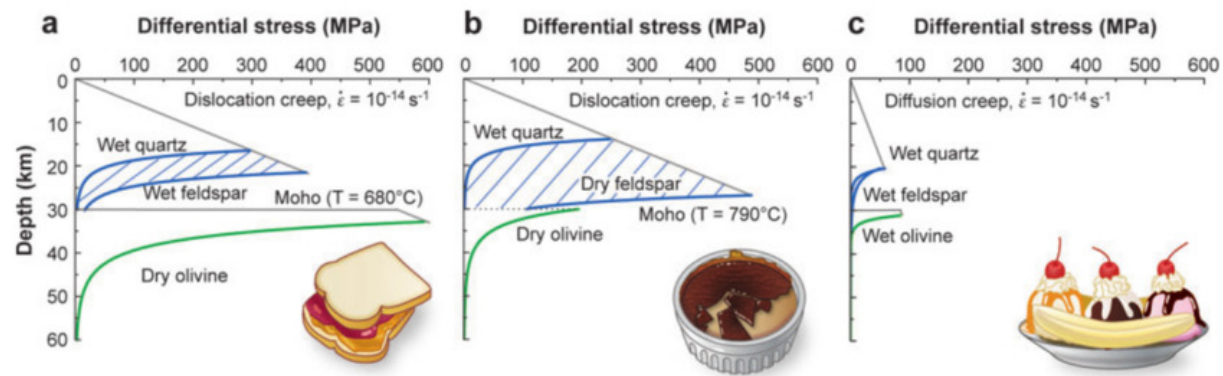
Figure 5. Material distribution in models with different configurations of initial anisotropic fabric taken after the fastest downwelling sinks over half the box depth. Panel A shows the results for an isotropic model. The black cosine curve at a depth of 0.15 marks the original interface between the dense and buoyant layers. The vertical dashed black lines show the deepest points of the original density interface, where the dense layer was thickest. Red material starts with a horizontal fabric; Yellow material starts with a fabric dipping at 45° . Blue materials are isotropic. Interestingly both panels (b) and (g), which start with distinctly different material arrangements, show large downwellings comprised of both anisotropic materials, while others do not.

Taken from Lev & Hager (2008) [776].

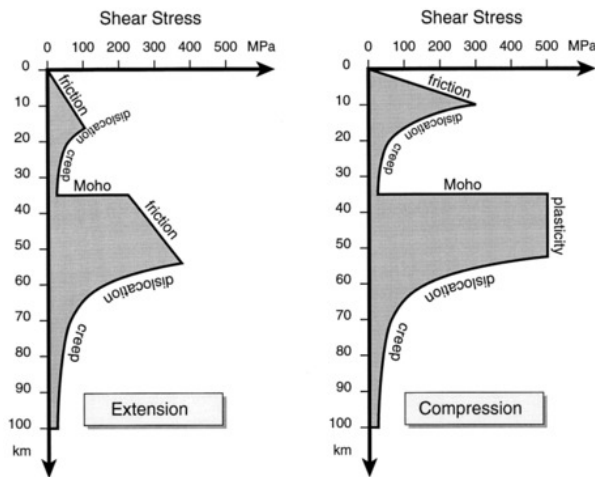
Relevant Literature:

- F.M. Richter and S.F. Daly. “Convection models having a multiplicity of large horizontal scales”. In: *J. Geophys. Res.* 83 (1978), pp. 4951–4956
- M. Saito and Y. Abe. “Consequences of anisotropic viscosity in the Earth’s mantle (in Japanese, with English Abstract)”. In: *Zisin* 37 (1984), pp. 237–245
- A. Vauchez, A. Tomassi, and G. Barroul. “Rheological heterogeneity, mechanical anisotropy and deformation of the continental lithosphere”. In: *Tectonophysics* 296 (1998), pp. 61–86. DOI: 10.1016/S0040-1951(98)00137-1
- H.-B. Mühlhaus, L. Moresi, B. Hobbs, and F. Dufour. “Large amplitude folding in finely layered viscoelastic rock structures”. In: *Pure appl. Geophys.* 159 (2002), pp. 2311–2333
- H-B Mühlhaus, L Moresi, and M Čada. “Anisotropy model for mantle convection”. In: *Computational Fluid and Solid Mechanics 2003*. 2003, pp. 1044–1046. DOI: 10.1016/B978-008044046-0.50255-4
- H-B Mühlhaus, L Moresi, and Miroslav Čada. “Emergent anisotropy and flow alignment in viscous rock”. In: *pure and applied geophysics* 161.11-12 (2004), pp. 2451–2463
- K. Michibayashi and D. Mainprice. “The role of pre-existing mechanical anisotropy on shear zone development within oceanic mantle lithosphere: an example from the Oman ophiolite”. In: *J. Petrol.* 45(2) (2004), pp. 405–414
- L Moresi and H-B Mühlhaus. “Anisotropic viscous models of large-deformation Mohr–Coulomb failure”. In: *Philosophical Magazine* 86.21-22 (2006), pp. 3287–3305. DOI: 10.1080/14786430500255419
- Hans Mühlhaus, Louis Moresi, Lutz Gross, and Joseph Grotowski. “The influence of non-coaxiality on shear banding in viscous-plastic materials”. In: *Granular Matter* 12.3 (2010), pp. 229–238. DOI: 10.1007/s10035-010-0176-9
- Hans B Mühlhaus, Jingyu Shi, Louise Olsen-Kettle, and Louis Moresi. “Effects of a non-coaxial flow rule on shear bands in viscous-plastic materials”. In: *Granular Matter* 13.3 (2011), pp. 205–210
- W. Sharples, L.N. Moresi, M. Velic, M.A. Jadamec, and D.A. May. “Simulating faults and plate boundaries with a transversely isotropic plasticity model”. In: *Phys. Earth. Planet. Inter.* 252 (2016), pp. 77–90. DOI: 10.1016/j.pepi.2015.11.007
- J. Perry-Houts and L. Karlstrom. “Anisotropic viscosity and time-evolving lithospheric instabilities due to aligned igneous intrusions”. In: *Geophysical Journal International* 216.2 (2018), pp. 794–802. DOI: 10.1093/gji/ggy466
- Á Király, Clinton P Conrad, and LN Hansen. “Evolving viscous anisotropy in the upper mantle and its geodynamic implications”. In: *Geochem. Geophys. Geosyst.* 21 (2020), e2020GC009159. DOI: 10.1029/2020GC009159

2.27.24 Rheology of the lithosphere




Schematic view of the three most common first order rheological models of the continental lithosphere under a strain rate of 10^{-14}s^{-1} . In all three models the upper crust has its frictional strength increased with pressure and depth. (a) The jelly sandwich model has a weak mid-lower crust and a strong mantle composed of dry olivine. (b) The crème brûlée model assumes that the mantle is weak, due to the presence of water and high temperature deformation, and the dry and brittle crust determines the strength of the lithosphere. (c) The banana split model assumes that the lithosphere as a whole has its strength greatly reduced due to various strain weakening and feedback processes [179]



Taken from [92]. Typical vertical distribution of maximum shear stress in continental lithosphere undergoing compressional (right) or extensional (left) strain at 10^{-15}s . Friction controls level of shear stress in upper part of crust and sometimes in mantle lithosphere; then, below brittle/ductile transition, shear stress is controlled by thermally-activated dislocation creep.

Molnar [889] discusses the validity of the Brace-Goetze strength profiles. In particular, he has this to say about the power law parameters: *The uncertainty alone in Q alone renders calculated strengths uncertain by 10 times at temperatures of about 700C. Correspondingly, that uncertainty in Q is approximately equivalent to an uncertainty of about 100C in temperature.*

| | |
|----------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Wet Quartzite | upper crust [632, 1343] upper continental crust [697, 295] lower crust [632, 1343] ocean sediment [697] |
| Dry Olivine | lithosphere [614] sublithospheric mantle [614] |
| Dry Maryland Diabase | lower crust [1343, 1344] lower continental crust [697, 295] oceanic crust [1343, 697, 1344, 295] |
| Wet Olivine | continental mantle lithosphere [1343, 1344] oceanic mantle lithosphere [1343, 1344] sublithospheric mantle [1343, 697, 1344] mantle lithosphere [697] |

 Relevant Literature[186, 179, 1043, 1042]

I need to talk about Byerlee's law. [196]

2.28 The Perzyna model

perzyna.tex

In what follows I make use of the approach and notations of Zienkiewicz and Corneau [1423] (and all the 1974-75 papers that follow) and Owen and Hinton [967].

The total strain (rate) is divided into two parts⁴¹:

$$\dot{\boldsymbol{\epsilon}} = \dot{\boldsymbol{\epsilon}}^e + \dot{\boldsymbol{\epsilon}}^{vp}$$

where $\boldsymbol{\epsilon}^e$ stands for the elastic strain tensor and $\boldsymbol{\epsilon}^{vp}$ stands for the visco-plastic strain tensor.

The yield condition is given as

$$F(\boldsymbol{\sigma}, \kappa) = \Psi(\boldsymbol{\sigma}, \dot{\boldsymbol{\epsilon}}) - Y(\kappa) = 0$$

with $F < 0$ denoting the purely elastic region, κ is a history-dependent hardening/softening parameter and $Y(\kappa)$ is a static yield stress. Ψ is a function of the stress and/or strain rate invariants.

We borrow from classical viscoplasticity theory (Perzyna [994, 993]) the idea of a plastic potential defined as $Q(\boldsymbol{\sigma})$ and write

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma \left\langle \phi(F) \right\rangle \frac{\partial Q}{\partial \boldsymbol{\sigma}} \quad (2.340)$$

where γ is a positive, possibly time-dependent fluidity parameter. Note that sometimes the pseudo-viscosity $\bar{\eta} = \gamma^{-1}$ is defined [1425] so that the equation above writes:

$$\dot{\boldsymbol{\epsilon}}^{vp} = \frac{1}{\bar{\eta}} \left\langle \phi(F) \right\rangle \frac{\partial Q}{\partial \boldsymbol{\sigma}} \quad (2.341)$$

F represents the plastic yield condition. $\phi(x)$ is a positive scalar-valued monotonic increasing function in the range $x > 0$ such that $\phi^{-1}(x)$ exists and possess similar properties in the same range. The notation $\langle \rangle$ denotes the Macaulay brackets⁴² and stands for⁴³

$$\begin{aligned} \langle \phi(x) \rangle &= \phi(x) \quad \text{if } x > 0 \\ \langle \phi(x) \rangle &= 0 \quad \text{if } x \leq 0 \end{aligned}$$

If $Q = F$ then we speak of an associative law and if $Q \neq F$ we have a non-associative situation. The tensor $\frac{\partial Q}{\partial \boldsymbol{\sigma}}$ represents the direction of plastic flow and when $F = Q$ it is a vector directed normal to the yield surface at the stress point under consideration. This is potentially problematic in the case of the Tresca and Mohr-Coulomb yield surfaces since the normal is not well defined along the apices of the surfaces (see Section 7.6 of Owen and Hinton [967]). In the non-associative case, the direction of plastic flow in the principal stress space during plastic flow is not the same as the direction of the vector normal to the yield surface.

In what follows we concentrate our attention on isotropic materials for which both F and Q can be defined in terms of stress invariants.

According to Zienkiewicz, Humpheson, and Lewis [1427] (1975): "One of the main stumbling blocks of the classical plasticity theory lay in the universal assumption, based on Drucker's postulates (Drucker and Prager, 1952), that the plastic behaviour is 'associated'. With the use of Mohr-Coulomb type yield envelopes to define the limit between states of elasticity and of continuing

⁴¹Zienkiewicz and Corneau [1423] add a third term $\boldsymbol{\epsilon}^0$ which stands for initial/autogenous strain such as due to temperature changes but I neglect it in what follows.

⁴²https://en.wikipedia.org/wiki/Macaulay_brackets

⁴³there is a difference between Zienkiewicz and Corneau [1423](1974) and Zienkiewicz and Corneau [1424](1974) wrt $>$ and \geq , and also a difference with wikipedia!

irreversible deformation, the associated behaviour manifestly contradicted observation and gave excessive dilation. It became necessary therefore to extend plasticity ideas to a ‘non-associated’ form in which the plastic potential and yield surfaces are defined separately”. At the same time, it is worth remembering that these early studies mostly dealt with plasticity in metals, and later soils, but not kilometer-scale crustal layers.

Also, the Perzyna model is not the only one, see for instance the Duvaut-Lions viscoplastic model or the Consistency model [1339, 558].

We therefore need to look into the derivative of the plastic potential Q with respect to the stress tensor. Since the potential is expressed as a function of the stress invariants $\mathcal{I}_1(\boldsymbol{\sigma})$, $\mathcal{I}_2(\boldsymbol{\tau})$ and $\theta_L(\boldsymbol{\tau})$, we then have⁴⁴:

$$\begin{aligned}
\frac{\partial Q}{\partial \boldsymbol{\sigma}} &= \frac{\partial}{\partial \boldsymbol{\sigma}} Q(\mathcal{I}_1(\boldsymbol{\sigma}), \mathcal{I}_2(\boldsymbol{\tau}), \theta_L(\boldsymbol{\tau})) \\
&= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \frac{\partial \theta_L(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\
&= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} - \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \frac{\sqrt{3}}{2 \cos 3\theta_L} \left[-\frac{3}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{5/2}} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \\
&= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \left(\frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} + \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{3}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{5/2}} \right) \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\
&\quad - \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\
&= C_1 \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + C_2 \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + C_3 \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}}
\end{aligned}$$

i.e.

$$\frac{\partial Q}{\partial \boldsymbol{\sigma}} = C_1 \mathbf{a}_1 + C_2 \mathbf{a}_2 + C_3 \mathbf{a}_3 = C_1 \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + C_2 \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + C_3 \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \quad (2.343)$$

where the $C_{1,2,3}$ coefficients depend on the plastic potential Q and the stress invariants as follows:

$$C_1 = \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \quad (2.344)$$

$$\begin{aligned}
C_2 &= \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} + \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{3}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{5/2}} \\
&= \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} - \frac{1}{2} \frac{\tan 3\theta_L}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \\
&= \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \left(\frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} - \frac{\tan 3\theta_L}{\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \right) \quad (2.345)
\end{aligned}$$

$$C_3 = -\frac{\sqrt{3}}{2 \cos 3\theta_L} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial Q}{\partial \theta_L(\boldsymbol{\tau})} \quad (2.346)$$

These are identical to those of Eq. (7.71) in Owen & Hinton⁴⁵.

⁴⁴The derivative of the Lodé angle was obtained in Section ??

⁴⁵This is not exactly true: the factor $\frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}$ is absent in their Eq. (7.71) but it is to be found in their Eq. (7.70).

$$\begin{aligned} C_1 &= \frac{\partial F}{\partial J_1}, \quad C_2 = \left(\frac{\partial F}{\partial (J_2')^{1/2}} - \frac{\tan 3\theta}{(J_2')^{1/2}} \frac{\partial F}{\partial \theta} \right), \\ C_3 &= \frac{-\sqrt{3}}{2 \cos 3\theta} \frac{1}{(J_2')^{3/2}} \frac{\partial F}{\partial \theta}. \end{aligned} \quad (7.71)$$

Note that we already have established (see Section 2.26) that

$$\mathbf{a}_1 = \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \mathbf{1} \quad (2.347)$$

$$\mathbf{a}_2 = \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} = \boldsymbol{\tau} \quad (2.348)$$

$$\mathbf{a}_3 = \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} = \boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \quad (2.349)$$

with

$$\text{Tr}[\mathbf{a}_1] = \text{tr} \left[\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} \right] = 3 \quad (2.350)$$

$$\text{Tr}[\mathbf{a}_2] = \text{tr} \left[\frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] = 0 \quad (2.351)$$

$$\text{Tr}[\mathbf{a}_3] = \text{tr} \left[\frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] = \text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau}] - 2\mathcal{I}_2(\boldsymbol{\tau}) = 2\mathcal{I}_2(\boldsymbol{\tau}) - 2\mathcal{I}_2(\boldsymbol{\tau}) = 0 \quad (2.352)$$

Then the generic form of the plastic potential derivative also reads

$$\frac{\partial Q}{\partial \boldsymbol{\sigma}} = C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 \left(\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right) \quad (2.353)$$

The momentum conservation equation that we solve is

$$-\vec{\nabla} p + \vec{\nabla} \cdot [2\eta \dot{\boldsymbol{\epsilon}}^d] + \rho \vec{g} = \vec{0}$$

so we need the deviatoric strain rate tensor. We here assume for simplicity that there is only a visco-plastic element in the system, i.e. $\dot{\boldsymbol{\epsilon}} = \dot{\boldsymbol{\epsilon}}^{vp}$. Then

$$\begin{aligned} \dot{\boldsymbol{\epsilon}}^d &= \dot{\boldsymbol{\epsilon}}^{vp} - \frac{1}{3} \text{tr}[\dot{\boldsymbol{\epsilon}}^{vp}] \mathbf{1} \\ &= \gamma \langle \phi(F) \rangle \left\{ \left(C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 (\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1}) \right) - \frac{1}{3} \text{tr} \left[C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 (\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1}) \right] \mathbf{1} \right\} \\ &= \gamma \langle \phi(F) \rangle \left\{ \left(C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 (\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1}) \right) - \frac{1}{3} 3C_1 \mathbf{1} \right\} \\ &= \gamma \langle \phi(F) \rangle \left(C_2 \boldsymbol{\tau} + C_3 (\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1}) \right) \end{aligned} \quad (2.354)$$

The $C_{1,2,3}$ coefficients have been computed in Sections 2.27.10, 2.27.11, 2.27.13 and 2.27.12, and are summarized below:

| | C_1 | $C_2 (\times \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}})$ | C_3 |
|----------------|---------------------------------------------|-------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------|
| Tresca | 0 | $2 \cos \theta_L (1 + 2 \tan \theta_L \tan 3\theta_L)$ | $\frac{\sqrt{3}}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\sin \theta_L}{\cos 3\theta_L}$ |
| von Mises | 0 | 1 | 0 |
| Mohr-Coulomb | $\frac{1}{3} \sin \phi \quad \cos \theta_L$ | $\left[(1 + 2 \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (2 \tan 3\theta_L - \tan \theta_L) \right]$ | $\frac{\sqrt{3} \sin \theta_L + \sin \phi \cos \theta_L}{2\mathcal{I}_2(\boldsymbol{\tau}) \cos 3\theta_L}$ |
| Drucker-Prager | α | 1 | 0 |

The differences with the table below taken from Owen and Hinton [967] are highlighted in blue. The difference in the von Mises simply comes from the definition of the yield value.

| Table 7.1 Constants defining the yield surface in a form suitable for numerical analysis. | | | |
|--------------------------------------------------------------------------------------------------|-------------------------|----------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------|
| Yield Criterion | C_1 | C_2 | C_3 |
| Tresca | 0 | $2 \cos \theta (1 + \tan \theta \tan 3\theta)$ | $\frac{\sqrt{3}}{J_2'} \frac{\sin \theta}{\cos 3\theta}$ |
| Von Mises | 0 | $\sqrt{3}$ | 0 |
| Mohr–Coulomb | $\frac{1}{3} \sin \phi$ | $\cos \theta [(1 + \tan \theta \tan 3\theta) + \sin \phi (\tan 3\theta - \tan \theta) / \sqrt{3}]$ | $\frac{(\sqrt{3} \sin \theta + \cos \theta \sin \phi)}{(2J_2' \cos 3\theta)}$ |
| Drucker–Prager | α | 1.0 | 0 |

Taken from Owen and Hinton [967]. This table supposedly presents all three $C_{1,2,3}$ coefficients for all four plastic potentials/yield functions (associative plasticity). This is however not the case: the C_2 column is not C_2 but $\partial F / \partial \sqrt{I_2}$!

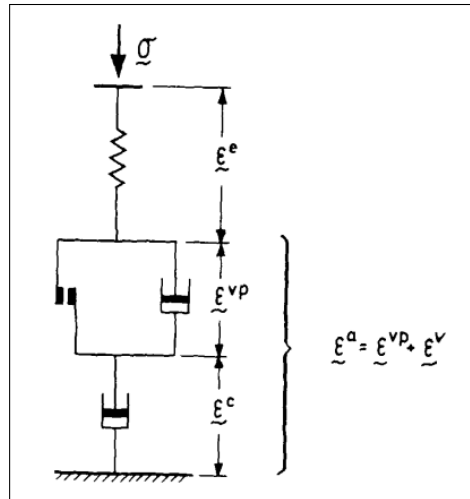
Also, the continuity equation for incompressible flow contains the divergence of the velocity field, and in this case

$$\vec{\nabla} \cdot \vec{v} = \text{tr}[\dot{\epsilon}^{vp}] = \gamma \langle \phi(F) \rangle 3C_1$$

I find it difficult to wrap my head around this as the continuity equation is usually derived by other means. If C_1 is not zero, then dilation occurs, the material is not incompressible so density should also change...

2.28.1 von Mises plasticity following Zienkiewicz (1975)

What follows is borrowed from Zienkiewicz (1975) [1422].



Taken from Zienkiewicz [1422].

We start from section 13.4.2 of the paper with the Perzyna formulation of the plastic strain ⁴⁶.

$$\dot{\epsilon}^{vp} = \gamma \langle \phi(F) \rangle \frac{\partial Q}{\partial \sigma}$$

⁴⁶from which I have removed the unnecessary/uncommon $\sqrt{3}$ terms

Associative plasticity is used, i.e. $F^{\text{vM}} = Q^{\text{vM}}$, and the von Mises yield criterion is $F^{\text{vM}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y$ so

$$\begin{aligned}\dot{\boldsymbol{\epsilon}}^{vp} &= \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{\partial(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y)}{\partial \boldsymbol{\sigma}} \\ &= \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}{\partial \boldsymbol{\sigma}} \\ &= \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}}\end{aligned}\quad (2.355)$$

Using results of Section 2.26 for the partial derivative of the second invariant we find⁴⁷:

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \quad (2.356)$$

which we can also write⁴⁸

$$\dot{\boldsymbol{\epsilon}}^{vp} = \frac{1}{2\eta} \boldsymbol{\tau} \quad \text{with} \quad \frac{1}{2\eta} = \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}$$

Note that it here follows that the flow is incompressible since the visco-plastic strain rate tensor is proportional to the deviatoric stress tensor so it is deviatoric itself!

From the definition of the second moment invariant:

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2} \boldsymbol{\tau} : \boldsymbol{\tau} = \frac{1}{2} (2\eta \dot{\boldsymbol{\epsilon}}^{vp}) : (2\eta \dot{\boldsymbol{\epsilon}}^{vp}) = 4\eta^2 \frac{1}{2} \dot{\boldsymbol{\epsilon}}^{vp} : \dot{\boldsymbol{\epsilon}}^{vp} = 4\eta^2 \mathcal{I}_2(\dot{\boldsymbol{\epsilon}}^{vp})$$

from which η can be found as a function of strain rates and hence $\boldsymbol{\Gamma}(\dot{\boldsymbol{\epsilon}})$ becomes available. Note that annoyingly the author defines the second invariant as $2\dot{\boldsymbol{\epsilon}} : \dot{\boldsymbol{\epsilon}}$ in Eq. (13.50) of the paper.

It follows that

$$\tau_e = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} = 2\eta \dot{\epsilon}_e^{vp}$$

Then we drop the $\langle \cdot \rangle$ as we assume to be above yield and we also assume a power-law form $\phi(F) = F^n$ so that we can solve explicitly for η :

$$\begin{aligned}\frac{1}{2\eta} &= \gamma \left\langle \phi \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \right\rangle \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \\ \frac{1}{2\eta} &= \gamma \left(2\eta \sqrt{\mathcal{I}_2(\dot{\boldsymbol{\epsilon}})} - Y \right)^n \frac{1}{2 \cdot 2\eta \sqrt{\mathcal{I}_2(\dot{\boldsymbol{\epsilon}})}} \\ \frac{1}{2\eta} &= \gamma (2\eta \dot{\epsilon}_e - Y)^n \frac{1}{2 \cdot 2\eta \dot{\epsilon}_e} \\ 2\dot{\epsilon}_e &= \gamma (2\eta \dot{\epsilon}_e - Y)^n \\ 2\dot{\epsilon}_e / \gamma &= (2\eta \dot{\epsilon}_e - Y)^n \\ (2\dot{\epsilon}_e / \gamma)^{1/n} &= 2\eta \dot{\epsilon}_e - Y \\ \eta &= \frac{Y + (2\dot{\epsilon}_e / \gamma)^{1/n}}{2\dot{\epsilon}_e}\end{aligned}\quad (2.357)$$

⁴⁷This is the same equation as Eq. (14) of Zienkiewicz, Jain, and Oñate [1428]

⁴⁸There most likely is a confusion in the paper between σ and τ there.

This form is convenient for plastic, visco-plastic and creep phenomena⁴⁹. This can be re-written

$$\begin{aligned}
\eta &= \frac{Y + (2\dot{\epsilon}_e/\gamma)^{1/n}}{2\dot{\epsilon}_e} \\
&= \frac{Y}{2\dot{\epsilon}_e} + \frac{(2\dot{\epsilon}_e/\gamma)^{1/n}}{2\dot{\epsilon}_e} \\
&= \frac{Y}{2\dot{\epsilon}_e} + \frac{(2/\gamma)^{1/n}}{2} \dot{\epsilon}_e^{\frac{1}{n}-1} \\
&= \frac{Y}{2\dot{\epsilon}_e} + \frac{1}{2}(\gamma/2)^{-1/n} \dot{\epsilon}_e^{\frac{1}{n}-1}
\end{aligned} \tag{2.358}$$

We often use dislocation creep/power-law rheologies and these yield an effective viscosity $\frac{1}{2}A^{-1/n}\dot{\epsilon}_e^{\frac{1}{n}-1}$ (the temperature and/or pressure-dependent exponential has been omitted for simplicity - it is a power law rheology). The equation above is then the sum of the 'plastic viscosity' and the 'viscous creep viscosity' – which corresponds to a dashpot and a plastic element in parallel. Note that the expression above is very similar to the one for Bingham or Herschel-Bulkley visco-plastic models.

If $n = 1$ then we find (as in Vilotte *et al.* (1982) [1322])

$$\eta = \frac{Y}{2\dot{\epsilon}_e} + \frac{1}{\gamma} \tag{2.359}$$

and γ is then the inverse of the (linear) viscosity of the dashpot. Also if $n = 1$ then $\bar{\eta} = \gamma^{-1}$.

For pure plasticity then $\gamma \rightarrow \infty$ and we have here simply

$$\eta = \frac{Y}{2\dot{\epsilon}_e}$$

As stated in Vilotte et al (1982) [1322]: “The plastic flow law of Eq. (2.358) permits us to represent in a single expression both the rigid-perfectly plastic flow ($\gamma \rightarrow \infty$) and the common power creep law without plastic limit ($Y = 0$) (usually referred to as the Norton-Hoff law)⁵⁰”. Vilotte, Daignieres, and Madariaga [1322] then explain that the fluidity γ can depend on temperature T in the form

$$\gamma = \gamma_0 \exp(-Q/RT)$$

In conclusion we find that this formulation allows us to represent linear, power law, perfectly plastic and visco-plastic materials. Also, we know that the additional term $\bar{\eta} = \gamma^{-1}$ introduces a length scale in the shear bands by limiting the viscosity value in said shear bands.

Remarks:

If we use a non-associative plasticity then often $Q = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}$. and the formulation of the previous section remains valid. In that case we have $C_1 = 0$ and $C_3 = 0$ which allows to easily arrive at a relationship of the type $\dot{\epsilon}^{vp} = \frac{1}{2\eta}\boldsymbol{\tau}$ where η is a scalar viscosity.

However, if Q is such that $C_3 \neq 0$, then we have a problem because even by setting $n = 1$ I do not know how to arrive at a scalar viscosity, and even thinking of η as a tensor then I am stuck, see Section about Choi & Petersen (2015).

⁴⁹Down to various $\sqrt{2}$ or $\sqrt{3}$ coefficients here or there, it is also to be found in Vilotte et al [1322, 1323, 1324]

⁵⁰<https://en.wikipedia.org/wiki/Viscoplasticity>

2.28.2 Dissecting Choi & Petersen (2015)

For implementation details, please look at [STONE](#) 39.

The original paper [237] is in 2D and focuses on the MC criterion. The authors state that the conservation of mass equation should be

$$\frac{\partial v_x}{\partial x} + \frac{\partial v_y}{\partial y} = R = 2 \sin \psi \dot{\epsilon}^p$$

where where R is the dilation rate, Ψ is the dilation angle and $\dot{\epsilon}^p$ is the square root of the second invariant of the deviatoric plastic strain rate tensor.

After multiple reads, I originally had many questions:

- where does this dilation rate R come from ?
- after reading *many* papers or textbooks on plasticity I cannot see a factor 2 in an equation anymore without re-deriving it from scratch with a coherent set of notations (preferably mine in fieldstone).
- is this relationship still valid in 3D?
- is it the same term for Drucker-Prager ?

Let us first look at their Eq. (3) in which the MC yield function is given by the function f :

$$f = \sigma_1 - N_\phi \sigma_3 - 2\sqrt{N_\phi} c$$

where σ_1 and σ_3 are the greatest and the least principal stress, $N_\phi = (1 + \sin \phi)/(1 - \sin \phi)$. This is a somewhat unusual formulation in the geodynamics community.

Let us then start with the MC yield criterion⁵¹

$$\tau_m = \sigma_m \sin \phi + c \cos \phi \quad (2.360)$$

which means that compression is assumed to be positive (the opposite as in fieldstone) and where τ_m is the magnitude of the shear stress, σ_m is the normal stress, c is the intercept of the failure envelope with the τ axis, and ϕ is the slope of the failure envelope. The quantity c is called the cohesion and the angle ϕ is called the angle of internal friction. We have

$$\sigma_m = \frac{1}{2}(\sigma_1 + \sigma_3)$$

and

$$\tau_m = \frac{1}{2}(\sigma_1 - \sigma_3)$$

Inserting these into Eq. (2.360)

$$\frac{1}{2}(\sigma_1 - \sigma_3) = \frac{1}{2}(\sigma_1 + \sigma_3) \sin \phi + c \cos \phi \quad (2.361)$$

which can be reworked as follows:

$$\sigma_1 - \frac{1 + \sin \phi}{1 - \sin \phi} \sigma_3 - 2 \frac{\cos \phi}{1 - \sin \phi} c = 0$$

⁵¹https://en.wikipedia.org/wiki/Mohr-Coulomb_theory

The third term can further be modified as follows:

$$\frac{\cos \phi}{1 - \sin \phi} = \frac{\sqrt{1 - \sin^2 \phi}}{\sqrt{(1 - \sin \phi)^2}} = \frac{\sqrt{(1 - \sin \phi)(1 + \sin \phi)}}{\sqrt{(1 - \sin \phi)^2}} = \sqrt{\frac{1 + \sin \phi}{1 - \sin \phi}}$$

Finally, we define N_ϕ as follows

$$N_\phi = \frac{1 + \sin \phi}{1 - \sin \phi}$$

so that the yield condition becomes:

$$\sigma_1 - N_\phi \sigma_3 - 2\sqrt{N_\phi} c = 0$$

which is Eq. 3 of the article by Choi & Petersen [237].

They also define the plastic potential as

$$g = \sigma_1 - N_\psi \sigma_3 = \tau_m - \sigma_m \sin \psi$$

We start again from the M-C criterion (in this case σ_2 replaces σ_3):

$$\frac{1}{2}(\sigma_1 - \sigma_2) = -\frac{1}{2}(\sigma_1 + \sigma_2) \sin \phi + c \cos \phi \quad (2.362)$$

In the case of incompressible flow I have established in Section 2.24 that

$$\frac{\sigma_1 + \sigma_2}{2} = \frac{\sigma_{xx} + \sigma_{yy}}{2} = \frac{1}{2} \mathcal{I}_1(\boldsymbol{\sigma}) \quad (2.363)$$

$$\frac{\sigma_1 - \sigma_2}{2} = \sqrt{\left(\frac{\sigma_{xx} - \sigma_{yy}}{2}\right)^2 + \sigma_{xy}^2} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad (2.364)$$

so that we now have

$$F^{\text{MC}} = \frac{1}{2} \mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - c \cos \phi$$

and then the plastic potential Q is given by

$$Q^{\text{MC}} = \frac{1}{2} \mathcal{I}_1(\boldsymbol{\sigma}) \sin \psi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}$$

We will need $\partial Q / \partial \boldsymbol{\sigma}$. By applying the chain rule we can write

$$\begin{aligned} \frac{\partial Q}{\partial \boldsymbol{\sigma}} &= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\ &= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\ &= \frac{1}{2} \sin \psi \mathbf{1} + \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \end{aligned} \quad (2.365)$$

Ultimately we would like to be able to write $\dot{\boldsymbol{\epsilon}}^{vp} = \boldsymbol{\tau} / (2\eta)$ where η is the 'viscoplastic' viscosity. However, as opposed to Zienkiewicz (1975) in the previous section, the term $\partial Q / \partial \boldsymbol{\sigma}$ is not directly/only proportional to the deviatoric stress $\boldsymbol{\tau}$ and we have instead:

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma \langle \phi(F^{\text{MC}}) \rangle \left(\frac{1}{2} \sin \psi \mathbf{1} + \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \right) \quad (2.366)$$

Right away we note that the strain rate tensor above is not deviatoric, i.e. the flow is not incompressible. Rather conveniently, the M-C criterion in plane strain can also be cast $F^{\text{MC}} = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y = \tau_e - Y$ as in the von Mises case, albeit with $Y = -\frac{1}{2}\mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + c \cos \phi$. Assuming $\phi(x) = x$ for convenience here and the argument of the brackets is positive,

$$\begin{aligned}\dot{\boldsymbol{\epsilon}}^{vp} &= \gamma \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \left(\frac{1}{2} \sin \psi \mathbf{1} + \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \right) \\ &= \gamma \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \frac{1}{2} \sin \psi \mathbf{1} + \gamma \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \\ &= \gamma (\tau_e - Y) \frac{1}{2} \sin \psi \mathbf{1} + \gamma (\tau_e - Y) \frac{1}{2\tau_e} \boldsymbol{\tau}\end{aligned}\quad (2.367)$$

If we follow the procedure of Zienkiewicz (1975), then the deviatoric part of the equation above would yield a viscosity

$$\eta = \frac{Y}{2\dot{\epsilon}_e^{vp}} + \frac{1}{\gamma} \quad \Rightarrow \quad \tau_e = 2\eta\dot{\epsilon}_e^{vp} = Y + \gamma^{-1}2\dot{\epsilon}_e^{vp} \quad \Rightarrow \quad \tau_e - Y = \gamma^{-1}2\dot{\epsilon}_e^{vp} \quad (2.368)$$

If we insert this in Eq. (2.367):

$$\begin{aligned}\dot{\boldsymbol{\epsilon}}^{vp} &= \gamma \gamma^{-1} 2\dot{\epsilon}_e^{vp} \frac{1}{2} \sin \psi \mathbf{1} + \frac{1}{2\eta} \boldsymbol{\tau} \\ &= \dot{\epsilon}_e^{vp} \sin \psi \mathbf{1} + \frac{1}{2\eta} \boldsymbol{\tau}\end{aligned}\quad (2.369)$$

Assuming that the total strain rate is the sum of the strain rates associated to the various deformation mechanisms, and that all other deformation mechanisms are deviatoric, then

$$\text{div}(\vec{\nabla}) = \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} = 2\dot{\epsilon}_e^{vp} \sin \psi$$

This is identical to the dilation rate of Choi and Petersen [237]!

2.28.3 my take on this in 3D for Drucker-Prager

I have established in Section ?? that in the general 3D case

$$F^{\text{DP}} = \alpha(\phi, c)\mathcal{I}_1(\boldsymbol{\sigma}) + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} + k(\phi, c) = \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \quad (2.370)$$

with α and k being functions of the cohesion c and angle of friction ϕ (but not from the stress). Then the plastic potential is

$$Q^{\text{DP}} = \alpha(\psi, c)\mathcal{I}_1(\boldsymbol{\sigma}) + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \quad (2.371)$$

where ψ is the dilation angle. We then have

$$\frac{\partial Q}{\partial \boldsymbol{\sigma}} = \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + \frac{\partial Q}{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \frac{\partial \sqrt{\mathcal{I}_2(\boldsymbol{\tau})}}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} = \alpha(\psi, c) \mathbf{1} + \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \quad (2.372)$$

Then

$$\begin{aligned}\dot{\boldsymbol{\epsilon}}^{vp} &= \gamma \left(\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y \right) \left(\alpha(\psi, c) \mathbf{1} + \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \right) \\ &= \gamma (\tau_e - Y) \alpha(\psi, c) \mathbf{1} + \gamma (\tau_e - Y) \frac{1}{2\tau_e} \boldsymbol{\tau}\end{aligned}\quad (2.373)$$

Using again $\tau_e - Y = \gamma^{-1}2\dot{\epsilon}_e^{vp}$ as in the 2D case we arrive finally

$$\text{div}(\vec{\nabla}) = \text{tr}[\dot{\boldsymbol{\epsilon}}^{vp}] = \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} = 6\alpha(\psi, c)\dot{\epsilon}_e^{vp}$$

Since k does not depend on stress, the only difference between the associative and the non-associative case is whether $\phi = \psi$ or not.

2.28.4 my take on this in 3D for MC

I have established in fieldstone that in the general 3D case

$$F^{\text{MC}} = \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \sin \phi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta_L(\boldsymbol{\tau}) - \frac{1}{\sqrt{3}} \sin \theta_L(\boldsymbol{\tau}) \sin \phi \right) - c \cos \phi \quad (2.374)$$

Note that since $p = -\mathcal{I}_1(\boldsymbol{\sigma})/3$ then we recover the usual ' $p \sin \phi + c \cos \phi$ '.

Following Eq. (4) of the paper the plastic potential would be given by

$$Q^{\text{MC}} = \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \sin \psi + \sqrt{\mathcal{I}_2(\boldsymbol{\tau})} \left(\cos \theta_L(\boldsymbol{\tau}) - \frac{1}{\sqrt{3}} \sin \theta_L(\boldsymbol{\tau}) \sin \psi \right)$$

The visco-plastic strain rate would then write

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma \langle \phi(F^{\text{MC}}) \rangle \frac{\partial Q^{\text{MC}}}{\partial \boldsymbol{\sigma}}$$

We have established that

$$\begin{aligned} \frac{\partial Q}{\partial \boldsymbol{\sigma}} &= C_1 \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} + C_2 \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + C_3 \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\ &= C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 \left(\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right) \end{aligned} \quad (2.375)$$

and in the case of the Mohr-Coulomb criterion:

$$C_1^{\text{MC}} = \frac{1}{3} \sin \phi \quad (2.376)$$

$$C_2^{\text{MC}} = \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \cos \theta_L \left[(1 + 2 \tan \theta_L \tan 3\theta_L) + \frac{1}{\sqrt{3}} \sin \phi (2 \tan 3\theta_L - \tan \theta_L) \right]$$

$$C_3^{\text{MC}} = \frac{\sqrt{3} \sin \theta_L + \sin \phi \cos \theta_L}{2\mathcal{I}_2(\boldsymbol{\tau}) \cos 3\theta_L} \quad (2.377)$$

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma \langle \phi(F) \rangle \left(C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 \left(\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right) \right)$$

Assuming brackets ok, and $\phi(x) = x^n$:

$$\dot{\boldsymbol{\epsilon}}^{vp} = \gamma (\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y)^n \left(C_1 \mathbf{1} + C_2 \boldsymbol{\tau} + C_3 \left(\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right) \right)$$

Taking the deviatoric part of this:

$$\dot{\boldsymbol{\epsilon}}^{vp,d} = \gamma (\sqrt{\mathcal{I}_2(\boldsymbol{\tau})} - Y)^n \left(C_2 \boldsymbol{\tau} + C_3 \left(\boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right) \right)$$

so I cannot find a scalar η such that

$$\dot{\boldsymbol{\epsilon}}^{vp,d} = \frac{1}{2\eta} \boldsymbol{\tau}$$

I am STUCK here?!

2.28.5 Revisiting Lemiale et al (2008) and Spiegelman et al (2016)

The authors postulate that the total strain rate is the sum of the viscous deformation and plastic deformation:

$$\dot{\epsilon} = \dot{\epsilon}^v + \dot{\epsilon}^p$$

Then

$$\sigma = -p\mathbf{1} + 2\eta(\dot{\epsilon} - \dot{\epsilon}^p) \quad (2.378)$$

Immediately we see that they implicitly assume that the flow is incompressible. Upon yielding a flow rule is needed to specify the plastic behaviour. The plastic strain rate is written as

$$\dot{\epsilon}^p = \dot{\lambda} \frac{\partial Q}{\partial \sigma} \quad (2.379)$$

where $\dot{\lambda}$ is a scalar plastic flow rate and Q is the so-called plastic potential. Note that in Heeres, Suiker, and de Borst [558] the authors define $\dot{\lambda} = \langle \phi(x) \rangle / \eta$ so that the equation above is the Perzyna model. A classical choice for Q , in conjunction with the incompressibility constraint, is:

$$Q = \sqrt{\mathcal{I}_2(\tau)}$$

We notice that Eq. 2.379 is different (although obviously not unrelated) than the Perzyna approach above, although in the end they arrive at a similar expression as we did before for the von Mises case.

If we consider only the deviatoric part of the stress tensor in Eq. (2.378), we thus obtain⁵²

$$\tau = 2\eta(\dot{\epsilon} - \dot{\epsilon}^p) = 2\eta \left(\dot{\epsilon} - \dot{\lambda} \frac{\partial Q}{\partial \sigma} \right) = 2\eta \left(\dot{\epsilon} - \dot{\lambda} \frac{1}{2\sqrt{\mathcal{I}_2(\tau)}} \frac{\partial \mathcal{I}_2(\tau)}{\partial \sigma} \right) = 2\eta \left(\dot{\epsilon} - \dot{\lambda} \frac{1}{2\sqrt{\mathcal{I}_2(\tau)}} \tau \right) \quad (2.380)$$

This equation can be written as

$$\left(1 + \dot{\lambda} \frac{\eta}{\sqrt{\mathcal{I}_2(\tau)}} \right) \tau = 2\eta \dot{\epsilon}$$

One can then take the square root of the second invariant of this equation:

$$\left(1 + \dot{\lambda} \frac{\eta}{\sqrt{\mathcal{I}_2(\tau)}} \right) \sqrt{\mathcal{I}_2(\tau)} = 2\eta \sqrt{\mathcal{I}_2(\dot{\epsilon})}$$

Then

$$\left(1 + \dot{\lambda} \frac{\eta}{\tau_e} \right) \tau_e = 2\eta \dot{\epsilon}_e$$

so that

$$\dot{\lambda} = \frac{2\eta \dot{\epsilon}_e - \tau_e}{\eta} = 2\dot{\epsilon}_e - \frac{\tau_e}{\eta}$$

⁵²do they assume varepsilon deviatoric too?

Finally we can insert this expression of $\dot{\lambda}$ in Eq. (2.380)

$$\boldsymbol{\tau} = 2\eta \left(\dot{\boldsymbol{\epsilon}} - \dot{\lambda} \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \right) \quad (2.381)$$

$$= 2\eta \left(\dot{\boldsymbol{\epsilon}} - \left(2\dot{\epsilon}_e - \frac{\tau_e}{\eta} \right) \frac{1}{2\sqrt{\mathcal{I}_2(\boldsymbol{\tau})}} \boldsymbol{\tau} \right) \quad (2.382)$$

$$= 2\eta \left(\dot{\boldsymbol{\epsilon}} - \left(2\dot{\epsilon}_e - \frac{\tau_e}{\eta} \right) \frac{1}{2\tau_e} \boldsymbol{\tau} \right) \quad (2.383)$$

$$= 2\eta \dot{\boldsymbol{\epsilon}} - 2\eta \dot{\epsilon}_e \frac{1}{\tau_e} \boldsymbol{\tau} + \eta \frac{\tau_e}{\eta} \frac{1}{\tau_e} \boldsymbol{\tau} \quad (2.384)$$

$$= 2\eta \dot{\boldsymbol{\epsilon}} - 2\eta \dot{\epsilon}_e \frac{1}{\tau_e} \boldsymbol{\tau} + \boldsymbol{\tau} \quad (2.385)$$

$$(2.386)$$

The term $\boldsymbol{\tau}$ is present on both sides of the equal sign so it cancels out and we are left with:

$$\mathbf{0} = 2\eta \dot{\boldsymbol{\epsilon}} - 2\eta \dot{\epsilon}_e \frac{1}{\tau_e} \boldsymbol{\tau}$$

or,

$$\boldsymbol{\tau} = \frac{\tau_e}{\dot{\epsilon}_e} \dot{\boldsymbol{\epsilon}}$$

On yield we have $\tau_e = Y(c, \phi)$ so in the end:

$$\boldsymbol{\tau} = 2 \underbrace{\frac{1}{2} \frac{Y(c, \phi)}{\dot{\epsilon}_e}}_{\eta_p} \dot{\boldsymbol{\epsilon}}$$

That last step is poorly documented in the paper! This is a cumbersome exercise and it relies heavily on the choice of $Q = \tau_e$.

Remark. *Lemiale et al. (2008) define $\bar{\tau} = \sqrt{\tau_{ij}\tau_{ij}}/2$ but define $\dot{\gamma} = \sqrt{2D_{ij}D_{ij}}$!*

Let us now turn to Spiegelman, May, and Wilson [1187] (2016). In Section 2.1.1 of this paper the authors follow the same path as above. They assume $Q = \tau_e$ but justify their choice by stating that “The use of incompressible materials mandates that we use a plastic potential g which is not a function of the pressure p ”

This is indeed very important in the context of our incompressible calculations in geodynamics.

They define the yield surface, $F(\boldsymbol{\sigma})$ which is a scalar function defining the failure (yield) state of a material. Yield surfaces are assumed to be of the following form

$$F(\boldsymbol{\sigma}) = \tau_e - Y(\boldsymbol{\sigma})$$

where Y is the yield criterion. The authors state that “it is common practice in geodynamics to define the plastic multiplier $\dot{\lambda}$ which exactly satisfies $F = 0$, or equivalently $\tau_e = Y$ [764]”. We see that it is then the same as the Lemiale *et al.* paper.

2.29 Moment of inertia

Consider a rigid body rotating with fixed angular velocity ω about an axis which passes through the origin. Let \mathbf{r}_i be the position vector of the i th mass element, whose mass is m_i . We expect this position vector to precess about the axis of rotation (which is parallel to ω) with angular velocity ω .

$$\frac{d\mathbf{r}_i}{dt} = \omega \times \mathbf{r}_i.$$

Thus, the above equation specifies the velocity, $\mathbf{v}_i = d\mathbf{r}_i/dt$, of each mass element as the body rotates with fixed angular velocity ω about an axis passing through the origin.

The total angular momentum of the body (about the origin) is written

$$\mathbf{L} = \sum_{i=1,N} m_i \mathbf{r}_i \times \frac{d\mathbf{r}_i}{dt} = \sum_{i=1,N} m_i \mathbf{r}_i \times (\omega \times \mathbf{r}_i) = \sum_{i=1,N} m_i [r_i^2 \omega - (\mathbf{r}_i \cdot \omega) \mathbf{r}_i]$$

The above formula can be written as a matrix equation of the form

$$\begin{pmatrix} L_x \\ L_y \\ L_z \end{pmatrix} = \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} \begin{pmatrix} \omega_x \\ \omega_y \\ \omega_z \end{pmatrix}$$

where

$$\begin{aligned} I_{xx} &= + \sum_{i=1,N} (y_i^2 + z_i^2) m_i = \int (y^2 + z^2) dm = \int_V (y^2 + z^2) \rho(x, y, z) dV \\ I_{yy} &= + \sum_{i=1,N} (x_i^2 + z_i^2) m_i = \int (x^2 + z^2) dm = \int_V (x^2 + z^2) \rho(x, y, z) dV \\ I_{zz} &= + \sum_{i=1,N} (x_i^2 + y_i^2) m_i = \int (x^2 + y^2) dm = \int_V (x^2 + y^2) \rho(x, y, z) dV \\ I_{xy} = I_{yx} &= - \sum_{i=1,N} x_i y_i m_i = - \int x y dm = - \int x y \rho(x, y, z) dV \\ I_{yz} = I_{zy} &= - \sum_{i=1,N} y_i z_i m_i = - \int y z dm = - \int y z \rho(x, y, z) dV \\ I_{xz} = I_{zx} &= - \sum_{i=1,N} x_i z_i m_i = - \int x z dm = - \int x z \rho(x, y, z) dV \end{aligned}$$

Here, I_{xx} is called the moment of inertia about the x -axis, I_{yy} the moment of inertia about the y -axis, I_{xy} the xy product of inertia, I_{yz} the yz product of inertia, etc. The matrix of the I_{ij} values is known as the moment of inertia tensor.

In general, the angular momentum vector, \mathbf{L} points in a different direction to the angular velocity vector, ω . In other words, \mathbf{L} is generally not parallel to ω .

Finally, although the above results were obtained assuming a fixed angular velocity, they remain valid at each instant in time if the angular velocity varies.

In the simplified case of a spherically symmetric planet, it is easy to see that $I_{xx} = I_{yy} = I_{zz}$ so that $I = \frac{1}{3}(I_{xx} + I_{yy} + I_{zz})$, and $\rho = \rho(r)$ with $dV = 4\pi r^2 dr$, leading to

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr$$

Assuming further that the planet has a constant density ρ_0 , we obtain

$$I = \frac{8\pi}{3} \rho_0 \int_0^R r^4 dr = \frac{8\pi}{3} \rho_0 \frac{R^5}{5} = \frac{2}{5} M R^2$$

where M is the mass of the planet and R is its radius.

Assuming now that the planet is composed of a core of radius R_c and density ρ_c surrounded by a mantle of density ρ_m , we have

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr = \frac{8\pi}{3} \left(\int_0^{R_c} \rho_c r^4 dr + \int_{R_c}^R \rho_m r^4 dr \right) = \frac{8\pi}{15} (\rho_c R_c^5 + \rho_m (R^5 - R_c^5))$$

The moment of inertia of the core is given in Table 2 of "Core Dynamics", Treatise on Geophysics, edited by Peter Olson: $I_{core} = 9.2 \times 10^{36} kg.m^2$. The total moment of inertia for the Earth is then given by $I = I_{core} + I_{mantle}$.

2.30 The need for numerical modelling

The governing equations we have seen in this chapter require the use of numerical solution techniques for three main reasons:

- the advection term in the energy equation couples velocity and temperature;
- the constitutive law (the relationship between stress and strain rate) often depends on velocity (or rather, strain rate), temperature, pressure, ...
- Even when the coefficients of the PDE's are linear, often their spatial variability, coupled to potentially complex domain geometries prevent arriving at the analytical solution.

Also we often have to deal with additional challenges:

- Complex geometries
- Multiphysics
- Many scales in space and time

Note that in CFD one makes a distinction between verification and validation. Simply put [1076]:

- verification: "solving the equations right"
- validation: "solving the right equations"

2.31 Important mathematical concepts and equations

mathematics.tex

2.31.1 Taylor expansion

$$f(a+h) = f(a) + hf'(a) + \frac{h^2}{2!}f''(a) + \cdots + \frac{h^{n-1}}{(n-1)!}f^{(n-1)}(a) + \frac{h^n}{n!}f^{(n)}(a) + \cdots$$

2.31.2 Divergence theorem

This is also coined the Green-Ostrogradski theorem. For a volume V bound by a surface Γ :

$$\iint_{\Gamma} \vec{V} \cdot d\vec{S} = \iiint_V \vec{\nabla} \cdot \vec{V} \, dV \quad (2.387)$$

Chapter 3

The Finite Difference Method

3.1 Back to basics: what is a derivative?

Before we start with the basics of the Finite Difference method, we should quickly recall the definition of the derivative of a function.

The derivative of a function $y = f(x)$ of a variable x is a measure of the rate at which the value y of the function changes with respect to the change of the variable x . It is called "the derivative of f with respect to x ". If x and y are real numbers, and if the graph of f is plotted against x , the derivative is the slope of this graph at each point.

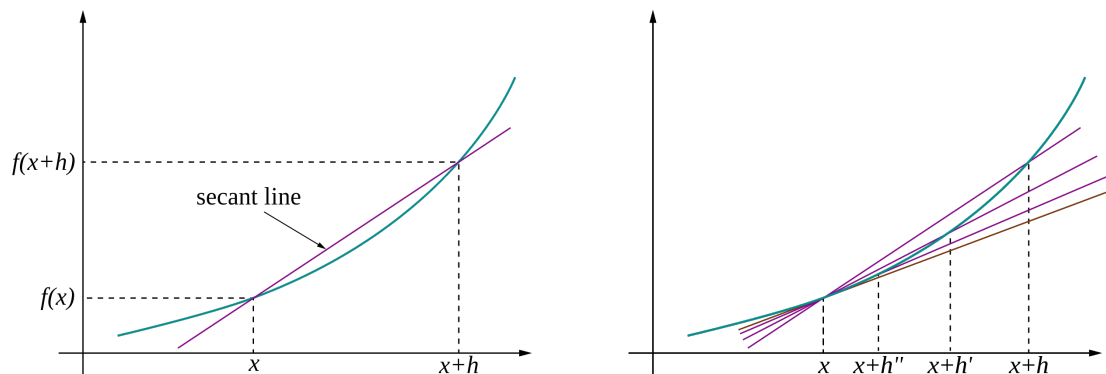
There are two standard notations:

$$\frac{df}{dx}(x) \quad \text{and} \quad f'(x) \quad (3.1)$$

The mathematical definition is¹:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{(x+h) - x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (3.2)$$

i.e., how much does the function grow between x and $x+h$, divided by the length h . On the following left plot the function is in green while the line joining the points $x, f(x)$ and $x+h, f(x+h)$ is shown in purple. On the right figure we see that when h becomes smaller and smaller this line does indeed get closer and closer to the real tangent line in brown.



Left: The secant to curve $y = f(x)$ determined by points $(x, f(x))$ and $(x+h, f(x+h))$;

Right: The tangent line as limit of secants ($h'' < h' < h$). Taken from Wikipedia².

¹if the limit exists

²<https://en.wikipedia.org/wiki/Derivative>

Also, one can rewrite the formula above as

$$f(x+h) \simeq f(x) + f'(x)h \quad (3.3)$$

which is in fact the beginning of the Taylor expansion of the function:

$$f(x+h) \simeq f(x) + f'(x)h + \frac{1}{2!}f''(x)h^2 + \dots \quad (3.4)$$

We will see in what follows that the Taylor expansion and the concept of derivation is central in the Finite Difference method.

3.2 Welcome to the discrete world

discrete.tex

In mathematics and in physics we assume that space and time form a continuum, i.e. a segment can be divided in two indefinitely, or in other words one can 'zoom in' on a part of space or time as much as needed. However computers are binary machines with a finite amount of memory so they cannot represent a continuum (the transistors in a processor are either 'on' or 'off', nothing in between). As a consequence, in the context of solving PDEs describing physical phenomena, computers will only allow us to compute the solution at certain discrete locations and at certain discrete times.

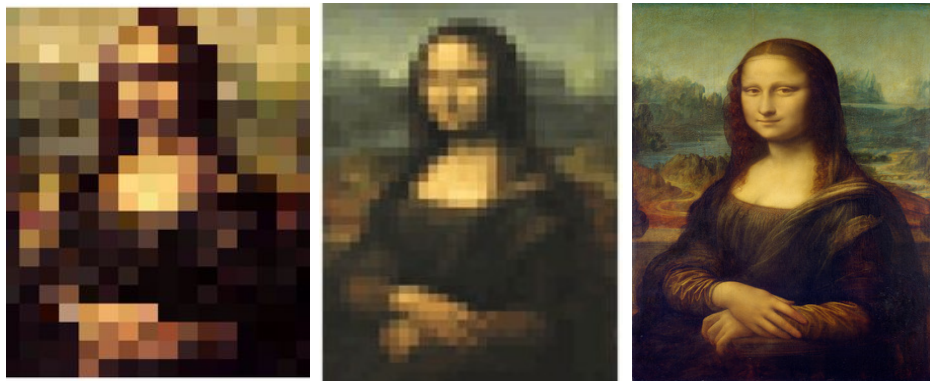


Illustration of space discretisation. If the picture is cut up in a small number of cells covering it we obtain the picture on the left. If more cells are used, we obtain the picture in the middle. When a (very) large number of cells is used we finally see all details and it approaches a continuum.

3.3 FDM basics in 1D

In what follows we suppose that we have a function $f(x)$, which is continuous and differentiable over the range of interest. Let us also assume that we know the value $f(x_0)$ and all the derivatives at $x = x_0$.

First order derivatives

The forward Taylor-series expansion for $f(x_0 + h)$, away from the point x_0 by a small amount h is given by

$$f(x_0 + h) = f(x_0) + h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \dots + \frac{h^n}{n!} \frac{\partial^n f}{\partial x^n}(x_0) + \mathcal{O}(h^{n+1}) \quad (3.5)$$

We can subtract $f(x_0)$ to each side of the equation and divide by h :

$$\frac{1}{h} (f(x_0 + h) - f(x_0)) = \frac{\partial f}{\partial x}(x_0) + \frac{h}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \dots \quad (3.6)$$

and we can then express the first derivative of f as follows:

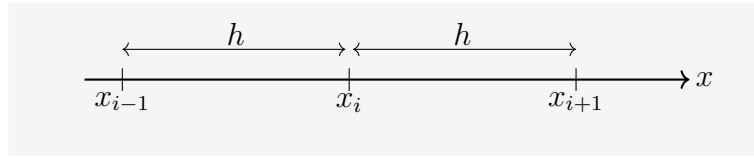
$$\frac{\partial f}{\partial x}(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} - \frac{h}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) \dots \quad (3.7)$$

or, replacing the term in h by $\mathcal{O}(h)$:

$$\boxed{\frac{\partial f}{\partial x}(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} + \mathcal{O}(h)} \quad (3.8)$$

$\mathcal{O}(h)$ indicates that the full solution would require additional terms of order h , h^2 , and so on. \mathcal{O} is called the **truncation error**: if the distance h is made smaller and smaller, the (numerical approximation) error decreases $\propto h$ in this case.

Let us assume that the 1D domain on which a given ODE/PDE is to be solved has been discretised and let us zoom in on three consecutive points ($x_0 = x_i$ here):



In the context of a discrete calculation on a set of discrete points x_i we can compute the first order derivative of f at point x_i as an approximation:

$$\boxed{\frac{\partial f}{\partial x}(x_i) = \frac{f_{i+1} - f_i}{h} + \mathcal{O}(h)} \quad (\text{forward difference}) \quad (3.9)$$

where functions $f_i = f(x_i)$ are evaluated at discretely spaced x_i with $x_{i+1} = x_i + h$ (i.e. $h = x_{i+1} - x_i$), where the node spacing, or resolution, h is assumed constant. We also introduce the notation $f'_i = f'(x_i) = \frac{\partial f}{\partial x}(x_i)$.

The **forward FD derivative** as expressed above is called **first order accurate**, and this means that very small h is required for an accurate solution.

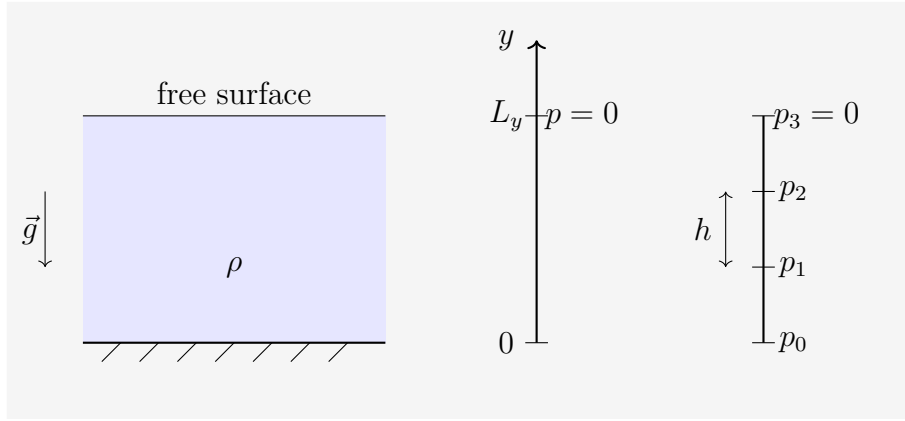
Example FDM-1: Before we go any further with the theory, let us look at a very simple example. Let us consider the Stokes equations in the absence of fluid motion (i.e. $\vec{v} = \vec{0}$). Then the strain rate tensor components are identically zero and the equation simply is

$$-\vec{\nabla}p + \rho\vec{g} = \vec{0}. \quad (3.10)$$

where we assume that 1) density is constant in space for simplicity and 2) the domain is infinite in the x -direction. The gravity vector is $\vec{g} = -g\vec{e}_y$ so that the above equation becomes:

$$-\frac{dp}{dy} - \rho g = 0 \quad (3.11)$$

This is a first-order ODE. It needs to be supplemented by a single boundary condition, which in this case constrains the pressure to be zero at the surface, i.e. $p(y = L_y) = 0$.



We can write the discretised ODE at node 2 (since we know $p_3 = 0$):

$$\frac{dp}{dy}(x_2) \simeq \frac{p_3 - p_2}{h} = -\rho g \quad (3.12)$$

or, $p_2 = \rho gh$. Having obtained p_2 , we write the ODE at node 1:

$$\frac{dp}{dy}(x_1) \simeq \frac{p_2 - p_1}{h} = -\rho g \quad (3.13)$$

or, $p_1 = \rho gh + p_2 = \rho g 2h$. And finally we obtain as expected

$$p_0 = \rho g 3h = \rho g L_y. \quad (3.14)$$

This brings us to our first exercise:



Exercise FDM-1

We will now put the previous example into practice and write a python code which uses forward differences to compute the 1D pressure field inside the crust.

→ `Exercise_1_FDM.ipynb`

We can also expand the Taylor series backward (i.e. looking 'left' of x_0)

$$f(x_0 - h) = f(x_0) - h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) - \dots \quad (3.15)$$

The **backward FD derivative** then writes:

$$\boxed{\frac{\partial f}{\partial x}(x_i) = \frac{f_i - f_{i-1}}{h} + \mathcal{O}(h)} \quad (\text{backward difference}) \quad (3.16)$$

Alternatively, we can subtract the backward formula from the forward one and divide by two. Concretely, we start from

$$f(x_0 + h) = f(x_0) + h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \dots \quad (3.17)$$

and subtract the following from it

$$f(x_0 - h) = f(x_0) - h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \dots \quad (3.18)$$

to obtain:

$$f(x_0 + h) - f(x_0 - h) = 2h \frac{\partial f}{\partial x}(x_0) + \mathcal{O}(h^3) \quad (3.19)$$

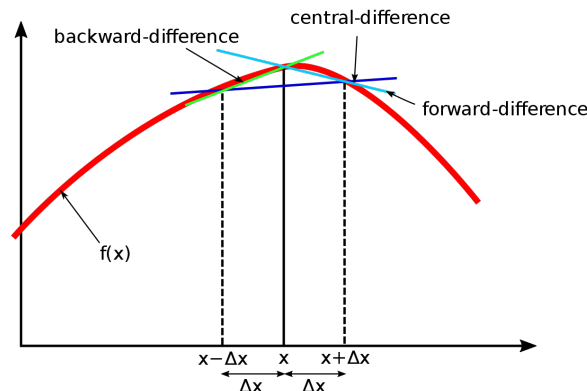
or,

$$\frac{\partial f}{\partial x}(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} + \mathcal{O}(h^2)$$

We see that the resulting **central difference** approximation is **second order accurate**. In the discrete world one then write

$$\boxed{\frac{\partial f}{\partial x}(x_i) = \frac{f_{i+1} - f_{i-1}}{2h} + \mathcal{O}(h^2)} \quad (\text{central difference}) \quad (3.20)$$

Simply put, the denominator is $2h$ because it is the distance between point x_{i-1} and x_{i+1} .



Can we do better than $\mathcal{O}(h^2)$? The answer is yes, and I list hereunder the formula:

$$f'_i = \frac{2f_{i+1} + 3f_i - 6f_{i-1} + f_{i-2}}{6h} + \mathcal{O}(h^3) \quad \text{backward difference} \quad (3.21)$$

$$f'_i = \frac{-f_{i+2} + 6f_{i+1} - 3f_i - 2f_{i-1}}{6h} + \mathcal{O}(h^3) \quad \text{forward difference} \quad (3.22)$$

$$f'_i = \frac{-f_{i+2} + 8f_{i+1} - 8f_{i-1} + f_{i-2}}{12h} + \mathcal{O}(h^4) \quad \text{central difference} \quad (3.23)$$

Looking at these formula it is obvious that the cost of forming the derivative is larger than before (more multiplications, additions, ...) which translates to longer calculations and, in the case of implicit methods, much denser matrices.



Exercise FDM-1 (bonus)

Prove the formula above.

Second-order derivatives

Many PDEs contain second order derivatives (typically diffusion equations) so we now turn to these and define $f''_i = f''(x_i) = \frac{\partial^2 f}{\partial x^2}(x_i)$.

Second order forward Let us define a function $g(x)$ such that $g = f'$. Then we have seen that the forward difference formula leads to write:

$$g'_i = \frac{g_{i+1} - g_i}{h} \quad (3.24)$$

On the one hand, we have $g'_i = g'(x_i) = f''(x_i) = f''_i$ and on the other hand

$$\frac{g_{i+1} - g_i}{h} = \frac{f'_{i+1} - f'_i}{h} \quad (3.25)$$

We can then use the forward derivative formula for f'_{i+1} and f'_i and obtain the following second order derivatives of f :

$$f''_i = \frac{f'_{i+1} - f'_i}{h} = \frac{\frac{f_{i+2} - f_{i+1}}{h} - \frac{f_{i+1} - f_i}{h}}{h} = \frac{f_{i+2} - 2f_{i+1} + f_i}{h^2} \quad (3.26)$$

which is the **first order accurate, forward difference** approximation for second order derivatives at x_i . In order to compute $f''(x_i)$ we need the value of f at x_i but also at two other locations right of this location.

Second order backward Likewise, we obtain the following formula when using the backward derivative twice:

$$f''_i = \frac{f'_i - f'_{i-1}}{h} = \frac{\frac{f_i - f_{i-1}}{h} - \frac{f_{i-1} - f_{i-2}}{h}}{h} = \frac{f_i - 2f_{i-1} + f_{i-2}}{h^2} \quad (3.27)$$

This time we need the value of f at x_i but also at two other locations left of this location.

³https://en.wikipedia.org/wiki/Finite_difference

Second order central By adding the Taylor expansions (with $+h$ and $-h$) a **second order accurate** approximation of the second derivative is obtained. We start from

$$\begin{aligned} f(x_0 + h) &= f(x_0) + h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \frac{h^3}{3!} \frac{\partial^3 f}{\partial x^3}(x_0) + \cdots + \frac{h^n}{n!} \frac{\partial^n f}{\partial x^n}(x_0) + \mathcal{O}(h^{n+1}) \\ f(x_0 - h) &= f(x_0) - h \frac{\partial f}{\partial x}(x_0) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) - \frac{h^3}{3!} \frac{\partial^3 f}{\partial x^3}(x_0) + \cdots + \frac{(-h)^n}{n!} \frac{\partial^n f}{\partial x^n}(x_0) + \mathcal{O}(h^{n+1}) \end{aligned}$$

and we see that adding the first equation to the second yields

$$f(x_0 + h) + f(x_0 - h) = 2f(x_0) + \underbrace{h \frac{\partial f}{\partial x}(x_0) - h \frac{\partial f}{\partial x}(x_0)}_{=0} + 2 \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(x_0) + \underbrace{\frac{h^3}{3!} \frac{\partial^3 f}{\partial x^3}(x_0) - \frac{h^3}{3!} \frac{\partial^3 f}{\partial x^3}(x_0)}_{=0} + \mathcal{O}(h^4) \quad (3.28)$$

or,

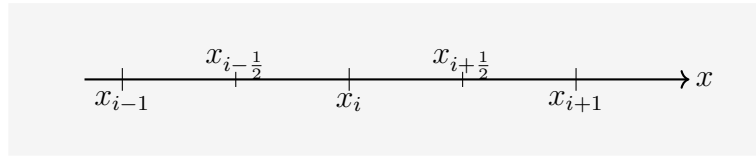
$$\frac{\partial^2 f}{\partial x^2}(x_0) = \frac{f(x_0 + h) - 2f(x_0) + f(x_0 - h)}{h^2} + \mathcal{O}(h^2) \quad (3.29)$$

which translates into

$$\boxed{f''_i = \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} + \mathcal{O}(h^2)} \quad (\text{second order central difference}) \quad (3.30)$$

Note that this formula requires one value left and one value right of the point under consideration.

Another way to arrive at the same expression is to write the expansion at $x_0 \pm h/2$, i.e. at the (convenient, yet nonexistent) half points $i \pm 1/2$:



$$f'_{i+1/2} = \frac{f_{i+1} - f_i}{h} \quad f'_{i-1/2} = \frac{f_i - f_{i-1}}{h} \quad (3.31)$$

$$f''_i = \frac{f'_{i+1/2} - f'_{i-1/2}}{h} = \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} \quad (3.32)$$

Note that derivatives of the form (see heat transport equation in Section 2.6)

$$\frac{\partial}{\partial x} \left(k \frac{\partial f}{\partial x} \right) \quad (3.33)$$

where k is a function of space, should be formed as follows

$$\left. \frac{\partial}{\partial x} \left(k \frac{\partial f}{\partial x} \right) \right|_i = \frac{k_{i+1/2} \frac{f_{i+1} - f_i}{h} - k_{i-1/2} \frac{f_i - f_{i-1}}{h}}{h} + \mathcal{O}(h^2) \quad (3.34)$$

where $k_{i \pm 1/2}$ is evaluated between the points to maintain the second order accuracy.

Remark. If the heat conductivity k shows strong jumps from one grid point to another that are not aligned with the grid-nodes, most second-order methods will show first order accuracy at best.

Can we do better than $\mathcal{O}(h^2)$? The answer is yes again:

$$f''_i = \frac{-f_{i+2} + 16f_{i+1} - 30f_i + 16f_{i-1} - f_{i-2}}{12h^2} + \mathcal{O}(h^4)$$

3.4 Solving the 1D diffusion equation

Consider the one-dimensional, transient (i.e. time-dependent) heat conduction equation without heat generating sources

$$\rho C_p \frac{\partial T}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) \quad (3.35)$$

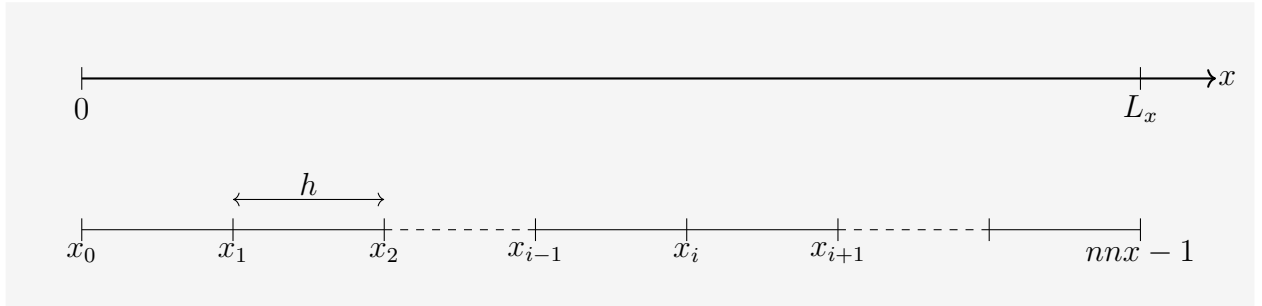
where ρ is density, C_p heat capacity, k thermal conductivity, T temperature, x distance, and t time.

If the thermal conductivity, density and heat capacity are constant over the model domain, the equation can be simplified to a diffusion equation:

$$\frac{\partial T}{\partial t} = \kappa \frac{\partial^2 T}{\partial x^2} \quad (3.36)$$

where $\kappa = k/\rho C_p$ is the heat diffusivity.

We wish to solve this PDE in time and space (provided the appropriate boundary conditions have been given). The domain is $[0, L_x]$ and it is discretised by means of nnx points as depicted hereunder:



The derivative of temperature with regards to time can be approximated with a forward finite difference approximation *in time* as

$$\frac{\partial T}{\partial t} \simeq \frac{T_i^{n+1} - T_i^n}{t^{n+1} - t^n} = \frac{T_i^{n+1} - T_i^n}{\delta t} \quad (3.37)$$

where δt is the time step, i.e. the time between two consecutive measurements (the equivalent of h in space). In all that follows the subscript will always refer to space indices while the superscript will always refer to time indices. To be clear: n represents the current time step whereas $n+1$ represents the next time step.

Both n and i are integers; n varies from 0 to $nstep - 1$ (total number of time steps) and i varies from 0 to $nnx - 1$ (where nnx is the total number of grid points in x -direction).

The spatial derivative is replaced by a central FD approximation

$$\frac{\partial^2 T}{\partial x^2} \simeq \frac{T_{i+1}^n - 2T_i^n + T_{i-1}^n}{h^2} \quad (3.38)$$

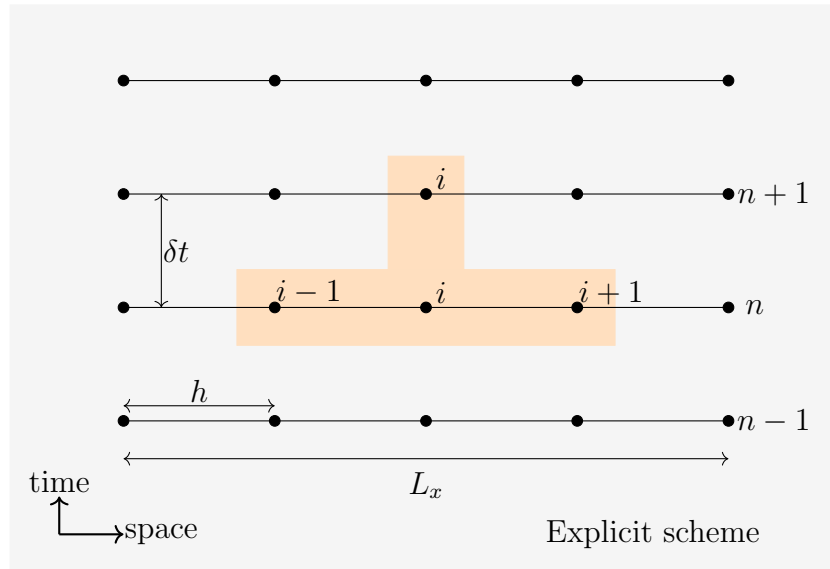
We obtain

$$\frac{T_i^{n+1} - T_i^n}{\delta t} = \kappa \frac{T_{i+1}^n - 2T_i^n + T_{i-1}^n}{h^2} \quad (3.39)$$

and finally

$$\boxed{T_i^{n+1} = T_i^n + \delta t \kappa \frac{T_{i+1}^n - 2T_i^n + T_{i-1}^n}{h^2}} \quad (3.40)$$

Because the temperature at the current time step n is known, we can compute the new temperature without solving any additional equations. Such a scheme is an **explicit** finite difference method and was made possible by the choice to evaluate the temporal derivative with forward differences.



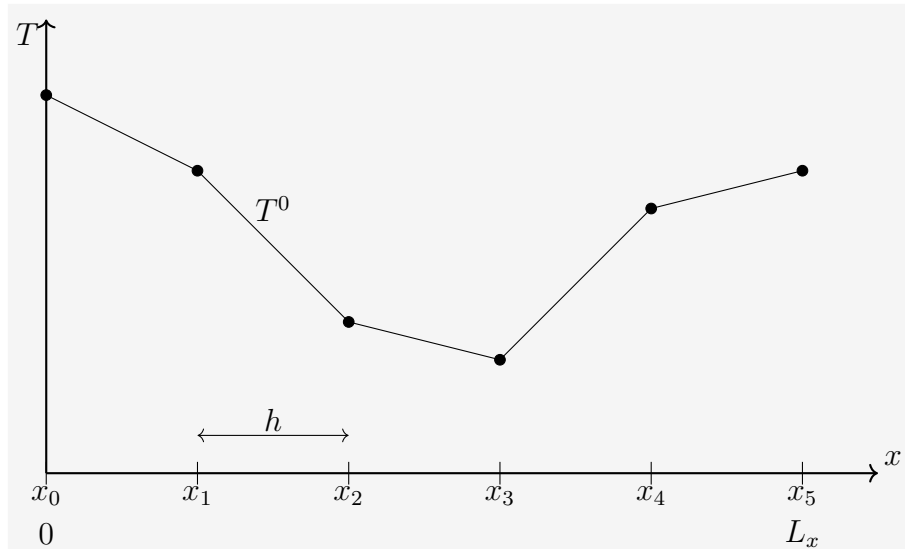
In order to solve the original PDE equation we need to

- prescribe an initial temperature field
- prescribe two boundary conditions

Such requirements hold also in the discrete world.

We know that this numerical scheme will converge to the exact solution for small h and δt because it has been shown to be **consistent** - that its discretization process can be reversed, through a Taylor series expansion, to recover the governing partial differential equation - and because it is **stable** for certain values of h and δt : any spontaneous perturbations in the solution (such as round-off error) will either be bounded or will decay.

Example FDM-2: let us prescribe an initial temperature field T_i^0 for $i = 0, \text{nnx} - 1$. For example:



Then, we will be able to compute the new temperature of (for example) node 3 at time $t = 1 \cdot \delta t$ (i.e. T_3^1) with

$$T_3^1 = T_3^0 + \delta t \kappa \frac{T_4^0 - 2T_3^0 + T_2^0}{h^2} \quad (3.41)$$

Note that T_0 and T_5 cannot be computed by means of the above equation, which is not a problem because both these values are actually the prescribed boundary conditions.

The main drawback of the explicit approach is that stable solutions are obtained *only* when

$$0 < \frac{2\kappa\delta t}{h^2} \leq 1 \quad \text{or,} \quad \delta t \leq \frac{h^2}{2\kappa} \quad (3.42)$$

If this condition is not satisfied, the solution becomes **unstable**, starts to wildly oscillate and ultimately 'blows up'. We will observe this during the practicals.

The stability condition means that the maximum time step needs to be smaller than the time it takes for an anomaly to diffuse across the grid (nodal) spacing h . The explicit solution is an example of a **conditionally stable method** that only leads to well behaved solutions if a criterion like the one above is satisfied.



Exercise FDM-2

We are going to solve the 1D diffusion equation with the explicit method for the following physical setup: The domain is $L_x = 1\text{km}$ long, it is maintained at a temperature $T = 100^\circ\text{C}$ at $x = 0$ and at a temperature $T = 200^\circ\text{C}$ at $x = L_x$. The initial temperature is $T(x, t = 0) = 123$ and $\kappa = 10^{-6}$. Time stepping will be carried out until steady state is reached.

→ `Exercise_2.FDM.ipynb`

An alternative approach is an **implicit** finite difference scheme, where the spatial derivatives of the Laplacian are evaluated (at least partially) at the new time step. We then use the backward difference for the time derivative:

$$\frac{\partial T}{\partial t} = \frac{T_i^n - T_i^{n-1}}{\delta t} \quad (3.43)$$

so that

$$\frac{T_i^n - T_i^{n-1}}{\delta t} = \kappa \frac{T_{i+1}^n - 2T_i^n + T_{i-1}^n}{h^2} \quad (3.44)$$

Note that this is often rewritten as follows in order to keep the unknowns at time $n + 1$:

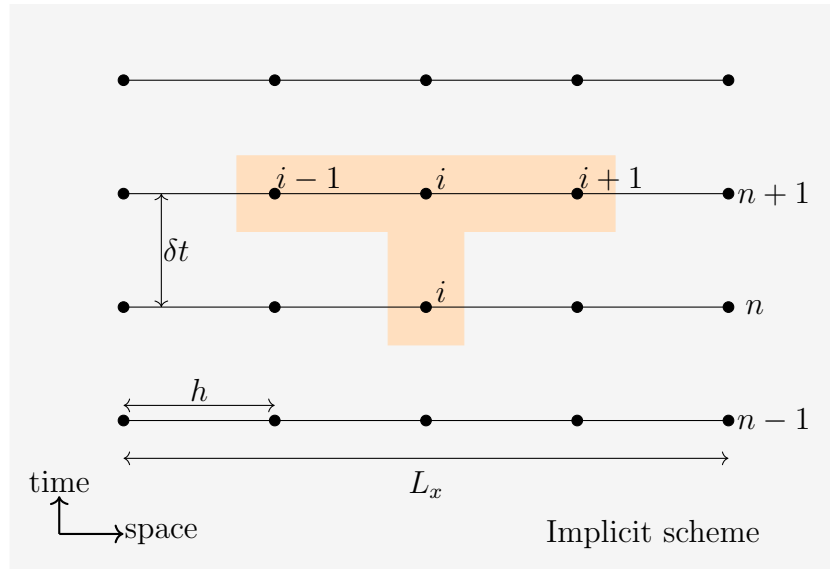
$$\frac{T_i^{n+1} - T_i^n}{\delta t} = \kappa \frac{T_{i+1}^{n+1} - 2T_i^{n+1} + T_{i-1}^{n+1}}{h^2} \quad (3.45)$$

It is a fully implicit scheme where the time derivative is taken backward. Let us define the dimensionless parameter s as follows:

$$s = \frac{\kappa \delta t}{h^2} \quad (3.46)$$

The previous equation can be rearranged as follows:

$$\boxed{-s T_{i+1}^{n+1} + (1 + 2s) T_i^{n+1} - s T_{i-1}^{n+1} = T_i^n} \quad (3.47)$$



Note that in this case we no longer have an explicit relationship for T_{i-1}^{n+1} , T_i^{n+1} and T_{i+1}^{n+1} . Instead, we have to solve a **linear system of equations**, which is discussed further below.

The main advantage of implicit methods is that there are no restrictions on the time step, the fully implicit scheme is **unconditionally stable**. This does not mean that it is accurate! Stability and accuracy are two different things! Taking large time steps may result in an inaccurate solution for features with small spatial scales!

For any application, it is therefore always a good idea to check the results by decreasing the time step until the solution does not change anymore (this is called a **convergence check**), and to ensure the method can deal with small and large scale features robustly at the same time.

Example FDM-3: Once again let us look at things with a very concrete approach. Let us discretise the domain of length L_x with 6 cells, i.e. $i = 0, \dots, 6$ ($nnx = 7$). We also prescribe the following boundary conditions (remember it is a 2nd order derivative in space, so we need two of them): $T(x = 0) = T_0 = 0$ and $T(x = L_x) = T_6 = 100$ (we assume that they do not change with time for simplicity). Finally we assume that we know T_i^0 for all i and we wish to compute T_i^1 .

We then have:

$$\begin{aligned}
T_0^1 &= 0 \\
-sT_2^1 + (1 + 2s)T_1^1 - sT_0^1 &= T_1^0 \\
-sT_3^1 + (1 + 2s)T_2^1 - sT_1^1 &= T_2^0 \\
-sT_4^1 + (1 + 2s)T_3^1 - sT_2^1 &= T_3^0 \\
-sT_5^1 + (1 + 2s)T_4^1 - sT_3^1 &= T_4^0 \\
-sT_6^1 + (1 + 2s)T_5^1 - sT_4^1 &= T_5^0 \\
T_6^1 &= 100
\end{aligned} \tag{3.48}$$

or,

$$\underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -s & 1+2s & -s & 0 & 0 & 0 & 0 \\ 0 & -s & 1+2s & -s & 0 & 0 & 0 \\ 0 & 0 & -s & 1+2s & -s & 0 & 0 \\ 0 & 0 & 0 & -s & 1+2s & -s & 0 \\ 0 & 0 & 0 & 0 & -s & 1+2s & -s \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{A}} \cdot \underbrace{\begin{pmatrix} T_0^1 \\ T_1^1 \\ T_2^1 \\ T_3^1 \\ T_4^1 \\ T_5^1 \\ T_6^1 \end{pmatrix}}_{\vec{T}} = \underbrace{\begin{pmatrix} 0 \\ T_1^0 \\ T_2^0 \\ T_3^0 \\ T_4^0 \\ T_5^0 \\ 100 \end{pmatrix}}_{\vec{b}}$$

As opposed to the explicit approach we must solve a linear system which size is given by the total number of nodes/points nnx in order to compute a new temperature field.

In summary, an implicit method requires us to solve $\mathbf{A} \cdot \vec{T} = \vec{b}$ with

- \mathbf{A} is a $nnx \times nnx$ **sparse** matrix (i.e. mostly empty),
- \vec{b} is a known vector of size nnx (often called the 'right-hand side', or **rhs**)
- \vec{T} the vector of unknowns.

A word about solvers There are two main approaches to solving such a linear system: one can use a **direct** approach or an **iterative** approach. In a nutshell, a direct solver will 'manipulate' the matrix lines and columns so as to arrive at the solution (like you would do on paper yourself for a small system). A simple example of such an approach is the technique of elimination of variables of the following example.

Example FDM-4: Consider the following system:

$$\begin{pmatrix} 1 & 3 & -2 \\ 3 & 5 & 6 \\ 2 & 4 & 3 \end{pmatrix} \cdot \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 5 \\ 7 \\ 8 \end{pmatrix} \quad (3.49)$$

which is of course equivalent to

$$\begin{aligned} x + 3y - 2z &= 5 \\ 3x + 5y + 6z &= 7 \\ 2x + 4y + 3z &= 8 \end{aligned} \quad (3.50)$$

Solving the first equation for x gives $x = 5 + 2z - 3y$, and plugging this into the second and third equation yields (or take the second line of the matrix and remove 3 times the first line from it, etc ...)

$$\begin{pmatrix} 1 & 3 & -2 \\ 0 & -4 & 12 \\ 0 & -2 & 7 \end{pmatrix} \cdot \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 5 \\ -8 \\ -2 \end{pmatrix} \quad (3.51)$$

Solving the second line for y yields $y = 2 + 3z$, and plugging this into the second equation yields $z = 2$. We now have:

$$\begin{pmatrix} 1 & 3 & -2 \\ 0 & -4 & 12 \\ 0 & 0 & 2 \end{pmatrix} \cdot \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 5 \\ -8 \\ 4 \end{pmatrix} \quad (3.52)$$

Substituting $z = 2$ into the second equation gives $y = 8$, and substituting $z = 2$ and $y = 8$ into the first equation yields $x = -15$. Therefore, the solution set is the single point $(x, y, z) = (-15, 8, 2)$.

Taken from https://en.wikipedia.org/wiki/System_of_linear_equations

This example is of course very naive and direct solvers often come in the form of very large numerical libraries which have been highly optimised to take advantage of the sparsity of the matrix in order to arrive at the solution in the lowest number of operations possible. In reality techniques such as LU decomposition or Cholesky decomposition are used.

Iterative solvers on the other hand compute the solution of the system by first postulating an initial guess for the solution and then by improving this guess iteratively until the **termination criterion** is met (see Section 9.33). There are two classes of iteratives methods in this context: **stationary iterative methods** (e.g. Jacobi, Gauss-Seidel, SSOR) and **Krylov subspace methods** (e.g. CG, GMRES, BiCG). At this stage things get real complicated and the details of iterative solvers are vastly out of the scope of this course⁴ (see for instance the book by Saad [1092]).

⁴https://en.wikipedia.org/wiki/Iterative_method

Example FDM-5: the stationary Jacobi method. The matrix \mathbf{A} is decomposed as follows:

$$\mathbf{A} = \mathbf{D} + \mathbf{L} + \mathbf{U} \quad (3.53)$$

where \mathbf{D} is the diagonal of the matrix \mathbf{A} , \mathbf{L} is the strict lower triangular part of \mathbf{A} and \mathbf{U} is the strict upper triangular part of \mathbf{A} . The iterative method is defined by:

$$\mathbf{D} \cdot \vec{T}^{k+1} = -(\mathbf{L} + \mathbf{U}) \cdot \vec{T}^k + \vec{b} \quad k = 0, 1, \dots \quad (3.54)$$

where \vec{T}^0 is the initial guess (often taken to be zero). Note that the superscript denotes the iteration number and has nothing to do with the time step in this context. This method is trivial to implement since the linear system on the left side of the equal sign involves a diagonal matrix. This can also be written

$$T_i^{k+1} = \frac{1}{A_{ii}} \left(b_i - \sum_{j \neq i} A_{ij} T_j^k \right) \quad i = 1, 2, \dots, n \quad (3.55)$$

Looking at the previous example, we have

$$\mathbf{D} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad \text{and} \quad \mathbf{L} + \mathbf{U} = \begin{pmatrix} 0 & 3 & -2 \\ 3 & 0 & 6 \\ 2 & 4 & 0 \end{pmatrix} \quad (3.56)$$

We then start with the guess $\vec{T}^0 = \vec{0}$, so that for $k = 0$:

$$\vec{T}^1 = \mathbf{D}^{-1} \cdot \vec{b} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1/5 & 0 \\ 0 & 0 & 1/3 \end{pmatrix} \cdot \begin{pmatrix} 5 \\ 7 \\ 8 \end{pmatrix} = \begin{pmatrix} 5 \\ 7/5 \\ 8/3 \end{pmatrix} \quad (3.57)$$

and then we obtain \vec{T}^2 by solving

$$\begin{aligned} \vec{T}^2 &= \mathbf{D}^{-1} \cdot \left[-(\mathbf{L} + \mathbf{U}) \cdot \vec{T}^1 + \vec{b} \right] \\ &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1/5 & 0 \\ 0 & 0 & 1/3 \end{pmatrix} \cdot \left[- \begin{pmatrix} 0 & 3 & -2 \\ 3 & 0 & 6 \\ 2 & 4 & 0 \end{pmatrix} \cdot \begin{pmatrix} 5 \\ 7/5 \\ 8/3 \end{pmatrix} + \begin{pmatrix} 5 \\ 7 \\ 8 \end{pmatrix} \right] \\ &= \dots \end{aligned} \quad (3.58)$$

We keep iterating until two consecutively obtained temperature vectors are nearly identical, or,

$$\|\vec{T}^{k+1} - \vec{T}^k\| < \epsilon \quad (3.59)$$

where ϵ is a carefully chosen small enough number.

A sufficient (but not necessary) condition for the method to converge is that the matrix \mathbf{A} is strictly or irreducibly diagonally dominant^a. Strict row diagonal dominance means that for each row, the absolute value of the diagonal term is greater than the sum of absolute values of other terms $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$. Also, this algorithm will fail if one or more diagonal terms of \mathbf{A} is nul.

^ahttps://en.wikipedia.org/wiki/Diagonally_dominant_matrix



Exercise FDM-3

Implement example FDM-5 from scratch in a new code. Write a separate function for the Jacobi solver. What do you observe ? why does it explode? Multiply all diagonal values by 10 and re-run it. What do you observe now? Bonus: implement the Gauss-Seidel method. Which of the two methods converges the fastest?

→https://en.wikipedia.org/wiki/Iterative_method.

Finally, looking at

$$-s T_{i+1}^{n+1} + (1 + 2s) T_i^{n+1} - s T_{i-1}^{n+1} = T_i^n \quad (3.60)$$

and dividing by $-s$ and letting $\delta t \rightarrow \infty$, we obtain:

$$T_{i+1}^{n+1} - 2T_i^{n+1} + T_{i-1}^{n+1} = 0 \quad (3.61)$$

which is a central difference approximation of the steady state solution

$$\frac{\partial^2 T}{\partial x^2} = 0 \quad (3.62)$$

Therefore, the fully implicit scheme will always yield the right equilibrium solution but may not capture small scale, transient features.



Exercise FDM-4

This is exactly the same exercise as Exercise 2 but we are going to solve the 1D diffusion equation with the implicit method this time. First use a solver from scipy to solve the system, then implement your own Jacobi solver. Note that the Jacobi solver must be implemented as a function which is to be called inside the time loop.

In your code use an if statement which allows to choose between explicit and implicit, and another if statement which allows to choose between scipy solver and iterative solver.

Bonus: implement the SSOR method in another function and compare Jacobi, Gauss-Seidel and SSOR.

→https://en.wikipedia.org/wiki/Iterative_method

Crank-Nicolson scheme It turns out that this fully implicit method is second order accurate in space but only first order accurate in time, i.e. the error goes as $\mathcal{O}(h^2, \delta t)$.

It is possible to write down a scheme which is second order accurate both in time and in space (i.e. $\mathcal{O}(h^2, \delta t^2)$), e.g. the **Crank-Nicolson**⁵ scheme which is unconditionally stable.

⁵The method was developed by John Crank and Phyllis Nicolson in the mid 20th century. https://en.wikipedia.org/wiki/Crank-Nicolson_method

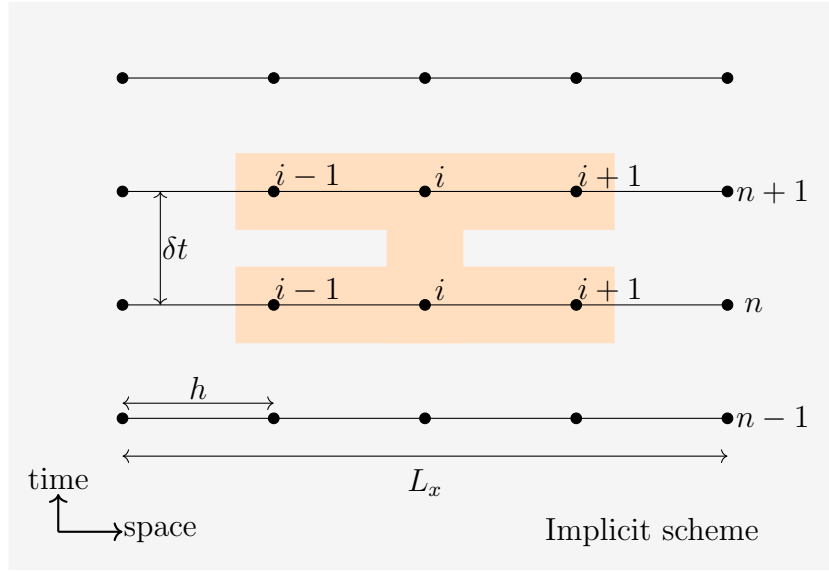
The Crank-Nicolson method is the time analog of central spatial differences and is given by

$$\frac{T_i^{n+1} - T_i^n}{\delta t} = \kappa \frac{1}{2} \left[\underbrace{\frac{T_{i+1}^n - 2T_i^n + T_{i-1}^n}{h^2}}_{\text{at time } n} + \underbrace{\frac{T_{i+1}^{n+1} - 2T_i^{n+1} + T_{i-1}^{n+1}}{h^2}}_{\text{at time } n+1} \right] \quad (3.63)$$

We define $s = \kappa \delta t / 2h^2$ so that the equation above can be rearranged as follows :

$$\boxed{-s T_{i+1}^{n+1} + (1 + 2s) T_i^{n+1} - s T_{i-1}^{n+1} = s T_{i+1}^n + (1 - 2s) T_i^n + s T_{i-1}^n} \quad (3.64)$$

Any partially implicit method is more complicated to compute as we need to infer the future solution at time $n + 1$ by solution (inversion) of a system of linear equations based on the known solution at time n .



Exercise FDM-5

Modify the code of Exercise FDM-4 to implement the Crank-Nicolson method.

3.5 Solving the 1D advection equation

fdm_adv1D.tex

The 1D hyperbolic advection equation is:

$$\rho C_p \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} \right) = 0 \quad (3.65)$$

or simply

$$\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} = 0 \quad (3.66)$$

We have seen how to deal with the time derivative (explicit, implicit) and with the first order space derivative (forward, backward or central). Let us consider the FTCS scheme (Forward in Time, Central in Space).

$$\frac{T_i^{n+1} - T_i^n}{\delta t} + u_i \frac{T_{i+1}^n - T_{i-1}^n}{2h} = 0$$

Note that although the velocity u is prescribed, it can vary in space, hence the subscript i .

There is however a major problem: the FTCS method is in this case **unconditionally unstable** (see Section 6.2.1 of [582], section 4.3.1 of [985]), i.e., it blows up for any δt . The instability is related to the fact that this scheme produces negative diffusion, which is numerically unstable. We could also consider the FTFS method:

$$\frac{T_i^{n+1} - T_i^n}{\delta t} + u_i \frac{T_{i+1}^n - T_i^n}{h} = 0$$

but it is also **unconditionally unstable** (see Section 6.2.1 of [582]).

We will now look at to methods which alleviate this problem:

- The **Lax-Friedrichs method**⁶ consists of replacing the T_i^n in the time derivative term with $(T_{i+1}^n + T_{i-1}^n)/2$ (see for instance Section 4.3.1 of [985] in the context of surface processes). The resulting equation is

$$\frac{T_i^{n+1} - (T_{i+1}^n + T_{i-1}^n)/2}{\delta t} = -u_i \frac{T_{i+1}^n - T_{i-1}^n}{2h}$$

or,

$$T_i^{n+1} = \frac{1}{2}(T_{i+1}^n + T_{i-1}^n) - \frac{u_i \delta t}{h} \frac{1}{2}(T_{i+1}^n - T_{i-1}^n)$$

von Neumann stability analysis indicates that this method is stable when $C = u\delta t/h \leq 1$ where C is the Courant number.

- In the **Streamline upwind** method the spatial finite difference scheme depends on the sign of the velocity:

$$\frac{T_i^{n+1} - (T_{i+1}^n + T_{i-1}^n)/2}{\delta t} = \begin{cases} -u_i \frac{T_i^n - T_{i-1}^n}{h_x} & \text{if } u_i < 0 \\ -u_i \frac{T_{i+1}^n - T_i^n}{h_x} & \text{if } u_i > 0 \end{cases}$$

In fact, we have replaced central with forward or backward derivatives, depending on the flow direction. This method is stable when $C = u\delta t/h \leq 1$.

These are not the only possibilities, see for instance the **leapfrog method** or the **Lax-Wendroff method** [582].

Finally, The Crank-Nicolson implicit scheme for solving the diffusion equation can be adapted to solve the advection equation:

$$T_i^{n+1} + \frac{u\delta t}{4h}(T_{i+1}^{n+1} - T_{i-1}^{n+1}) = T_i^n - \frac{u\delta t}{4h}(T_{i+1}^n - T_{i-1}^n)$$

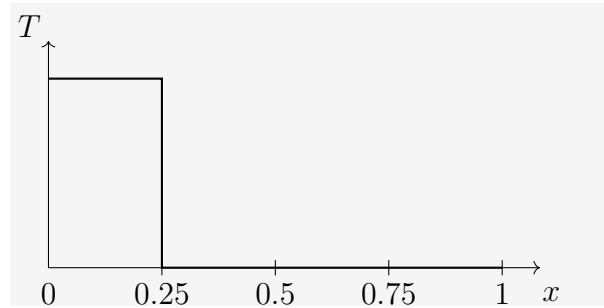
TODO: write about how we obtain this.

⁶Named after Peter Lax and Kurt O. Friedrichs



Exercise FDM-6

Let us consider the domain $[0, 1]$. The temperature field at $t = 0$ is given by $T = 1$ for $x < 0.25$ and $T = 0$ otherwise. The prescribed velocity is $u = 1$ and we set $nnx = 51$. Boundary conditions are $T = 1$ at $x = 0$ and $T = 0$ at $x = 1$.

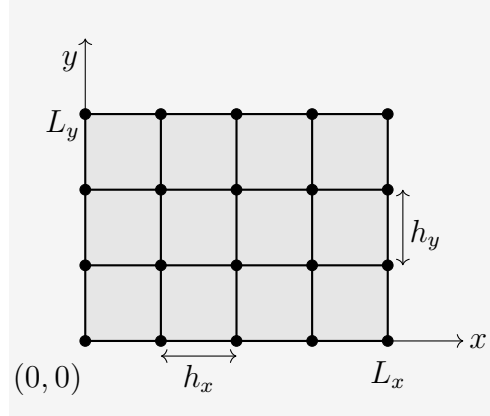


Program the above FTCS method. Run the model for 250 time steps with $\delta t = 0.002$. Program the Lax-Friedrichs method by modifying the previous code.

Bonus: Program the upwind method and/or the Crank-Nicolson method.

3.6 FDM basics in 2D

In a 2D Cartesian domain overlain by a $nnx \times nny$ grid, the spacing between nodes in the x and y direction is h_x and h_y respectively.



We have seen in Section 3.3 how to discretise second-order derivatives in 1D. In 2D, we then logically have for a function $f(x, y)$

$$\frac{\partial^2 f}{\partial x^2}(x_0, y_0) = \frac{f(x_0 + h_x, y_0) - 2f(x_0, y_0) + f(x_0 - h_x, y_0)}{h_x^2} + \mathcal{O}(h_x^2) \quad (3.67)$$

$$\frac{\partial^2 f}{\partial y^2}(x_0, y_0) = \frac{f(x_0, y_0 + h_y) - 2f(x_0, y_0) + f(x_0, y_0 - h_y)}{h_y^2} + \mathcal{O}(h_y^2) \quad (3.68)$$

What about mixed derivatives? Since these are combinations of first-order derivatives, we can straightforwardly discretise them:

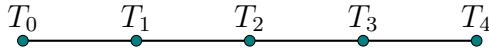
$$\begin{aligned} & \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0) \\ &= \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial y} \right) (x_0, y_0) \\ &= \frac{\partial}{\partial x} \left(\frac{f(x_0, y_0 + h_y) - f(x_0, y_0 - h_y)}{2h_y} \right) \\ &= \frac{1}{2h_y} \frac{\partial f}{\partial x}(x_0, y_0 + h_y) - \frac{1}{2h_y} \frac{\partial f}{\partial x}(x_0, y_0 - h_y) \\ &= \frac{1}{2h_y} \frac{f(x_0 + h_x, y_0 + h_y) - f(x_0 - h_x, y_0 + h_y)}{2h_x} - \frac{1}{2h_y} \frac{f(x_0 + h_x, y_0 - h_y) - f(x_0 - h_x, y_0 - h_y)}{2h_x} \\ &= \frac{f(x_0 + h_x, y_0 + h_y) - f(x_0 - h_x, y_0 + h_y) - f(x_0 + h_x, y_0 - h_y) + f(x_0 - h_x, y_0 - h_y)}{2h_x h_y} + \mathcal{O}(h_x^2, h_y^2) \end{aligned}$$

From 1D to 2D

INSERT TEXT

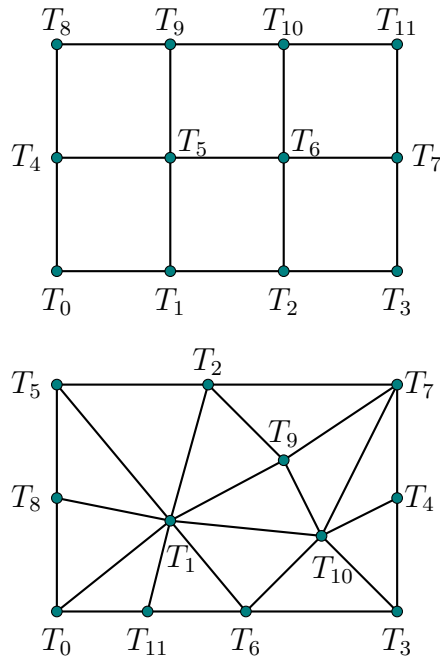
(tikz_needicon.tex)

1D



$$\vec{T} = \begin{pmatrix} T_0 \\ T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix}$$

2D



$$\vec{T} = \begin{pmatrix} T_0 \\ T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \\ T_6 \\ T_7 \\ T_8 \\ T_9 \\ T_{10} \\ T_{11} \end{pmatrix}$$

Also, here is a rather handy code snippet which should allow you to make nice plots of the coming exercises.

```
filename = 'solution_{:04d}.pdf'.format(istep)
fig = plt.figure ()
#ax = fig.gca(projection='3d')
ax = fig.add_subplot(projection='3d')
ax.plot_surface(x.reshape ((nny,nnx)),y.reshape ((nny,nnx)),T.reshape ((nny,nnx)
),color = 'darkseagreen')
ax.set_xlabel ( 'X[_m_] ' )
ax.set_ylabel ( 'Y[_m_] ' )
ax.set_zlabel ( 'Temperature[_C_] ' )
plt.title('Timestep_{:.2d}' %(istep),loc='right')
plt.grid ()
plt.savefig(filename)
#plt.show ()
plt.close()
```

3.7 Solving the 2D diffusion equation

We now revisit the transient heat equation, this time with sources/sinks for 2D problems. In the absence of advective heat transport, the heat equation is

$$\rho C_p \frac{\partial T}{\partial t} = \vec{\nabla} \cdot k \vec{\nabla} T + Q \quad (3.69)$$

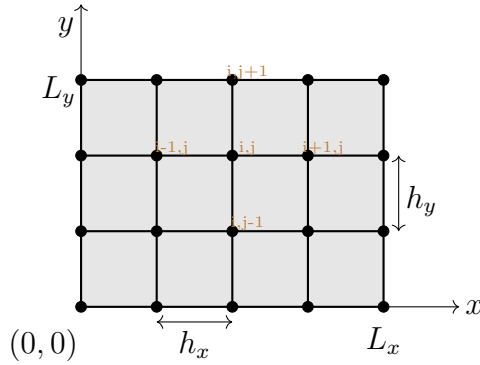
where Q is the radiogenic heat production. It simply writes as follows when Cartesian coordinates are used:

$$\rho C_p \frac{\partial T}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) + \frac{\partial}{\partial y} \left(k \frac{\partial T}{\partial y} \right) + Q \quad (3.70)$$

If the heat conductivity is constant in space (and so are the other coefficients), it writes:

$$\frac{\partial T}{\partial t} = \kappa \left(\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} \right) + \tilde{Q} \quad (3.71)$$

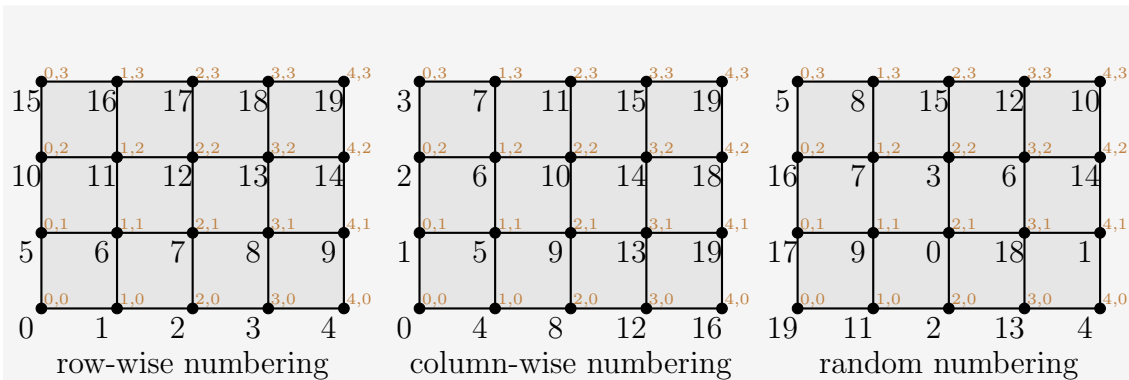
with $\tilde{Q} = Q/\rho C_p$. In order to solve this equation over the Cartesian domain of size $L_x \times L_y$ we need to generate a mesh as shown hereunder:



The spacing between the nodes in the x -direction is h_x and h_y is the spacing between the nodes in the y direction. There are now $nnp = nnx \times nny$ nodes in total. The above grid is characterised by $i = 0, 1, 2, 3, 4$ and $j = 0, 1, 2, 3$ and counts in total 20 nodes.

In one dimension, the subscript indicated the node i . In two dimensions we therefore need two indices i and j to identify a node, so that the temperature at node i, j at time n is denoted $T_{i,j}^n$.

One question remains: should we number nodes row by row ? column by column ? randomly ? These three approaches are shown hereunder:



This is a critical point because the discretised PDE is formulated as a function of $T_{i,j}$ with $i = 0, \dots, nnx - 1$ and $j = 0, \dots, nny - 1$ but the vector \vec{T} containing all these values (encountered in implicit methods) is indexed by a single index $k = 0, \dots, nnp - 1$. The numbering strategy determines

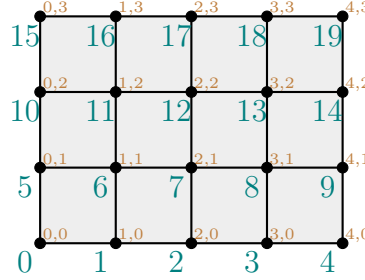
how easy it is to go from (i, j) to k and vice versa. Very concretely again, where should $T_{3,4}$ be placed in the global vector of unknowns \vec{T} ?

At the same time we cannot do away with i, j indices because these are needed to locate the direct neighbours of any node and allow to form discrete derivatives.

We then need a (preferably simple/straightforward) 'function' which associates to every (i, j) a global index k . For the first grid with row-wise numbering, we have $0 \leq i \leq 4$, $0 \leq j \leq 3$ and $0 \leq k \leq 19$. It follows that

$$k(i, j) = j \cdot nnx + i \quad (3.72)$$

This is easy to verify: $i = 3$ and $j = 2$ indeed corresponds to node # 13, $i = 4$ and $j = 1$ corresponds to node # 9, etc ...



Exercise FDM-7

In a new code declare and assign values to nnx and nnx . Compute nnp . Set $L_x = 7$ and $L_y = 6$. Compute h_x and h_y .

Declare two arrays $xcoords$ and $ycoords$ which will contain the x and y coordinates of all nnp nodes.

By means of two imbricated for loops compute these coordinates & fill both arrays.

Visualise the nodes with matplotlib.

Tip: Make sure your code works for various combinations of nnx and nnx .

3.7.1 Explicit scheme

The simplest approach is an FTCS (forward time, centered space) explicit method like in 1D:

$$\frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} = \kappa \left(\frac{T_{i-1,j}^n - 2T_{i,j}^n + T_{i+1,j}^n}{h_x^2} + \frac{T_{i,j-1}^n - 2T_{i,j}^n + T_{i,j+1}^n}{h_y^2} \right) + \tilde{Q}_{i,j}^n \quad (3.73)$$

where we have assumed that the source term \tilde{Q} can depend of space coordinates and therefore appears as $\tilde{Q}_{i,j}$ in the equation. We define s_x and s_y as follows:

$$s_x = \frac{\kappa \delta t}{h_x^2} \quad s_y = \frac{\kappa \delta t}{h_y^2} \quad (3.74)$$

so that

$$T_{i,j}^{n+1} = T_{i,j}^n + s_x(T_{i-1,j}^n - 2T_{i,j}^n + T_{i+1,j}^n) + s_y(T_{i,j-1}^n - 2T_{i,j}^n + T_{i,j+1}^n) + \tilde{Q}_{i,j}^n \delta t \quad (3.75)$$

or,

$$T_{k(i,j)}^{n+1} = T_{k(i,j)}^n + s_x(T_{k(i-1,j)}^n - 2T_{k(i,j)}^n + T_{k(i+1,j)}^n) + s_y(T_{k(i,j-1)}^n - 2T_{k(i,j)}^n + T_{k(i,j+1)}^n) + \tilde{Q}_{k(i,j)}^n \delta t \quad (3.76)$$

The scheme is stable for

$$\delta t \leq \frac{\min(h_x^2, h_y^2)}{2\kappa} \quad (3.77)$$

Boundary conditions can be set the usual way: for example a constant (Dirichlet) temperature at node (i, j) is given by

$$T_{i,j} = T_{bc} \quad (3.78)$$

where T_{bc} is the prescribed temperature.

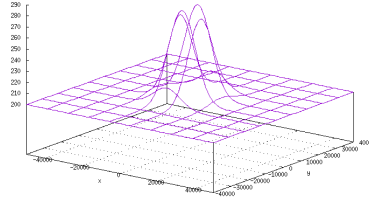


Exercise FDM-8

A simple (time-dependent) analytical solution for the temperature equation exists for the case that the initial temperature field is

$$T(x, y, t = 0) = T_0 + T_{max} \exp \left[-\frac{x^2 + y^2}{\sigma^2} \right] \quad (3.79)$$

where T_{max} is the maximum amplitude of the temperature perturbation at $(x, y) = (0, 0)$ and σ its half-width.



initial temperature field

The solution of the time-dependent PDE is

$$T(x, y, t) = T_0 + \frac{T_{max}}{1 + 4t\kappa/\sigma^2} \exp \left[-\frac{x^2 + y^2}{\sigma^2 + 4t\kappa} \right] \quad (3.80)$$

Set $L_x=100\text{km}$ and $L_y = 80\text{km}$, $\kappa = 10^{-6}$, $\tilde{Q} = 0$, $T_{max} = 100^\circ$, $T_0 = 200^\circ$, and $\sigma = 10^4\text{m}$.

Use the previous exercise to generate a $nnx \times nny$ grid in the $[-L_x/2, L_x/2] \times [-L_y/2, L_y/2]$ domain.

Write a function which takes x , y , t , T_0 , T_{max} , κ and σ as argument and returns the analytical temperature value.

Write a an explicit FDM code which solves the 2D diffusion equation. At each time step prescribe on the boundary the analytical solution.

3.7.2 Implicit scheme

If we now employ a fully implicit, unconditionally stable discretization scheme, the discretised PDE becomes:

$$\frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} = \kappa \left(\frac{T_{i-1,j}^{n+1} - 2T_{i,j}^{n+1} + T_{i+1,j}^{n+1}}{h_x^2} + \frac{T_{i,j-1}^{n+1} - 2T_{i,j}^{n+1} + T_{i,j+1}^{n+1}}{h_y^2} \right) + \frac{Q_{i,j}^n}{\rho C_p} \quad (3.81)$$

Rearranging terms with $n + 1$ on the left and terms with n on the right hand side gives

$$-s_x T_{i+1,j}^{n+1} - s_y T_{i,j+1}^{n+1} + (1 + 2s_x + 2s_y) T_{i,j}^{n+1} - s_x T_{i-1,j}^{n+1} - s_y T_{i,j-1}^{n+1} = T_{i,j}^n + \tilde{Q}_{i,j}^n \delta t \quad (3.82)$$

or

$$-s_x T_{k(i+1,j)}^{n+1} - s_y T_{k(i,j+1)}^{n+1} + (1 + 2s_x + 2s_y) T_{k(i,j)}^{n+1} - s_x T_{k(i-1,j)}^{n+1} - s_y T_{k(i,j-1)}^{n+1} = T_{k(i,j)}^n + \tilde{Q}_{k(i,j)}^n \delta t \quad (3.83)$$

which here again yields a linear system of equations written $\mathbf{A} \cdot \vec{T} = \vec{b}$ where \mathbf{A} is a $(nnp \times nnp)$ matrix.

Boundary conditions are $T(x, y) = 0$ on all sides, so all nodes on the boundary have a prescribed zero temperature⁷:

$$\begin{aligned} T_{0,0} = T_0 &= 0 \\ T_{1,0} = T_1 &= 0 \\ T_{2,0} = T_2 &= 0 \\ T_{3,0} = T_3 &= 0 \\ T_{4,0} = T_4 &= 0 \\ T_{0,1} = T_5 &= 0 \\ T_{4,1} = T_9 &= 0 \\ T_{0,2} = T_{10} &= 0 \\ T_{4,2} = T_{14} &= 0 \\ T_{0,3} = T_{15} &= 0 \\ T_{1,3} = T_{16} &= 0 \\ T_{2,3} = T_{17} &= 0 \\ T_{3,3} = T_{18} &= 0 \\ T_{4,3} = T_{19} &= 0 \end{aligned}$$

In what follows we assume for simplicity and conciseness of notation that $h_x = h_y = h$ so that $s_x = s_y = s$. The discretised PDE equation will now be applied to the interior nodes:

- For node $k = 6$ ($i = 1, j = 1$):

$$\begin{aligned} &-sT_{2,1}^{n+1} - sT_{1,2}^{n+1} + (1 + 4s)T_{1,1}^{n+1} - sT_{0,1}^{n+1} - sT_{1,0}^{n+1} = T_{1,1}^n + \tilde{Q}_{1,1}^n \delta t \\ \Rightarrow &-sT_7^{n+1} - sT_{11}^{n+1} + (1 + 4s)T_6^{n+1} - sT_5^{n+1} - sT_1^{n+1} = T_6^n + \tilde{Q}_6^n \delta t \end{aligned} \quad (3.84)$$

- For node $k = 7$ ($i = 2, j = 1$):

$$\begin{aligned} &-sT_{3,1}^{n+1} - sT_{2,2}^{n+1} + (1 + 4s)T_{2,1}^{n+1} - sT_{1,1}^{n+1} - sT_{2,0}^{n+1} = T_{2,1}^n + \tilde{Q}_{2,1}^n \delta t \\ \Rightarrow &-sT_8^{n+1} - sT_{12}^{n+1} + (1 + 4s)T_7^{n+1} - sT_6^{n+1} - sT_2^{n+1} = T_7^n + \tilde{Q}_7^n \delta t \end{aligned} \quad (3.85)$$

- For node $k = 8$ ($i = 3, j = 1$):

$$\begin{aligned} &-sT_{4,1}^{n+1} - sT_{3,2}^{n+1} + (1 + 4s)T_{3,1}^{n+1} - sT_{2,1}^{n+1} - sT_{3,0}^{n+1} = T_{3,1}^n + \tilde{Q}_{3,1}^n \delta t \\ \Rightarrow &-sT_9^{n+1} - sT_{13}^{n+1} + (1 + 4s)T_8^{n+1} - sT_7^{n+1} - sT_3^{n+1} = T_8^n + \tilde{Q}_8^n \delta t \end{aligned} \quad (3.86)$$

⁷We assume here again that these boundary conditions do not change with time.

- For node $k = 11$ ($i = 1, j = 2$):

$$\begin{aligned} & -sT_{2,2}^{n+1} - sT_{1,3}^{n+1} + (1 + 4s)T_{1,2}^{n+1} - sT_{0,2}^{n+1} - sT_{2,1}^{n+1} = T_{1,2}^n + \tilde{Q}_{1,2}^n \delta t \\ \Rightarrow & -sT_{12}^{n+1} - sT_{16}^{n+1} + (1 + 4s)T_{11}^{n+1} - sT_{10}^{n+1} - sT_6^{n+1} = T_{11}^n + \tilde{Q}_{11}^n \delta t \end{aligned} \quad (3.87)$$

- For node $k = 12$ ($i = 2, j = 2$):

$$\begin{aligned} & -sT_{3,2}^{n+1} - sT_{2,3}^{n+1} + (1 + 4s)T_{2,2}^{n+1} - sT_{1,2}^{n+1} - sT_{2,1}^{n+1} = T_{2,2}^n + \tilde{Q}_{2,2}^n \delta t \\ \Rightarrow & -sT_{13}^{n+1} - sT_{17}^{n+1} + (1 + 4s)T_{12}^{n+1} - sT_{11}^{n+1} - sT_7^{n+1} = T_{12}^n + \tilde{Q}_{12}^n \delta t \end{aligned} \quad (3.88)$$

- For node $k = 13$ ($i = 3, j = 2$):

$$\begin{aligned} & -sT_{4,2}^{n+1} - sT_{3,3}^{n+1} + (1 + 4s)T_{3,2}^{n+1} - sT_{2,2}^{n+1} - sT_{3,1}^{n+1} = T_{3,2}^n + \tilde{Q}_{3,2}^n \delta t \\ \Rightarrow & -sT_{14}^{n+1} - sT_{18}^{n+1} + (1 + 4s)T_{13}^{n+1} - sT_{12}^{n+1} - sT_8^{n+1} = T_{13}^n + \tilde{Q}_{13}^n \delta t \end{aligned} \quad (3.89)$$

Putting it all together yields the following linear system:

$$\underbrace{\begin{pmatrix}
1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & 1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & . & 1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & . & . & 1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & . & . & . & 1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & . & . & . & . & 1 & . & . & . & . & . & . & . & . & . & . & . & . & . & . \\
. & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . & . & . & . & . & . & . & . \\
. & . & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . & . & . & . & . & . & . \\
. & . & . & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . & . & . & . & . & . \\
. & . & . & . & . & . & . & . & . & 1 & . & . & . & . & . & . & . & . & . & . \\
. & . & . & . & . & . & . & . & . & . & 1 & . & . & . & . & . & . & . & . & . \\
. & . & . & . & . & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . & . & . & . \\
. & . & . & . & . & . & -s & . & . & -s & 1+4s & -s & . & . & . & -s & . & . & . & . \\
. & . & . & . & . & . & . & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . & . \\
. & . & . & . & . & . & . & . & -s & . & . & . & -s & 1+4s & -s & . & . & . & -s & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & 1 & . & . & . & . & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & . & 1 & . & . & . & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & 1 & . & . & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & 1 & . & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & 1 & . \\
. & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & . & 1
\end{pmatrix}}_{\mathbf{A}} \cdot \underbrace{\begin{pmatrix} T_0^{n+1} \\ T_1^{n+1} \\ T_2^{n+1} \\ T_3^{n+1} \\ T_4^{n+1} \\ T_5^{n+1} \\ T_6^{n+1} \\ T_7^{n+1} \\ T_8^{n+1} \\ T_9^{n+1} \\ T_{10}^{n+1} \\ T_{11}^{n+1} \\ T_{12}^{n+1} \\ T_{13}^{n+1} \\ T_{14}^{n+1} \\ T_{15}^{n+1} \\ T_{16}^{n+1} \\ T_{17}^{n+1} \\ T_{18}^{n+1} \\ T_{19}^{n+1} \end{pmatrix}}_{\vec{T}} = \underbrace{\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ T_6^n + \tilde{Q}_6 \delta t \\ T_7^n + \tilde{Q}_7 \delta t \\ T_8^n + \tilde{Q}_8 \delta t \\ 0 \\ 0 \\ T_{11}^n + \tilde{Q}_{11} \delta t \\ T_{12}^n + \tilde{Q}_{12} \delta t \\ T_{13}^n + \tilde{Q}_{13} \delta t \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}}_{\vec{b}}$$

Note that we now have five 'diagonals' filled with non-zero entries as opposed to three diagonals in the 1D case.

replace zeros by some more random value corresponding to Tboundary

Note also that this is a simplified matrix since we assumed that $s_x = s_y$.



Exercise FDM-9

Same exercise as exercise FDM-8, but now with implicit method.

Looking at this matrix, it is clear that this approach is sub-optimal: for such a small grid counting 20 nodes, the boundary conditions enforce the temperature on 14 of them, so that these temperatures should/could be removed from the list of unknowns, leaving a vector of unknowns \vec{T} of size 6 (the number of nodes which are not on the boundary). As a consequence, we would have to solve a 6×6 linear system, as opposed to a 20×20 one!

In this case, we focus again on nodes 6,7,8,11,12,13. we start from

$$-sT_7^{n+1} - sT_{11}^{n+1} + (1 + 4s)T_6^{n+1} - sT_5^{n+1} - sT_1^{n+1} = T_6^n + \tilde{Q}_6^n \delta t \quad (3.90)$$

but we know that the boundary conditions impose that $T_1 = 0$ and $T_5 = 0$ so that the equation above simplifies to:

$$-sT_7^{n+1} - sT_{11}^{n+1} + (1 + 4s)T_6^{n+1} = T_6^n + \tilde{Q}_6^n \delta t \quad (3.91)$$

These 6 equations can finally be combined in the expected smaller linear system:

$$\underbrace{\begin{pmatrix} 1 + 4s & -s & . & -s & . & . \\ -s & 1 + 4s & -s & . & -s & . \\ . & -s & 1 + 4s & . & . & -s \\ -s & . & -s & 1 + 4s & -s & . \\ . & -s & . & -s & 1 + 4s & -s \\ . & . & -s & . & -s & 1 + 4s \end{pmatrix}}_{\mathbf{A}} \cdot \underbrace{\begin{pmatrix} T_6^{n+1} \\ T_7^{n+1} \\ T_8^{n+1} \\ T_{11}^{n+1} \\ T_{12}^{n+1} \\ T_{13}^{n+1} \end{pmatrix}}_{\vec{T}} = \underbrace{\begin{pmatrix} T_6^n + \tilde{Q}_6^n \delta t \\ T_7^n + \tilde{Q}_7^n \delta t \\ T_8^n + \tilde{Q}_8^n \delta t \\ T_{11}^n + \tilde{Q}_{11}^n \delta t \\ T_{12}^n + \tilde{Q}_{12}^n \delta t \\ T_{13}^n + \tilde{Q}_{13}^n \delta t \end{pmatrix}}_{\vec{b}} \quad (3.92)$$

Note that is the boundary values had not been zero they would have found their way to the right hand side vector.

The Crank-Nicolson version of the implicit scheme is then as follows:

$$\begin{aligned} \frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} &= \frac{1}{2}\kappa \left(\frac{T_{i-1,j}^{n+1} - 2T_{i,j}^{n+1} + T_{i+1,j}^{n+1}}{h_x^2} + \frac{T_{i,j-1}^{n+1} - 2T_{i,j}^{n+1} + T_{i,j+1}^{n+1}}{h_y^2} \right) \\ &+ \frac{1}{2}\kappa \left(\frac{T_{i-1,j}^n - 2T_{i,j}^n + T_{i+1,j}^n}{h_x^2} + \frac{T_{i,j-1}^n - 2T_{i,j}^n + T_{i,j+1}^n}{h_y^2} \right) \end{aligned} \quad (3.93)$$

The implementation of this method will require from you to bring all the terms in T^{n+1} to the left of the equal sign while all the terms in T^n are assumed to be known and therefore find their way into the right hand side.

Likewise, the Lax-Friedrichs method is as follows:

$$\frac{T_{i,j}^{n+1} - \frac{1}{4}(T_{i-1,j}^n + T_{i+1,j}^n + T_{i,j-1}^n + T_{i,j+1}^n)}{\delta t} = \kappa \left(\frac{T_{i-1,j}^{n+1} - 2T_{i,j}^{n+1} + T_{i+1,j}^{n+1}}{h_x^2} + \frac{T_{i,j-1}^{n+1} - 2T_{i,j}^{n+1} + T_{i,j+1}^{n+1}}{h_y^2} \right) + \frac{Q_{i,j}^n}{\rho C_p} \quad (3.94)$$

Rearranging terms with $n + 1$ on the left and terms with n on the right hand side gives

3.7.3 The 9-point stencil for the Laplace operator

What follows is mostly borrowed from Wikipedia⁸.

If we discretize the 2D Laplacian by using central-difference methods, we obtain the commonly used five-point stencil, represented by the following convolution kernel:

$$D = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

or,

$$DT_{i,j} = \frac{1}{h^2}(T_{i-1,j} + T_{i+1,j} + T_{i,j-1} + T_{i,j+1} - 4T_{i,j})$$

Even though it is simple to obtain and computationally lighter, the central difference kernel possess an undesired intrinsic anisotropic property, since it doesn't take into account the diagonal neighbours.

The two most commonly used isotropic nine-point stencils are displayed below, in their convolution kernel forms. They can be obtained by the following formula

$$D = (1 - \gamma) \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix} + \gamma \begin{bmatrix} 1/2 & 0 & 1/2 \\ 0 & -2 & 0 \\ 1/2 & 0 & 1/2 \end{bmatrix}$$

The first one is known by Oono-Puri, and it is obtained when $\gamma = 1/2$:

$$D = \begin{bmatrix} 1/4 & 2/4 & 1/4 \\ 2/4 & -12/4 & 2/4 \\ 1/4 & 2/4 & 1/4 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 2 & -12 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

The second one is known by Patra-Karttunen or Mehrstellen, and it is obtained when $\gamma = 1/3$:

$$D = \begin{bmatrix} 1/6 & 4/6 & 1/6 \\ 4/6 & -20/6 & 4/6 \\ 1/6 & 4/6 & 1/6 \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 1 & 4 & 1 \\ 4 & -20 & 4 \\ 1 & 4 & 1 \end{bmatrix}$$

or,

$$\vec{\nabla}^2 T_{i,j} = \frac{1}{6h^2}(T_{i+1,j+1} + T_{i-1,j+1} + T_{i+1,j-1} + T_{i-1,j-1} + 4(T_{i+1,j} + T_{i-1,j}) + 4(T_{i,j+1} + T_{i,j-1}) - 20T_{i,j})$$

Both are isotropic forms of discrete Laplacian, and in the limit of small h , they all become equivalent. This form is the one we find in LeVeque [778, p64].

3.8 Solving the 2D advection-diffusion equation

So far, we have mainly focused on the diffusion equation in a non-moving flow (relevant for the case of a dike intrusion cooling off or for a lithosphere which remains undeformed).

We now want to consider problems where material moves during the time period under consideration and takes temperature anomalies with it (e.g. a plume rising through a convecting mantle). If the numerical grid remains fixed in the background, the hot temperatures should be moved to different grid points at each time step.

We start again from the heat transport equation of Section 2.6:

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T + Q \quad (3.95)$$

⁸https://en.wikipedia.org/wiki/Nine-point_stencil

We have previously dealt with the one-dimensional Cartesian coordinates equation:

$$\rho C_p \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} \right) = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) + Q \quad (3.96)$$

and we now turn to the two-dimensional equation:

$$\rho C_p \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} \right) = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) + \frac{\partial}{\partial y} \left(k \frac{\partial T}{\partial y} \right) + Q \quad (3.97)$$

As before, assuming that k is constant in space we can rewrite the equation as a function of the heat diffusivity κ :

$$\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} = \kappa \left(\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} \right) + Q \quad (3.98)$$

Since we have already seen how to deal with 'pure' diffusion equations in the previous section, let us now turn to 'pure' advection equations:

$$\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} = 0 \quad (3.99)$$

or

$$\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} = 0 \quad (3.100)$$

where we assume $\vec{v} = (u, v)$ known.

Even though the equations appear simple, it is quite tricky to solve them accurately, more so than for the diffusion problem. This is particularly the case if there are large gradients in the quantity that is to be advected.

We have seen how to deal with the time derivative (explicit, implicit) and with the first order space derivative (forward, backward or central). Let us consider again the FTCS scheme (Forward in Time, Central in Space).

$$\frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} + u_{i,j} \frac{T_{i+1,j}^n - T_{i-1,j}^n}{2h_x} + v_{i,j} \frac{T_{i,j+1}^n - T_{i,j-1}^n}{2h_y} = 0 \quad (3.101)$$

The fully implicit version is then as follows:

$$\frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} + u_{i,j} \frac{T_{i+1,j}^{n+1} - T_{i-1,j}^{n+1}}{2h_x} + v_{i,j} \frac{T_{i,j+1}^{n+1} - T_{i,j-1}^{n+1}}{2h_y} = 0 \quad (3.102)$$

or,

$$T_{i,j}^{n+1} + \frac{u_{i,j} \delta t}{2h_x} (T_{i+1,j}^{n+1} - T_{i-1,j}^{n+1}) + \frac{v_{i,j} \delta t}{2h_y} (T_{i,j+1}^{n+1} - T_{i,j-1}^{n+1}) = T_{i,j}^n$$

The terms on the left will form five diagonals in the matrix while the term on the right is the right hand side.

The Crank-Nicolson approach is then easily derived by taking:

$$\frac{T_{i,j}^{n+1} - T_{i,j}^n}{\delta t} + \frac{u_{i,j}}{2} \frac{T_{i+1,j}^n - T_{i-1,j}^n}{2h_x} + \frac{v_{i,j}}{2} \frac{T_{i,j+1}^n - T_{i,j-1}^n}{2h_y} + \frac{u_{i,j}}{2} \frac{T_{i+1,j}^{n+1} - T_{i-1,j}^{n+1}}{2h_x} + \frac{v_{i,j}}{2} \frac{T_{i,j+1}^{n+1} - T_{i,j-1}^{n+1}}{2h_y} = 0 \quad (3.103)$$



Exercise FDM-10

We wish to compute the advection of a product-cosine hill in a prescribed velocity field. The initial temperature is:

$$T_0(x, y) = \begin{cases} \frac{1}{4} \left(1 + \cos \pi \frac{x-x_c}{\sigma}\right) \left(1 + \cos \pi \frac{y-y_c}{\sigma}\right) & \text{if } (x - x_c)^2 + (y - y_c)^2 \leq \sigma^2 \\ 0 & \text{otherwise} \end{cases}$$

The boundary conditions are $T(x, y) = 0$ on all four sides of the unit square domain. In what follows we set $x_c = y_c = 2/3$ and $\sigma = 0.2$. The velocity field is analytically prescribed: $\vec{v} = (-(y - L_y/2), +(x - L_x/2))$. Resolution is set to 31×31 nodes.

The timestep is set to $\delta t = 2\pi/200$ and we wish to carry out 200 timesteps so that the cone does a 2π rotation.

See Stone 43 for results/figures of this experiment obtained with Finite Elements.

Implement this with the FTCS method. What do you observe? What happens when you decrease the value of δt ?

Bonus: Lax method, Crank-Nicolson method.



Exercise FDM-11

NOT FOR 2020!

Redo exercise FDM-6 in a unit square domain. The temperature field at $t = 0$ is given by $T(x, y) = 1$ for $x < 0.25$ and $T(x, y) = 0$ otherwise. The prescribed velocity is $\vec{v} = (1, 0)$ and we set $n_{nx} = n_{ny} = 51$. Boundary conditions are $T = 1$ at $x = 0$ and $T = 0$ at $x = 1$.

Program the above FTCS method. Run the model for 250 time steps with $\delta t = 0.002$. Compare the 2D solution with the previously obtained 1D solution of exercise FDM-6.

Make sure the code works in the y -direction too by rotating the initial temperature by 90° anti-clockwise, set $\vec{v} = (0, 1)$ and change boundary conditions accordingly.

Bonus: Lax method, Crank-Nicolson method.

Note to self: - CFL missing still - stencils next to boundaries missing too - add many more visual aids

3.9 FEM vs FDM?

Let us start with the 1D steady advection-diffusion equation:

$$\rho C_p u \frac{dT}{dx} - k \frac{d^2 T}{dx^2} = f \quad \text{in } [0, L_x] \quad (3.104)$$

with the boundary conditions $T(x=0) = 0$ and $T(x=L_x) = 0$.

We have seen before (see Section ??) that the elemental matrix \mathbf{K}_a for the advection and the elemental matrix \mathbf{K}_d for the diffusion terms are

$$\mathbf{K}_a^e = \frac{\rho C_p u}{2} \begin{pmatrix} -1 & 1 \\ -1 & 1 \end{pmatrix} \quad \mathbf{K}_d^e = \frac{k}{h_x} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

where h_x is the distance between nodes in the x -direction and e denotes the element number.

Assuming that we have 5 elements (i.e. 6 nodes), the assembled 6×6 advection and diffusion matrices (before boundary conditions are applied) are:

$$\mathbf{K}_a = \frac{\rho C_p u}{2} \begin{pmatrix} -1 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \\ 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix} \quad \mathbf{K}_d = \frac{k}{h_x} \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix}$$

The rhs is zero, so that we would have to solve $(\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = 0$, or:

$$\left[\frac{\rho C_p u}{2} \begin{pmatrix} -1 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \\ 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix} + \frac{k}{h_x} \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix} \right] \cdot \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \\ T_6 \end{pmatrix} = \vec{0}$$

Note that boundary conditions are not applied yet. Therefore the algebraic equation for an interior node i is

$$\rho C_p u \frac{T_{i+1} - T_{i-1}}{2} + \frac{k}{h_x} (-T_{i-1} + 2T_i - T_{i+1}) = 0$$

or,

$$\boxed{\rho C_p u \frac{T_{i+1} - T_{i-1}}{2h_x} - \frac{k}{h_x^2} (T_{i-1} - 2T_i + T_{i+1}) = 0} \quad (3.105)$$

However, we have seen in Section ?? that the second order accurate central differencing based approximate first and second derivatives written for an interior node i of a finite difference mesh with a constant node spacing of h is

$$\left. \frac{dT}{dx} \right|_i \simeq \frac{T_{i+1} - T_{i-1}}{2h_x} \quad \frac{d^2 T}{dx^2} \simeq \frac{T_{i+1} - 2T_i + T_{i-1}}{h_x^2}$$

Using these approximations, the discretised formulation of Eq. (3.104) is exactly the same as Eq. (3.105). This simple example proves that the FEM and the FDM share similarities!

It is also useful to introduce the elemental Peclet number

$$Pe = \frac{uh}{2\kappa} = \frac{uh\rho C_p}{2k}$$

and Eq. (3.104) becomes:

$$\frac{u}{2h_x} \left[\left(1 - \frac{1}{Pe} \right) T_{i+1} + \frac{2}{Pe} T_i - \left(1 + \frac{1}{Pe} \right) T_{i-1} \right] = f$$

CHECK!!!

Chapter 4

Numerical integration

As we will see later, using the Finite Element method to solve problems involves computing integrals which are more often than not too complex to be computed analytically/exactly. We will then need to compute them numerically.

[wiki] In essence, the basic problem in numerical integration is to compute an approximate solution to a definite integral

$$I = \int_a^b f(x)dx \quad (4.1)$$

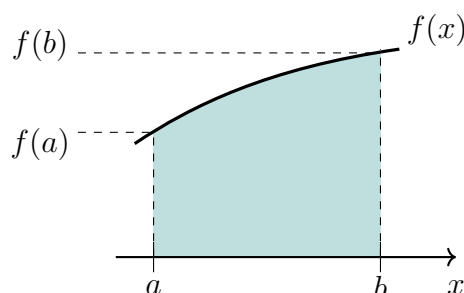
to a given degree of accuracy. This problem has been widely studied and we know that if $f(x)$ is a smooth function, and the domain of integration is bounded, there are many methods for approximating the integral to the desired precision.

There are several reasons for carrying out numerical integration.

- The integrand $f(x)$ may be known only at certain points, such as obtained by sampling. Some embedded systems and other computer applications may need numerical integration for this reason.
- A formula for the integrand may be known, but it may be difficult or impossible to find an antiderivative that is an elementary function. An example of such an integrand is $f(x) = \exp(-x^2)$, the antiderivative of which (the error function, times a constant) cannot be written in elementary form.
- It may be possible to find an antiderivative symbolically, but it may be easier to compute a numerical approximation than to compute the antiderivative. That may be the case if the antiderivative is given as an infinite series or product, or if its evaluation requires a special function that is not available.

Let us remember that the integral of Eq. (4.1) is in fact equal to the (signed) area between the x -axis and the curve $f(x)$ over the interval $[a, b]$:

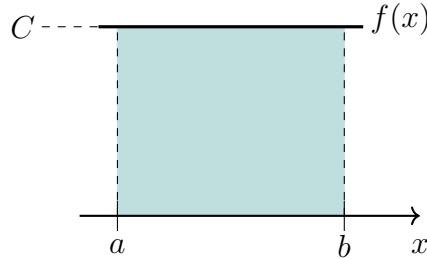
(tikz-quadrature_idf.tex)



Note that in the example above $f(x) > 0$ so the area of the gray domain is counted positive. For example, if the function $f(x)$ is a polynomial the integral can easily be computed analytically. In the case of a 0^{th} order polynomial, we have $f(x) = C$ where C is a constant. We then have

$$I = \int_a^b f(x)dx = \int_a^b C dx = C \int_a^b dx = C(b-a) \quad (4.2)$$

(tikz-quadrature_idf2.tex)



We see that the area of the gray domain is simply the product of its length $b-a$ by its height C and we indeed recover $I = C(b-a)$.

4.1 In 1 dimension

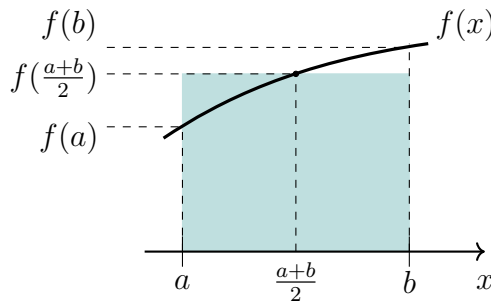
4.1.1 Midpoint and Trepezoidal rules

The simplest method of this type is to let the interpolating function be a constant function (a polynomial of degree zero) that passes through the point $((a+b)/2, f((a+b)/2))$. This is called the midpoint rule or rectangle rule. We then have

$$I = \int_a^b f(x)dx \simeq (b-a)f\left(\frac{a+b}{2}\right)$$

which is the area of this gray domain:

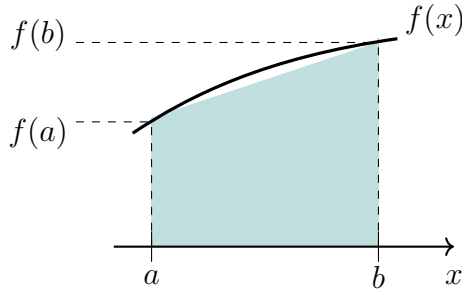
(tikz-quadrature_rectangle.tex)



We can do a little bit better at virtually no cost: we choose the interpolating function to be a straight line (an affine function, i.e. a polynomial of degree 1) passing through the points $(a, f(a))$ and $(b, f(b))$. This is called the trapezoidal rule. Then

$$I = \int_a^b f(x)dx \simeq (b-a)\frac{f(a)+f(b)}{2}$$

(tikz-quadrature_trapeze.tex)

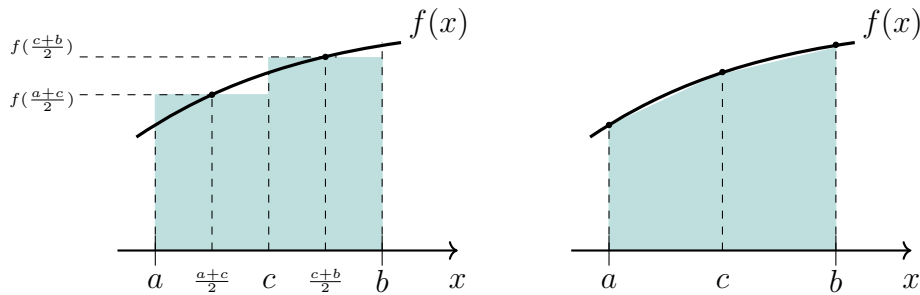


We see that if the function f is monotonous on the interval $[a, b]$ then the trapezoidal approach is likely to return a value close to the real value. However if the function f oscillates a lot in the interval, approximating it with a single rectangle or trapeze is not a sound assumption. We can then make use of the additive property of the integral: let c be the coordinate of the middle of the $[a, b]$ interval, i.e. $c = (a + b)/2$. Then we have

$$I = \int_a^b f(x)dx = \int_a^c f(x)dx + \int_c^b f(x)dx$$

We can then apply the midpoint rule or the trapezoidal rule over both segments $[a, c]$ and $[c, b]$:

(tikz_quadrature_both.tex)



In this case we would have

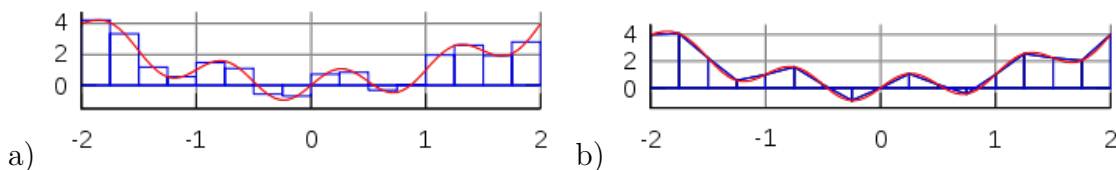
$$I_{\text{midpoint}} = (c-a)f\left(\frac{a+c}{2}\right) + (b-c)f\left(\frac{c+b}{2}\right)$$

$$I_{\text{trapeze}} = (c-a)\frac{f(a)+f(c)}{2} + (b-c)\frac{f(c)+f(b)}{2}$$

Of course we can repeat the process and for either one of these rules, we can make a more accurate approximation by breaking up the interval $[a, b]$ into some number n of subintervals, computing an approximation for each subinterval, then adding up all the results. For example, the composite trapezoidal rule can be stated as

$$\int_a^b f(x)dx \simeq \frac{b-a}{n} \left(\frac{f(a)}{2} + \sum_{k=1}^{n-1} f\left(a + k\frac{b-a}{n}\right) + \frac{f(b)}{2} \right) \quad (4.3)$$

where the subintervals have the form $[kh, (k+1)h]$, with $h = (b-a)/n$ and $k = 0, 1, 2, \dots, n-1$.



The interval $[-2, 2]$ is broken into 16 sub-intervals. The blue lines correspond to the approximation of the red curve by means of a) the midpoint rule, b) the trapezoidal rule.

There are several algorithms for numerical integration (also commonly called “**numerical quadrature**”, or simply “**quadrature**”) . Interpolation with polynomials evaluated at equally spaced points in $[a, b]$ yields the Newton-Cotes formulas, of which the rectangle rule and the trapezoidal rule are examples.

4.1.2 in 1D - Gauss-Legendre quadrature

If we allow the intervals between interpolation points to vary, we find another group of quadrature formulas, such as the Gauss(ian) quadrature formulas. A Gaussian quadrature rule is typically more accurate than a Newton-Cotes rule, which requires the same number of function evaluations, if the integrand is smooth (i.e., if it is sufficiently differentiable).

An n -point Gaussian quadrature rule, named after Carl Friedrich Gauss, is a quadrature rule constructed to yield an exact result for polynomials of degree $2n - 1$ or less by a suitable choice of the points x_i and weights w_i for $i = 1, \dots, n$.

The domain of integration for such a rule is conventionally taken as $[-1, 1]$, so the rule is stated as

$$\int_{-1}^{+1} f(x)dx = \sum_{i_q=1}^n w_{i_q} f(x_{i_q})$$

In this formula the x_{i_q} coordinate is the i -th root of the **Legendre polynomial**¹ $P_n(x)$.

It is important to note that a Gaussian quadrature will only produce good results if the function $f(x)$ is well approximated by a polynomial function within the range $[-1, 1]$. As a consequence, the method is not, for example, suitable for functions with singularities.

¹https://en.wikipedia.org/wiki/Legendre_polynomials

| n | x_{iq} | w_{iq} | x_{iq} (approx) | w_{iq} (approx) |
|----|-------------------------------------------------|-------------------------------|---------------------------------|-----------------------|
| 1 | 0 | 2 | 0 | 2 |
| 2 | $\pm\sqrt{1/3}$ | 1 | $\pm 0.577\ 350\ 269\ 189\ 626$ | 1 |
| 3 | 0 | 8/9 | 0 | 0.888 888 888 888 888 |
| | $\pm\sqrt{3/5}$ | 5/9 | $\pm 0.774\ 596\ 669\ 241\ 483$ | 0.555 555 555 555 555 |
| 4 | $\pm\sqrt{\frac{3}{7} - \frac{2}{7}\sqrt{6/5}}$ | $\frac{18+\sqrt{30}}{36}$ | $\pm 0.339\ 981\ 043\ 584\ 856$ | 0.652 145 154 862 546 |
| | $\pm\sqrt{\frac{3}{7} + \frac{2}{7}\sqrt{6/5}}$ | $\frac{18-\sqrt{30}}{36}$ | $\pm 0.861\ 136\ 311\ 594\ 953$ | 0.347 854 845 137 454 |
| 5 | 0 | 128/225 | 0 | 0.568 888 888 888 889 |
| | $\pm\frac{1}{3}\sqrt{5 - 2\sqrt{\frac{10}{7}}}$ | $\frac{322+13\sqrt{70}}{900}$ | $\pm 0.538\ 469\ 310\ 105\ 683$ | 0.478 628 670 499 366 |
| | $\pm\frac{1}{3}\sqrt{5 + 2\sqrt{\frac{10}{7}}}$ | $\frac{322-13\sqrt{70}}{900}$ | $\pm 0.906\ 179\ 845\ 938\ 664$ | 0.236 926 885 056 189 |
| 6 | ? | ? | $\pm 0.238\ 619\ 186\ 083\ 197$ | 0.467 913 934 572 691 |
| | | | $\pm 0.661\ 209\ 386\ 466\ 265$ | 0.360 761 573 048 139 |
| | | | $\pm 0.932\ 469\ 514\ 203\ 152$ | 0.171 324 492 379 170 |
| | | | $\pm 0.949\ 107\ 912\ 342\ 759$ | 0.129 484 966 168 870 |
| 7 | | | $\pm 0.741\ 531\ 185\ 599\ 394$ | 0.279 705 391 489 277 |
| | | | $\pm 0.405\ 845\ 151\ 377\ 397$ | 0.381 830 050 505 119 |
| | | | 0.000 000 000 000 000 | 0.417 959 183 673 469 |
| | | | $\pm 0.960\ 289\ 856\ 497\ 536$ | 0.101 228 536 290 376 |
| 8 | | | $\pm 0.796\ 666\ 477\ 413\ 627$ | 0.222 381 034 453 374 |
| | | | $\pm 0.525\ 532\ 409\ 916\ 329$ | 0.313 706 645 877 887 |
| | | | $\pm 0.183\ 434\ 642\ 495\ 650$ | 0.362 683 783 378 362 |
| | | | $\pm 0.968\ 160\ 239\ 507\ 626$ | 0.081 274 388 361 574 |
| 9 | | | $\pm 0.836\ 031\ 107\ 326\ 636$ | 0.180 648 160 694 857 |
| | | | $\pm 0.613\ 371\ 432\ 700\ 590$ | 0.260 610 696 402 935 |
| | | | $\pm 0.324\ 253\ 423\ 403\ 809$ | 0.312 347 077 040 003 |
| | | | 0.000 000 000 000 000 | 0.330 239 355 001 260 |
| 10 | | | $\pm 0.973\ 906\ 528\ 517\ 172$ | 0.066 671 344 308 688 |
| | | | $\pm 0.865\ 063\ 366\ 688\ 985$ | 0.149 451 349 150 581 |
| | | | $\pm 0.679\ 409\ 568\ 299\ 024$ | 0.219 086 362 515 982 |
| | | | $\pm 0.433\ 395\ 394\ 129\ 247$ | 0.269 266 719 309 996 |
| | | | $\pm 0.148\ 874\ 338\ 981\ 631$ | 0.295 524 224 714 753 |

Abcissae and weights for Gauss quadratures up to $n = 10$. See [779, p89]. Also check <https://pomax.github.io/bezierinfo/legendre-gauss.html>.

As shown in the above table, it can be shown that the weight values must fulfil the following condition:

$$\sum_{i_q} w_{i_q} = 2 \quad (4.4)$$

This simply comes from the requirement that when $f(x) = 1$ then $\int_{-1}^{+1} f(x)dx = 2 = \sum w_{i_q}$. It is also worth noting that all quadrature point coordinates are symmetrical around the origin.

Since most quadrature formula are only valid on a specific interval, we now must address the problem of their use outside of such intervals. The solution turns out to be quite simple: one must carry out a change of variables from the interval $[a, b]$ to $[-1, 1]$. We then consider the reduced coordinate $r \in [-1, 1]$ such that

$$r = \frac{2}{b-a}(x-a) - 1 \quad (4.5)$$

This relationship can be reversed such that when r is known, its equivalent coordinate $x \in [a, b]$ can

be computed:

$$x = \frac{b-a}{2}(1+r) + a \quad (4.6)$$

From this it follows that

$$dx = \frac{b-a}{2}dr \quad (4.7)$$

and then

$$\int_a^b f(x)dx = \frac{b-a}{2} \int_{-1}^{+1} f(r)dr \simeq \frac{b-a}{2} \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) \quad (4.8)$$

4.1.3 A probably naive way of finding the quadrature points coordinates and weights

We start from the assumption that the quadrature must be exact for polynomials $f(r)$, that it is written

$$I = \int_{-1}^{+1} f(r)dr = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q})$$

and that $n_q > 0$, $w_{i_q} \neq 0$ and $r_{i_q} \in [-1, 1]$.

Let us start with zero-th order polynomials, i.e. $f(r) = C$. Then $I = 2C$ and we must then have

$$2C = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} C$$

which imposes

$$\sum_{i_q=1}^{n_q} w_{i_q} = 2 \quad \forall n_q > 0 \quad (4.9)$$

As long as the sum of the weights is equal to 2, any n_q -point based quadrature can integrate exactly a zero-th order polynomial.

Let us move on with first-order polynomials. Since we have covered the constant term hereabove, we set $f(r) = ar$ where $a \neq 0$. We have $I = 0$ so

$$0 = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = a \sum_{i_q=1}^{n_q} w_{i_q} r_{i_q} \quad \Rightarrow \quad \sum_{i_q=1}^{n_q} w_{i_q} r_{i_q} = 0 \quad \forall n_q > 0 \quad (4.10)$$

In order to integrate exactly first-order polynomials an n_q -point based quadrature must fulfil Eqs.(4.9) and (4.10).

- If $n_q = 1$, then we automatically have $w_1 = 2$ and $w_1 r_1 = 0$, i.e. $r_1 = 0$.
- If $n_q = 2$, then $w_1 + w_2 = 2$ and $w_1 r_1 + w_2 r_2 = 0$. There are many solutions w_1, w_2, r_1, r_2 which can fulfil these two equations, so this is not enough to determine a unique set of coordinates and weights.

Let us now turn to second-order polynomials: as before, I choose $f(r) = ar^2$. We have $I = 2a/3$ and

$$\frac{2a}{3} = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = a \sum_{i_q=1}^{n_q} w_{i_q} r_{i_q}^2 \quad \sum_{i_q=1}^{n_q} w_{i_q} r_{i_q}^2 = \frac{2}{3} \quad \forall n_q > 0 \quad (4.11)$$

- If $n_q = 1$, we know that $w_1 = 2$ and $r_1 = 0$. This means that 1-point quadrature cannot exactly integrate polynomials higher than 1.
- If $n_q = 2$, then $w_1 + w_2 = 2$, $w_1 r_1 + w_2 r_2 = 0$ and now $w_1 r_1^2 + w_2 r_2^2 = 2/3$. We have three equations but still four unknowns. At this stage, we can do a simple additional assumption: common sense would have us realise that there is no reason why the (in this case) 2 quadrature point coordinates should be both negative or both positive. In light thereof we require that quadrature point coordinates are symmetric with respect with the origin $r = 0$, i.e. $r_1 = -r_2$ in this case. This yields to write: $w_1 r_1 + w_2 r_2 = w_1 r_1 + w_2 (-r_1) = r_1 (w_1 - w_2) = 0$. If $r_1 = 0$ then $r_2 = 0$ too and we do not have a 2-point quadrature. It must then follows that $w_1 = w_2$. And finally $w_1 r_1^2 + w_2 r_2^2 = w_1 r_1^2 + w_1 (-r_1)^2 = 2/3$, i.e. $r_1 = -1/\sqrt{3}$ and $r_2 = 1/\sqrt{3}$ since $r_1 < r_2$.

If we now turn to third-order polynomials, i.e. $f(r) = ar^3$, then $I = 0$ again. We then must have

$$\boxed{\sum_{i_q=1}^{n_q} w_{i_q} r_{i_q}^3 = 0 \quad \forall n_q > 0} \quad (4.12)$$

We see that the coordinates and weights obtained for a 2-point quadrature verify this equation, i.e. a 2-point quadrature can also exactly integrate a 3rd-order polynomial. However, it is equally easy to verify that the 2-point quadrature cannot exactly integrate a 4th-order polynomial since

$$I = \int_{-1}^{+1} r^4 dr = \frac{2}{5} \neq \sum_{i_q=1}^2 w_{i_q} r_{i_q}^4$$

A three-point quadrature will then be needed for those. Because of the symmetry, we know that the middle point will be at $r = 0$.

Remark. *This approach unfortunately does not shed any light on why the method is called Gauss-Legendre quadrature nor why the quadrature points are the zeros of the Legendre polynomials...*

4.1.4 Examples

Example 1

Since we know how to carry out any required change of variables, we choose for simplicity $a = -1$, $b = +1$. Let us take for example $f(r) = \pi$. Then we can compute the integral of this function over the interval $[a, b]$ exactly:

$$I = \int_{-1}^{+1} f(r) dr = \pi \int_{-1}^{+1} dr = 2\pi$$

We can now use a Gauss-Legendre formula to compute this same integral:

$$I_{gq} = \int_{-1}^{+1} f(r) dr = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} \pi = \pi \underbrace{\sum_{i_q=1}^{n_q} w_{i_q}}_{=2} = 2\pi$$

where we have used the property of the weight values of Eq.(4.4). Since the actual number of points was never specified, this result is valid for all quadrature rules.

Example 2

Let us now take $f(r) = mr + p$ and repeat the same exercise:

$$I = \int_{-1}^{+1} f(r)dr = \int_{-1}^{+1} (mr + p)dr = \left[\frac{1}{2}mr^2 + pr\right]_{-1}^{+1} = 2p$$

$$I_{gq} = \int_{-1}^{+1} f(r)dr = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} (mr_{i_q} + p) = m \underbrace{\sum_{i_q=1}^{n_q} w_{i_q} r_{i_q}}_{=0} + p \underbrace{\sum_{i_q=1}^{n_q} w_{i_q}}_{=2} = 2p$$

since the quadrature points are symmetric w.r.t. to zero on the r -axis. Once again the quadrature is able to compute the exact value of this integral: this makes sense since an n -point rule exactly integrates a $2n - 1$ order polynomial such that a 1 point quadrature exactly integrates a first order polynomial like the one above.

Example 3

Let us now take $f(r) = r^2$. We have

$$I = \int_{-1}^{+1} f(r)dr = \int_{-1}^{+1} r^2 dr = \left[\frac{1}{3}r^3\right]_{-1}^{+1} = \frac{2}{3}$$

and

$$I_{gq} = \int_{-1}^{+1} f(r)dr = \sum_{i_q=1}^{n_q} w_{i_q} f(r_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} r_{i_q}^2$$

- $n_q = 1$: $r_{i_q}^{(1)} = 0$, $w_{i_q} = 2$. $I_{gq} = 0$
- $n_q = 2$: $r_q^{(1)} = -1/\sqrt{3}$, $r_q^{(2)} = 1/\sqrt{3}$, $w_q^{(1)} = w_q^{(2)} = 1$. $I_{gq} = \frac{2}{3}$
- It also works $\forall n_q > 2$!

4.2 In 2 & 3 dimensions

4.2.1 On the reference square

Let us now turn to a two-dimensional integral of the form

$$I = \int_{-1}^{+1} \int_{-1}^{+1} f(r, s) dr ds$$

where $f(r, s)$ is again assumed to be continuous over the domain. The equivalent Gaussian quadrature writes:

$$I_{gq} \simeq \sum_{i_q=1}^{n_q} \sum_{j_q=1}^{n_q} f(r_{i_q}, s_{j_q}) w_{i_q} w_{j_q}$$

Finally we have

$$I = \int_a^{+b} \int_c^{+d} f(r, s) dr ds \simeq \frac{b-a}{2} \frac{d-c}{2} \sum_{i_q=1}^{n_q} \sum_{j_q=1}^{n_q} f(r_{i_q}, s_{j_q}) w_{i_q} w_{j_q} \quad (4.13)$$

4.2.2 On a generic quadrilateral

Let K be a quadrilateral element with straight boundary lines and with vertices arranged as follows:

IMAGE

We wish to evaluate

$$I = \iint_K f(x, y) dx dy$$

In order to do so we will first transform the element K to the reference square element and then apply the quadrature of the previous section. This transformation can be carried out by means of the Q_1 basis functions, see Section 5.3.1. We construct a linear mapping between the quadrilateral element K and the reference square element:

$$x(r, s) = \sum_{i=1}^4 \mathcal{N}_i(r, s) x_i \quad (4.14)$$

$$y(r, s) = \sum_{i=1}^4 \mathcal{N}_i(r, s) y_i \quad (4.15)$$

Then we have

$$I = \iint_K f(x, y) dx dy = \int_{-1}^{+1} \int_{-1}^{+1} f(x(r, s), y(r, s)) |\mathbf{J}(r, s)| dr ds$$

where $\mathbf{J}(r, s)$ is the Jacobian of the transformation defined by

$$\mathbf{J}(r, s) = \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial y}{\partial r} \\ \frac{\partial x}{\partial s} & \frac{\partial y}{\partial s} \end{pmatrix}$$

Finally applying the Gaussian quadrature yields:

$$I = \iint_K f(x, y) dx dy \simeq \sum_{i_q=1}^{n_q} \sum_{j_q=1}^{n_q} f(x(r_{i_q}, s_{j_q}), y(r_{i_q}, s_{j_q})) |\mathbf{J}(r_{i_q}, s_{j_q})| w_{i_q} w_{j_q}$$

4.2.3 Exercises



Exercise Quad-1

Write a program which uses the midpoint rule to compute (subdivide the interval in n subintervals)

$$I = \int_0^{\pi/2} f(x) dx \quad f(x) = x \quad \text{and} \quad f(x) = \cos(x)$$

Compute and plot the (absolute) error between the measured I_n and the analytical value I as a function of the subinterval size h .

Bonus: same as before with $I = \int_1^3 \int_2^4 (x^2 y^3 + xy + 1) dx dy$.



Exercise Quad-2

Same exercise as above but with the trapezoidal rule. Which method is the most accurate?



Exercise Quad-3

The following Fortran program is an example of how Gauss quadrature can be implemented:

```
program integration
implicit none
integer, parameter:: nq=2
real(8),dimension(nq),parameter:: xq=(-1.d0/sqrt(3.d0),+1.d0/sqrt(3.d0)/)
real(8),dimension(nq),parameter:: wq=(1.d0,1.d0/)
real(8) I
integer iq

I=0.d0
do iq=1,nq
  I=I+wq(iq)*fct(xq(iq))
end do

write(*,*) 'I=',I

contains

function fct(x)
implicit none
real(8) x,fct

!fct=3.14
!fct=x+1
fct=x**2

end function
end program
```

Modify/translate this previous program to use 5 quadrature points instead of two.

Integrate the functions

$$f_1(x) = \sin(x\pi + \pi/2) \quad f_2(x) = \sqrt{x+1} \quad f_3(x) = x^4 - x^3$$

with the 2-point and the 5-point quadrature rules.

Compare the results with the analytical values.



Exercise Quad-4

Compute analytically the integral of the function $f(x, y) = x^2 + 4y$ over the domain $\Omega = [11, 14] \times [7, 10]$.

Write a code which integrates this function by means of a 2×2 , 3×3 or 4×4 Gauss-Legendre quadrature algorithm.

4.2.4 Quadrature on triangles

quadrature_triangles.tex

Our goal is to develop a quadrature rule of the form

$$\iint_{\Delta} f(r, s) \, drds \simeq \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} f(r_{i_q}, s_{i_q})$$

We will here add two requirements: a) we would like to find quadrature rules which achieve the highest possible accuracy for the lowest possible number of quadrature points; b) we would like the quadrature points to possess some kind of symmetry. Note that the factor $1/2$ in the equation above is a convention. If $f = 1$ then the left hand term is the area of the triangle which is $1/2$. Since people usually require that $\sum_i \omega_i = 1$ then the factor $1/2$ is necessary.

Before we go any further, we need to establish that

$$\iint_{\Delta} \{1, r, s, r^2, rs, t^2, r^3, r^2s, rs^2, s^3\} \, drds = \left\{ \frac{1}{2}, \frac{1}{6}, \frac{1}{6}, \frac{1}{12}, \frac{1}{24}, \frac{1}{12}, \frac{1}{20}, \frac{1}{60}, \frac{1}{60}, \frac{1}{20} \right\}$$

where Δ stands for the reference triangle.

Gaussian quadrature of degree 1

This means that the quadrature should be accurate for $f(r, s) = \{1, r, s\}$. We then obtain:

$$\begin{aligned} \iint_{\Delta} 1 \, drds &= \frac{1}{2} &= \frac{1}{2} \sum_{i_q=1}^1 \omega_{i_q} f(r_{i_q}, s_{i_q}) &= \frac{1}{2} \omega_1 f(r_1, s_1) = \frac{1}{2} \omega_1 \\ \iint_{\Delta} r \, drds &= \frac{1}{6} &= \frac{1}{2} \sum_{i_q=1}^1 \omega_{i_q} r_{i_q} &= \frac{1}{2} \omega_1 r_1 \\ \iint_{\Delta} s \, drds &= \frac{1}{6} &= \frac{1}{2} \sum_{i_q=1}^1 \omega_{i_q} s_{i_q} &= \frac{1}{2} \omega_1 s_1 \end{aligned}$$

We then obtain $\omega_1 = 1$ and $r_1 = s_1 = \frac{1}{3}$.

Gaussian quadrature of degree 2

The quadrature should be accurate for $f(r, s) = \{1, r, s, r^2, rs, s^2\}$ so that

$$\begin{aligned}
\iint_{\Delta} 1 dr ds &= \frac{1}{2} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} \\
\iint_{\Delta} r dr ds &= \frac{1}{6} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q} \\
\iint_{\Delta} s dr ds &= \frac{1}{6} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} s_{i_q} \\
\iint_{\Delta} r^2 dr ds &= \frac{1}{12} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q}^2 \\
\iint_{\Delta} r s dr ds &= \frac{1}{24} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q} s_{i_q} \\
\iint_{\Delta} s^2 dr ds &= \frac{1}{12} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} s_{i_q}^2
\end{aligned} \tag{4.16}$$

If we set $n_q = 1$ then we have 6 equations and only three unknowns ω_1, r_1, s_1 so this will not work. If we set $n_q = 2$ then we have 6 equations and six unknowns $\omega_1, r_1, s_1, \omega_2, r_2, s_2$ but we find that the quadrature is not symmetric². If we set $n_q = 3$ then we have 6 equations and 9 unknowns and the solution is not unique. Because of symmetry we for instance have to impose $r_2 = s_3$ and $s_2 = r_3$ (points 2 and 3 are symmetric with respect to the $r = s$ line). Two common quadratures are found:

$$(r_1, s_1) = \left(\frac{1}{6}, \frac{1}{6}\right) \quad (r_2, s_2) = \left(\frac{2}{3}, \frac{1}{6}\right) \quad (r_3, s_3) = \left(\frac{1}{6}, \frac{2}{3}\right) \quad \omega_1 = \omega_2 = \omega_3 = \frac{1}{3}$$

and

$$(r_1, s_1) = \left(0, \frac{1}{2}\right) \quad (r_2, s_2) = \left(\frac{1}{2}, 0\right) \quad (r_3, s_3) = \left(\frac{1}{2}, \frac{1}{2}\right) \quad \omega_1 = \omega_2 = \omega_3 = \frac{1}{3}$$

²proof to do

Gaussian quadrature of degree 3

The quadrature should be accurate for $f(r, s) = \{1, r, s, r^2, rs, t^2, r^3, r^2s, rs^2, t^2\}$ so that

$$\begin{aligned}
\iint_{\triangle} 1 dr ds &= \frac{1}{2} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} \\
\iint_{\triangle} r dr ds &= \frac{1}{6} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q} \\
\iint_{\triangle} s dr ds &= \frac{1}{6} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} s_{i_q} \\
\iint_{\triangle} r^2 dr ds &= \frac{1}{12} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q}^2 \\
\iint_{\triangle} r s dr ds &= \frac{1}{24} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q} s_{i_q} \\
\iint_{\triangle} s^2 dr ds &= \frac{1}{12} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} s_{i_q}^2 \\
\iint_{\triangle} r^3 dr ds &= \frac{1}{20} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q}^3 \\
\iint_{\triangle} r^2 s dr ds &= \frac{1}{60} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q}^2 s_{i_q} \\
\iint_{\triangle} r s^2 dr ds &= \frac{1}{60} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} r_{i_q} s_{i_q}^2 \\
\iint_{\triangle} s^3 dr ds &= \frac{1}{20} = \frac{1}{2} \sum_{i_q=1}^{n_q} \omega_{i_q} s_{i_q}^3
\end{aligned} \tag{4.17}$$

This time $n_q = 3$ will not work since we have 10 equations. Switching to $n_q = 4$ we now have 12 unknowns and 10 equations so there are many possibilities. One of them is given in the table below. As a side note there could be cases where having a negative weight ω could be problematic.

Tabulated quadrature rules on the reference triangle

In what follows we use the following quadrature rule:

$$\iint_{\triangle} f(r, s) dr ds \simeq \sum_{i_q=1}^{n_q} \omega_{i_q} f(r_{i_q}, s_{i_q})$$

with $\sum_i \omega_i = 1/2$.

Quadrature rules for triangles can be found in Dunavant [350] (1985). The following ones are identical to those in the *ip_triangle.m* file of the MILAMIN code [299]. See also Lether [775] (1976) on the topic of computation of double integrals over a triangle.

| | r_q exact | s_q exact | w_q exact | r_q approx. | s_q approx. | w_q approx. |
|------------|----------------|----------------|----------------|--------------------------|--------------------------|--------------------------|
| $i_q = 1$ | 1/3 | 1/3 | 1/2 | | | |
| $i_q = 1$ | 1/6 | 1/6 | 1/6 | | | |
| $i_q = 2$ | 2/3 | 1/6 | 1/6 | | | |
| $i_q = 3$ | 1/6 | 2/3 | 1/6 | | | |
| $i_q = 1$ | 1/3 | 1/3 | -27/96 | | | |
| $i_q = 2$ | 0.6 | 0.2 | 25/96 | | | |
| $i_q = 3$ | 0.2 | 0.6 | 25/96 | | | |
| $i_q = 4$ | 0.2 | 0.2 | 25/96 | | | |
| $i_q = 1$ | $1 - 2g_1$ | g_1 | $w_1/2$ | 0.108103018168070 | 0.445948490915965 | |
| $i_q = 2$ | g_1 | $1 - 2g_1$ | $w_1/2$ | 0.445948490915965 | 0.108103018168070 | |
| $i_q = 3$ | g_1 | g_1 | $w_1/2$ | 0.445948490915965 | 0.445948490915965 | |
| $i_q = 4$ | $1 - 2g_2$ | g_2 | $w_2/2$ | 0.816847572980459 | 0.091576213509771 | |
| $i_q = 5$ | g_2 | $1 - 2g_2$ | $w_2/2$ | 0.091576213509771 | 0.816847572980459 | |
| $i_q = 6$ | g_2 | g_2 | $w_2/2$ | 0.091576213509771 | 0.091576213509771 | |
| $i_q = 1$ | | | | 0.091576213509771 | 0.091576213509771 | 0.109951743655322/2.0 |
| $i_q = 2$ | | | | 0.816847572980459 | 0.091576213509771 | 0.109951743655322/2.0 |
| $i_q = 3$ | | | | 0.091576213509771 | 0.816847572980459 | 0.109951743655322/2.0 |
| $i_q = 4$ | | | | 0.445948490915965 | 0.445948490915965 | 0.223381589678011/2.0 |
| $i_q = 5$ | | | | 0.108103018168070 | 0.445948490915965 | 0.223381589678011/2.0 |
| $i_q = 6$ | | | | 0.445948490915965 | 0.108103018168070 | 0.223381589678011/2.0 |
| $i_q = 1$ | | | | 0.1012865073235 | 0.1012865073235 | 0.0629695902724 |
| $i_q = 2$ | | | | 0.7974269853531 | 0.1012865073235 | 0.0629695902724 |
| $i_q = 3$ | | | | 0.1012865073235 | 0.7974269853531 | 0.0629695902724 |
| $i_q = 4$ | | | | 0.4701420641051 | 0.0597158717898 | 0.0661970763942 |
| $i_q = 5$ | | | | 0.4701420641051 | 0.4701420641051 | 0.0661970763942 |
| $i_q = 6$ | | | | 0.0597158717898 | 0.4701420641051 | 0.0661970763942 |
| $i_q = 7$ | | | | 0.3333333333333 | 0.3333333333333 | 0.1125000000000 |
| $i_q = 1$ | | | | $5.01426509658179e - 01$ | $2.49286745170910e - 01$ | $5.83931378631895e - 02$ |
| $i_q = 2$ | | | | $2.49286745170910e - 01$ | $5.01426509658179e - 01$ | $5.83931378631895e - 02$ |
| $i_q = 3$ | | | | $2.49286745170910e - 01$ | $2.49286745170910e - 01$ | $5.83931378631895e - 02$ |
| $i_q = 4$ | | | | $8.73821971016996e - 01$ | $6.30890144915020e - 02$ | $2.54224531851035e - 02$ |
| $i_q = 5$ | | | | $6.30890144915020e - 02$ | $8.73821971016996e - 01$ | $2.54224531851035e - 02$ |
| $i_q = 6$ | | | | $6.30890144915020e - 02$ | $6.30890144915020e - 02$ | $2.54224531851035e - 02$ |
| $i_q = 7$ | | | | $5.31450498448170e - 02$ | $3.10352451033784e - 01$ | $4.14255378091870e - 02$ |
| $i_q = 8$ | | | | $6.36502499121399e - 01$ | $5.31450498448170e - 02$ | $4.14255378091870e - 02$ |
| $i_q = 9$ | | | | $3.10352451033784e - 01$ | $6.36502499121399e - 01$ | $4.14255378091870e - 02$ |
| $i_q = 10$ | | | | $5.31450498448170e - 02$ | $6.36502499121399e - 01$ | $4.14255378091870e - 02$ |
| $i_q = 11$ | | | | $6.36502499121399e - 01$ | $3.10352451033784e - 01$ | $4.14255378091870e - 02$ |
| $i_q = 12$ | | | | $3.10352451033784e - 01$ | $5.31450498448170e - 02$ | $4.14255378091870e - 02$ |

where

$$g_1 = \left(8 - \sqrt{10} + \sqrt{38 - 44\sqrt{2/5}}\right) / 18 \quad g_2 = \left(8 - \sqrt{10} - \sqrt{38 - 44\sqrt{2/5}}\right) / 18$$

$$w_1 = \left(620 + \sqrt{213125 - 53320\sqrt{10}}\right) / 3720 \quad w_2 = \left(620 - \sqrt{213125 - 53320\sqrt{10}}\right) / 3720$$

All these are implemented in [STONE](#) 120 (see also [STONE](#) 112).

The case of a generic triangle

Let T be a triangular element with straight edges defined by the three vertices (x_i, y_i) ($i = 1, 2, 3$) arranged in the counter-clockwise order:

INSERT FIGURE

We now need to evaluate the following integral over T :

$$I = \iint_T f(x, y) \, dx dy$$

We will proceed similarly to the quadrilateral case: first transform ('map') this triangle into the reference triangle Δ and then use the adequate quadrature rule.

We can base this transformation ('mapping') on the linear (' P_1 ') basis functions³:

$$\begin{aligned}\mathcal{N}_1(r, s) &= 1 - r - s \\ \mathcal{N}_2(r, s) &= r \\ \mathcal{N}_3(r, s) &= s\end{aligned}$$

with

$$\begin{aligned}x = P(r, s) &= \sum_{i=1}^3 x_i \mathcal{N}_i(r, s) = x_1 \mathcal{N}_1(r, s) + x_2 \mathcal{N}_2(r, s) + x_3 \mathcal{N}_3(r, s) \\ y = Q(r, s) &= \sum_{i=1}^3 y_i \mathcal{N}_i(r, s) = y_1 \mathcal{N}_1(r, s) + y_2 \mathcal{N}_2(r, s) + y_3 \mathcal{N}_3(r, s)\end{aligned}$$

For example when $(r, s) \rightarrow (0, 0)$ we find $(x, y) \rightarrow (x_1, y_1)$, when $(r, s) \rightarrow (1, 0)$ we find $(x, y) \rightarrow (x_2, y_2)$, and when $(r, s) \rightarrow (0, 1)$ we find $(x, y) \rightarrow (x_3, y_3)$.

Then we have

$$I = \iint_T f(x, y) \, dx dy = \iint_{\Delta} f(P(r, s), Q(r, s)) |J(r, s)| \, dr ds$$

where $J(r, s)$ is the Jacobian of the transformation:

$$J(r, s) = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial y}{\partial r} \\ \frac{\partial x}{\partial s} & \frac{\partial y}{\partial s} \end{vmatrix} = \begin{vmatrix} \frac{\partial}{\partial r} \sum_{i=1}^3 x_i \mathcal{N}(r_i, s_i) & \frac{\partial}{\partial r} \sum_{i=1}^3 y_i \mathcal{N}(r_i, s_i) \\ \frac{\partial}{\partial s} \sum_{i=1}^3 x_i \mathcal{N}(r_i, s_i) & \frac{\partial}{\partial s} \sum_{i=1}^3 y_i \mathcal{N}(r_i, s_i) \end{vmatrix} = \begin{vmatrix} (-x_1 + x_2) & (-y_1 + y_2) \\ (-x_1 + x_3) & (-y_1 + y_3) \end{vmatrix} = 2\mathcal{A}_T$$

where \mathcal{A}_T is the area of triangle T . In the end we have

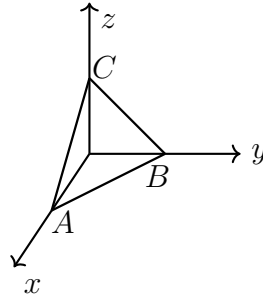
$$I = \iint_T f(x, y) \, dx dy = 2\mathcal{A}_T \iint_{\Delta} f(P(r, s), Q(r, s)) \, dr ds$$

and the rhs integral can then be computed by means of quadrature.

4.2.5 A mathematical recreation: computing the volume of a tetrahedron

Let us find the volume of tetrahedron bounded by the planes passing through the points A(1,0,0), B(0,1,0), C(0,0,1) and the coordinate planes Oxy, Oxz and Oyz.

³See Section 7.6



The equation of the plane is $x + y + z = 1$, or $z = 1 - x - y$. Hence, the limits of integration over the variable z range in the interval from $z = 0$ to $z = 1 - x - y$. Now we can calculate the volume of the tetrahedron:

$$\begin{aligned}
 V &= \iiint dx \, dy \, dz \\
 &= \int_0^1 dx \int_0^{1-x} dy \int_0^{1-x-y} dz \\
 &= \int_0^1 dx \int_0^{1-x} dy (1 - x - y) \\
 &= \int_0^1 dx \left[y - xy - \frac{1}{2}y^2 \right]_0^{1-x} \\
 &= \int_0^1 dx \left((1-x) - x(1-x) - \frac{1}{2}(1-x)^2 \right) \\
 &= \int_0^1 dx \left(\frac{1}{2} - x + \frac{1}{2}x^2 \right) \\
 &= \left[\frac{1}{2}x - \frac{1}{2}x^2 + \frac{1}{6}x^3 \right]_0^1 \\
 &= \frac{1}{6}
 \end{aligned} \tag{4.18}$$

We will use this result in the following section.

4.2.6 Quadrature on tetrahedra

quadrature_tetrahedra.tex

Remark. In what follows the coefficients in the tables are not the reduced coordinates of the quadrature points but the coefficients corresponding to the 4 nodes.

Quadrature rules on tetrahedra take the form:

$$\int \int \int_{el} f(x, y, z) dx dy dz = V_{el} \sum_{iq=1}^{n_{qel}} w_{iq} f(\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq})$$

or, that is to say:

$$\int \int \int_{el} f(x, y, z) dx dy dz = \sum_{iq=1}^{n_{qel}} (w_{iq} V_{el}) f(\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq})$$

with in our case $V_{el} = 1/6$.

In the literature it can be found that a one point quadrature is characterised by

$$w_{iq} = 1 \quad \xi_1^{iq} = \xi_2^{iq} = \xi_3^{iq} = \xi_4^{iq} = 0.25$$

i.e, the coordinates of the single point are given by:

$$x_{iq} = \sum_{i=1}^4 \xi_i^{iq} x_i = \frac{1}{4}(x_1 + x_2 + x_3 + x_4)$$

Same for y and z coordinates.

A four-point quadrature rule is characterised by $w_{iq} = V_{el} * 0.25 = 1/24 \simeq 0.0416666666666667$ and

| | ξ_1 | ξ_2 | ξ_3 | ξ_4 |
|------|-------------------|-------------------|-------------------|-------------------|
| iq=1 | 0.585410196624969 | 0.138196601125011 | 0.138196601125011 | 0.138196601125011 |
| iq=2 | 0.138196601125011 | 0.585410196624969 | 0.138196601125011 | 0.138196601125011 |
| iq=3 | 0.138196601125011 | 0.138196601125011 | 0.585410196624969 | 0.138196601125011 |
| iq=4 | 0.138196601125011 | 0.138196601125011 | 0.138196601125011 | 0.585410196624969 |

We then have:

$$r_{iq} = \sum_{i=1}^4 \xi_i^{iq} x_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (r_1, r_2, r_3, r_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 1, 0, 0) = \xi_2^{iq}$$

$$s_{iq} = \sum_{i=1}^4 \xi_i^{iq} y_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (s_1, s_2, s_3, s_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 0, 1, 0) = \xi_3^{iq}$$

$$t_{iq} = \sum_{i=1}^4 \xi_i^{iq} z_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (t_1, t_2, t_3, t_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 0, 0, 1) = \xi_4^{iq}$$

Finally:

| | r_q | s_q | t_q | w_q |
|------|-------------------|-------------------|-------------------|--------------------|
| iq=1 | 0.138196601125011 | 0.138196601125011 | 0.138196601125011 | 0.0416666666666667 |
| iq=2 | 0.585410196624969 | 0.138196601125011 | 0.138196601125011 | 0.0416666666666667 |
| iq=3 | 0.138196601125011 | 0.585410196624969 | 0.138196601125011 | 0.0416666666666667 |
| iq=4 | 0.138196601125011 | 0.138196601125011 | 0.585410196624969 | 0.0416666666666667 |

4.2.7 The Gauss-Lobatto approach

All what we have seen above falls under the Gauss-Legendre quadrature method. There is however another somewhat common quadrature method: the Gauss-Lobatto quadrature. . It is similar to Gaussian quadrature with the following important differences: 1) There are integration points in the interval but they also always include the end points of the integration interval; 2) It is accurate for polynomials up to degree $2n - 3$, where n is the number of integration points.

In 1D, it reads:

$$\int_{-1}^{+1} f(x)dx = \frac{2}{n(n-1)}[f(-1) + f(1)] + \sum_{i=2}^{n-1} w_i f(x_i)$$

The locations and weights of the integration points are as follows:

| n | x_{iq} | w_{iq} | x_{iq} (approx) | w_{iq} (approx) |
|---|-------------------------------------------------|--------------------------|-------------------|-------------------|
| 3 | 0 | 4/3 | | |
| | ± 1 | 1/3 | | |
| 4 | $\pm \sqrt{\frac{1}{5}}$ | 5/6 | | |
| | ± 1 | 1/6 | | |
| 5 | 0 | 32/45 | | |
| | $\pm \sqrt{\frac{3}{7}}$ | 49/90 | | |
| | ± 1 | 1/10 | | |
| 6 | $\pm \sqrt{\frac{1}{3} - \frac{2\sqrt{7}}{21}}$ | $\frac{14+\sqrt{7}}{30}$ | | |
| | $\pm \sqrt{\frac{1}{3} + \frac{2\sqrt{7}}{21}}$ | $\frac{14-\sqrt{7}}{30}$ | | |
| | ± 1 | 1/15 | | |
| | | | | |

4.2.8 Computing the 'real' coordinates of the quadrature points and other considerations

The quadrature point coordinates are always given in (what I call) reduced coordinates, i.e. between -1 and 1. However, one sometimes need their equivalent in the x, y Cartesian space. This is trivial once one remembers that within an element, a field f is represented as follow:

$$f(r, s) = \sum_{i=1}^m \mathcal{N}_i(r, s) f_i$$

where m is the number of nodes, r and s are the reduced coordinates and \mathcal{N}_i are the basis functions. The value of f at a quadrature point (r_q, s_q) is then simply

$$f(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) f_i$$

If we now take $f = x$, then

$$x_q = x(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) x_i$$

and

$$y_q = y(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) y_i$$

where x_i and y_i are the Cartesian coordinates of the nodes. This is then easily extended to three dimensions:

$$x_q = x(r_q, s_q, t_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q, t_q) x_i$$

$$y_q = y(r_q, s_q, t_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q, t_q) y_i$$

$$z_q = z(r_q, s_q, t_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q, t_q) z_i$$

or,

$$\vec{r}_q = \vec{r}(r_q, s_q, t_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q, t_q) \vec{r}_i$$

This also applies to other fields such as velocity, temperature, or even strain rate components (as long as the strain rate values have previously been computed on the nodes):

$$\vec{\mathbf{v}}_q = \vec{\mathbf{v}}(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) \vec{\mathbf{v}}_i$$

$$T_q = T(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) T_i$$

$$\dot{\epsilon}_{xy,q} = \dot{\epsilon}_{xy}(r_q, s_q) = \sum_{i=1}^m \mathcal{N}_i(r_q, s_q) \dot{\epsilon}_{xy,i}$$

Chapter 5

The building blocks of the Finite Element Method

chapter_fem0.tex

5.1 A bit of FE terminology

terminology.tex

We introduce here some terminology for efficient element descriptions [488]:

- For triangles/tetrahedra, the designation $P_m \times P_n$ means that each component of the velocity is approximated by continuous piecewise complete Polynomials of degree m and pressure by continuous piecewise complete Polynomials of degree n . For example $P_2 \times P_1$ means

$$u^h(x, y) \sim a_1 + a_2x + a_3y + a_4xy + a_5x^2 + a_6y^2$$

with similar approximations for v , and

$$p^h(x, y) \sim b_1 + b_2x + b_3y$$

Both velocity and pressure are continuous across element boundaries, and each triangular element contains 6 velocity nodes and three pressure nodes.

- For the same families, $P_m \times P_{-n}$ is as above, except that pressure is approximated via piecewise *discontinuous* polynomials of degree n . For instance, $P_2 \times P_{-1}$ is the same as P_2P_1 except that pressure is now an independent linear function in each element and therefore discontinuous at element boundaries.
- For quadrilaterals/hexahedra, the designation $Q_m \times Q_n$ means that each component of the velocity is approximated by a continuous piecewise polynomial of degree m *in each direction* on the quadrilateral and likewise for pressure, except that the polynomial is of degree n . For instance, $Q_2 \times Q_1$ means

$$u^h(x, y) \sim a_1 + a_2x + a_3y + a_4xy + a_5x^2 + a_6y^2 + a_7x^2y + a_8xy^2 + a_9x^2y^2$$

and

$$p^h(x, y) \sim b_1 + b_2x + b_3y + b_4xy$$

- For these same families, $Q_m \times Q_{-n}$ is as above, except that the pressure approximation is not continuous at element boundaries.

- Again for the same families, $Q_m \times P_{-n}$ indicates the same velocity approximation with a pressure approximation that is a discontinuous complete piecewise polynomial of degree n (not of degree n in each direction !)
- The designation P_m^+ or Q_m^+ means that some sort of bubble function was added to the polynomial approximation for the velocity. You may also find the term 'enriched element' in the literature.
- Finally, for $n = 0$, we have piecewise-constant pressure, and we omit the minus sign for simplicity.

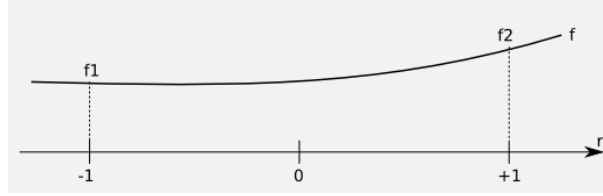
Another point which needs to be clarified is the use of so-called 'conforming elements' (or 'non-conforming elements'). Following again Gresho & Sani [488], conforming velocity elements are those for which the basis functions form a subset of H^1 for the continuous problem (the first derivatives and their squares are integrable in Ω). For instance, the rotated $Q_1 \times P_0$ element of Rannacher and Turek (see section 7.3.9) is such that the velocity is discontinuous across element edges, so that the derivative does not exist there. Another typical example of non-conforming element is the Crouzeix-Raviart element [290].

5.2 Elements and basis functions in 1D

elements1D.tex

5.2.1 Linear basis functions (Q_1)

Let $f(r)$ be a C^1 function on the interval $[-1 : 1]$ with $f(-1) = f_1$ and $f(1) = f_2$.



Let us assume that the function $f(r)$ is to be approximated on $[-1, 1]$ by the first order polynomial

$$f^h(r) = a + br \quad (5.1)$$

Then it must fulfil

$$\begin{aligned} f^h(r = -1) &= a - b = f_1 \\ f^h(r = +1) &= a + b = f_2 \end{aligned}$$

This leads to

$$\begin{aligned} a &= \frac{1}{2}(f_1 + f_2) \\ b &= \frac{1}{2}(-f_1 + f_2) \end{aligned} \quad (5.2)$$

and then replacing a, b in Eq. (5.1) by the above values one gets

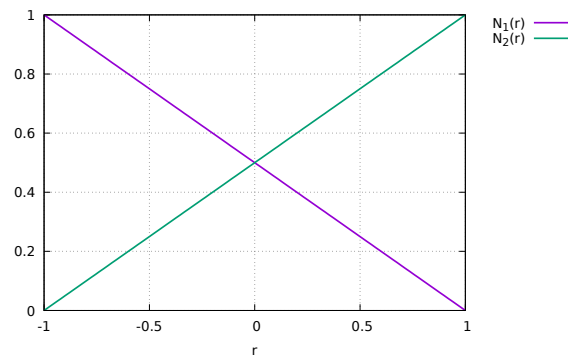
$$f^h(r) = \left[\frac{1}{2}(1 - r) \right] f_1 + \left[\frac{1}{2}(1 + r) \right] f_2$$

or

$$f^h(r) = \sum_{i=1}^2 N_i(r) f_i$$

with

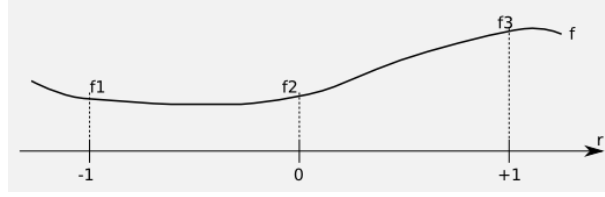
$$\begin{aligned} N_1(r) &= \frac{1}{2}(1 - r) \\ N_2(r) &= \frac{1}{2}(1 + r) \end{aligned} \quad (5.3)$$



Plot of the two linear functions $N_1(r)$ and $N_2(r)$.

5.2.2 Quadratic basis functions (Q_2)

Let $f(r)$ be a C^1 function on the interval $[-1 : 1]$ with $f(-1) = f_1$, $f(0) = f_2$ and $f(1) = f_3$.



Let us assume that the function $f(r)$ is to be approximated on $[-1, 1]$ by the second order polynomial $f^h(r)$:

$$f(r) = a + br + cr^2 \quad (5.4)$$

Then it must fulfil

$$\begin{aligned} f^h(r = -1) &= a - b + c = f_1 \\ f^h(r = 0) &= a = f_2 \\ f^h(r = +1) &= a + b + c = f_3 \end{aligned}$$

This leads to

$$\begin{aligned} a &= f_2 \\ b &= \frac{1}{2}(-f_1 + f_3) \\ c &= \frac{1}{2}(f_1 + f_3 - 2f_2) \end{aligned} \quad (5.5)$$

and then replacing a, b, c in Eq. (5.4) by the above values one gets

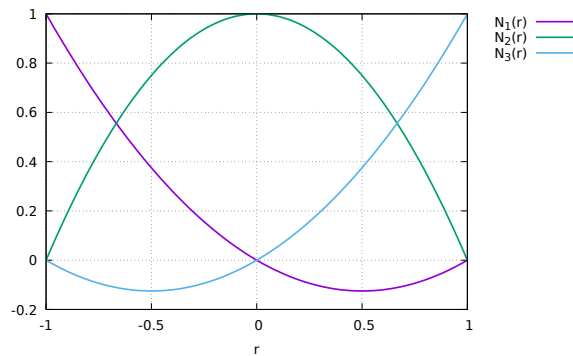
$$f^h(r) = \left[\frac{1}{2}r(r-1) \right] f_1 + (1-r^2)f_2 + \left[\frac{1}{2}r(r+1) \right] f_3$$

or,

$$f^h(r) = \sum_{i=1}^3 N_i(r) f_i$$

with

$$\begin{aligned} N_1(r) &= \frac{1}{2}r(r-1) \\ N_2(r) &= (1-r^2) \\ N_3(r) &= \frac{1}{2}r(r+1) \end{aligned} \quad (5.6)$$



Plot of the three quadratic functions $N_1(r)$, $N_2(r)$ and $N_3(r)$.

Note that Q_2 basis functions can take negative values.

We will later need the first-order derivatives of these functions:

$$\begin{aligned}\frac{\partial N_1}{\partial r} &= r - \frac{1}{2} \\ \frac{\partial N_2}{\partial r} &= -2r \\ \frac{\partial N_3}{\partial r} &= r + \frac{1}{2}\end{aligned}\tag{5.7}$$

5.2.3 Cubic basis functions (Q_3)

We proceed as previously by assuming that the third-order polynomial representation of function $f(r)$ is given by

$$f^h(r) = a + br + cr^2 + dr^3$$

with the nodes at position -1, -1/3, +1/3 and +1. It then must fulfil all four conditions:

$$\begin{aligned}f(-1) &= a - b + c - d = f_1 \\ f(-1/3) &= a - \frac{b}{3} + \frac{c}{9} - \frac{d}{27} = f_2 \\ f(+1/3) &= a - \frac{b}{3} + \frac{c}{9} - \frac{d}{27} = f_3 \\ f(+1) &= a + b + c + d = f_4\end{aligned}$$

Adding the first and fourth equation and the second and third, one arrives at

$$f_1 + f_4 = 2a + 2c \qquad f_2 + f_3 = 2a + \frac{2c}{9}$$

and finally:

$$\begin{aligned}a &= \frac{1}{16} (-f_1 + 9f_2 + 9f_3 - f_4) \\ c &= \frac{9}{16} (f_1 - f_2 - f_3 + f_4)\end{aligned}$$

Combining the original 4 equations in a different way yields

$$2b + 2d = f_4 - f_1 \qquad \frac{2b}{3} + \frac{2d}{27} = f_3 - f_2$$

so that

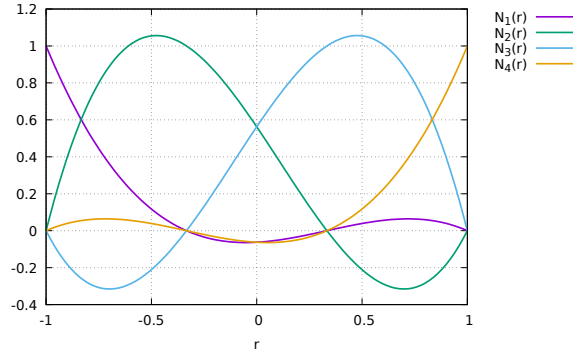
$$\begin{aligned}b &= \frac{1}{16} (f_1 - 27f_2 + 27f_3 - f_4) \\ d &= \frac{9}{16} (-f_1 + 3f_2 - 3f_3 + f_4)\end{aligned}$$

Finally,

$$\begin{aligned}
 f^h(r) &= a + b + cr^2 + dr^3 \\
 &= \frac{1}{16}(-1 + r + 9r^2 - 9r^3)f_1 \\
 &\quad + \frac{1}{16}(9 - 27r - 9r^2 + 27r^3)f_2 \\
 &\quad + \frac{1}{16}(9 + 27r - 9r^2 - 27r^3)f_3 \\
 &\quad + \frac{1}{16}(-1 - r + 9r^2 + 9r^3)f_4 \\
 &= \sum_{i=1}^4 N_i(r)f_i
 \end{aligned}$$

where (see also for example [779, p49])

$$\begin{aligned}
 N_1 &= \frac{1}{16}(-1 + r + 9r^2 - 9r^3) \\
 N_2 &= \frac{1}{16}(9 - 27r - 9r^2 + 27r^3) \\
 N_3 &= \frac{1}{16}(9 + 27r - 9r^2 - 27r^3) \\
 N_4 &= \frac{1}{16}(-1 - r + 9r^2 + 9r^3)
 \end{aligned}$$



Plot of the four cubic functions $N_1(r)$, $N_2(r)$, $N_3(r)$ and $N_4(r)$.

Let us now verify that these functions can represent any polynomial function up to third order:

- Let us assume $f(r) = C$, then

$$f^h(r) = \sum N_i(r)f_i = \sum_i N_i C = C \sum_i N_i = C$$

so that a constant function is exactly reproduced, as expected. This is a very important property of the N_i functions: They must fulfil $\sum_i N_i = 1$.

- Let us assume $f(r) = r$, then $f_1 = -1$, $f_2 = -1/3$, $f_3 = 1/3$ and $f_4 = +1$. We then have

$$\begin{aligned}
f^h(r) &= \sum N_i(r) f_i \\
&= -N_1(r) - \frac{1}{3}N_2(r) + \frac{1}{3}N_3(r) + N_4(r) \\
&= [-(-1 + r + 9r^2 - 9r^3) \\
&\quad - \frac{1}{3}(9 - 27r - 9r^2 - 27r^3) \\
&\quad + \frac{1}{3}(9 + 27r - 9r^2 + 27r^3) \\
&\quad + (-1 - r + 9r^2 + 9r^3)] / 16 \\
&= [-r + 9r + 9r - r] / 16 + \dots 0 \dots \\
&= r
\end{aligned} \tag{5.8}$$

- The cases $f(r) = r^2$ and $f(r) = r^3$ are left as exercise.

The basis functions first-order derivatives are given by

$$\begin{aligned}
\frac{\partial N_1}{\partial r} &= \frac{1}{16}(1 + 18r - 27r^2) \\
\frac{\partial N_2}{\partial r} &= \frac{1}{16}(-27 - 18r + 81r^2) \\
\frac{\partial N_3}{\partial r} &= \frac{1}{16}(+27 - 18r - 81r^2) \\
\frac{\partial N_4}{\partial r} &= \frac{1}{16}(-1 + 18r + 27r^2)
\end{aligned}$$

We can also verify that the derivatives are also properly approximated:

- Let us assume $f(r) = C$, then

$$\begin{aligned}
\frac{\partial f^h}{\partial r} &= \sum_i \frac{\partial N_i}{\partial r} f_i \\
&= C \sum_i \frac{\partial N_i}{\partial r} \\
&= \frac{C}{16} [(1 + 18r - 27r^2) + (-27 - 18r + 81r^2) + (+27 - 18r - 81r^2) + (-1 + 18r + 27r^2)] \\
&= 0
\end{aligned}$$

- Let us assume $f(r) = r$, then $f_1 = -1$, $f_2 = -1/3$, $f_3 = 1/3$ and $f_4 = +1$. We then have

$$\begin{aligned}
\frac{\partial f^h}{\partial r} &= \sum_i \frac{\partial N_i}{\partial r} f_i \\
&= \frac{1}{16} [-(1 + 18r - 27r^2) - \frac{1}{3}(-27 - 18r + 81r^2) + \frac{1}{3}(27 - 18r - 81r^2) + (-1 + 18r + 27r^2)] \\
&= \frac{1}{16} [-2 + 18 + 54r^2 - 54r^2] \\
&= 1
\end{aligned}$$

- Let us assume $f(r) = r^2$, then $f_1 = 1$, $f_2 = 1/9$, $f_3 = 1/9$ and $f_4 = 1$. We then have

$$\begin{aligned}
\frac{\partial f^h}{\partial r} &= \sum_i \frac{\partial N_i}{\partial r} f_i \\
&= \frac{1}{16} \left[(1 + 18r - 27r^2) + \frac{1}{9}(-27 - 18r + 81r^2) + \frac{1}{9}(27 - 18r - 81r^2) + (-1 + 18r + 27r^2) \right] \\
&= \frac{1}{16}(32r) \\
&= 2r
\end{aligned} \tag{5.9}$$

as expected.

5.2.4 Quartic basis functions (Q_4)

The 1D basis polynomial is given by

$$f_h(r) = a + br + cr^2 + dr^3 + er^4$$

with the nodes at position -1, -1/2, 0, +1/2 and +1. The function $f^h(r)$ must then fulfil

$$\begin{aligned}
f_h(-1) &= a - b + c - d + e = f_1 \\
f_h(-1/2) &= a - \frac{b}{2} + \frac{c}{4} - \frac{d}{8} + \frac{e}{16} = f_2 \\
f_h(0) &= a = f_3 \\
f_h(+1/2) &= a + \frac{b}{2} + \frac{c}{4} + \frac{d}{8} + \frac{e}{16} = f_4 \\
f_h(+1) &= a + b + c + d + e = f_5
\end{aligned}$$

or,

$$\begin{pmatrix} 1 & -1 & 1 & -1 & 1 \\ 1 & -1/2 & 1/4 & -1/8 & 1/16 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 1/2 & 1/4 & 1/8 & 1/16 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \\ d \\ e \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \end{pmatrix} \tag{5.10}$$

The third line gives $a = f_3$ so that

$$\underbrace{\begin{pmatrix} -1 & 1 & -1 & 1 \\ -1/2 & 1/4 & -1/8 & 1/16 \\ 1/2 & 1/4 & 1/8 & 1/16 \\ 1 & 1 & 1 & 1 \end{pmatrix}}_A \begin{pmatrix} b \\ c \\ d \\ e \end{pmatrix} = \begin{pmatrix} f_1 - f_3 \\ f_2 - f_3 \\ f_4 - f_3 \\ f_5 - f_3 \end{pmatrix} \tag{5.11}$$

The inverse of the matrix A is:

$$A^{-1} = \frac{1}{6} \begin{pmatrix} 1 & -8 & 8 & -1 \\ -1 & 16 & 16 & -1 \\ -4 & 8 & -8 & 4 \\ 4 & -16 & -16 & 4 \end{pmatrix}$$

so that

$$\begin{pmatrix} b \\ c \\ d \\ e \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 1 & -8 & 8 & -1 \\ -1 & 16 & 16 & -1 \\ -4 & 8 & -8 & 4 \\ 4 & -16 & -16 & 4 \end{pmatrix} \cdot \begin{pmatrix} f_1 - f_3 \\ f_2 - f_3 \\ f_4 - f_3 \\ f_5 - f_3 \end{pmatrix}$$

and then

$$b = \frac{1}{6} (f_1 - 8f_2 + 8f_4 - f_5) \quad (5.12)$$

$$c = \frac{1}{6} (-f_1 + 16f_2 - 30f_3 + 16f_4 - f_5) \quad (5.13)$$

$$d = \frac{1}{6} (-4f_1 + 8f_2 - 8f_4 + 4f_5) \quad (5.14)$$

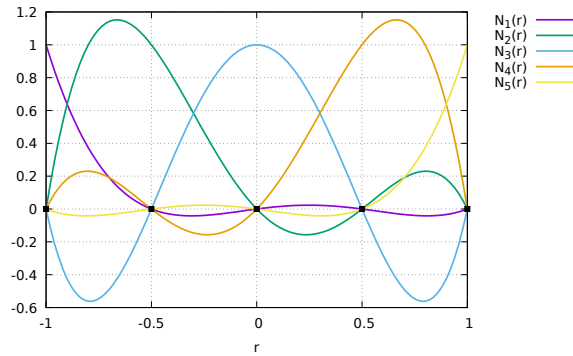
$$e = \frac{1}{6} (4f_1 - 16f_2 + 24f_3 - 16f_4 + 4f_5) \quad (5.15)$$

Finally

$$\begin{aligned} f_h(r) &= a + br + cr^2 + dr^3 + er^4 \\ &= f_3 + \frac{1}{6} (f_1 - 8f_2 + 8f_4 - f_5) r + \frac{1}{6} (-f_1 + 16f_2 - 30f_3 + 16f_4 - f_5) r^2 + \\ &\quad \frac{1}{6} (-4f_1 + 8f_2 - 8f_4 + 4f_5) r^3 + \frac{1}{6} (4f_1 - 16f_2 + 24f_3 - 16f_4 + 4f_5) r^4 \\ &= \frac{1}{6} (r - r^2 - 4r^3 + 4r^4) f_1 \\ &\quad + \frac{1}{6} (-8r + 16r^2 + 8r^3 - 16r^4) f_2 \\ &\quad + (1 - 5r^2 + 4r^4) f_3 \\ &\quad + \frac{1}{6} (8r + 16r^2 - 8r^3 - 16r^4) f_4 \\ &\quad + \frac{1}{6} (-r - r^2 + 4r^3 + 4r^4) f_5 \end{aligned}$$

with

$$\begin{aligned} N_1(r) &= \frac{1}{6} (r - r^2 - 4r^3 + 4r^4) \\ N_2(r) &= \frac{1}{6} (-8r + 16r^2 + 8r^3 - 16r^4) \\ N_3(r) &= (1 - 5r^2 + 4r^4) \\ N_4(r) &= \frac{1}{6} (8r + 16r^2 - 8r^3 - 16r^4) \\ N_5(r) &= \frac{1}{6} (-r - r^2 + 4r^3 + 4r^4) \end{aligned} \quad (5.16)$$



Plot of the 5 quartic basis functions.

The basis functions derivative are given by

$$\begin{aligned}
\frac{\partial N_1}{\partial r} &= \frac{1}{6}(1 - 2r - 12r^2 + 16r^3) \\
\frac{\partial N_2}{\partial r} &= \frac{1}{6}(-8 + 32r + 24r^2 - 64r^3) \\
\frac{\partial N_3}{\partial r} &= -10r + 16r^3 \\
\frac{\partial N_4}{\partial r} &= \frac{1}{6}(8 + 32r - 24r^2 - 64r^3) \\
\frac{\partial N_5}{\partial r} &= \frac{1}{6}(-1 - 2r + 12r^2 + 16r^3)
\end{aligned} \tag{5.17}$$

5.2.5 Fifth-order basis functions (Q_5)

Following the methodology presented hereafter for Q_6 , we arrive at

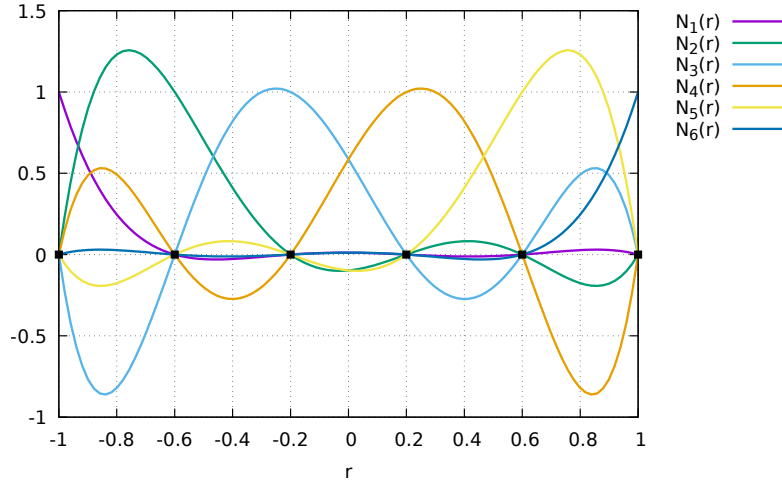
$$\begin{aligned}
\mathcal{N}_1(r) &= -\frac{625}{768}(r + \frac{3}{5})(r + \frac{1}{5})(r - \frac{1}{5})(r - \frac{3}{5})(r - 1) \\
\mathcal{N}_2(r) &= \frac{3125}{768}(r + 1)(r + \frac{1}{5})(r - \frac{1}{5})(r - \frac{3}{5})(r - 1) \\
\mathcal{N}_3(r) &= -\frac{3125}{384}(r + 1)(r + \frac{3}{5})(r - \frac{1}{5})(r - \frac{3}{5})(r - 1) \\
\mathcal{N}_4(r) &= \frac{3125}{384}(r + 1)(r + \frac{3}{5})(r + \frac{1}{5})(r - \frac{3}{5})(r - 1) \\
\mathcal{N}_5(r) &= -\frac{3125}{768}(r + 1)(r + \frac{3}{5})(r + \frac{1}{5})(r - \frac{1}{5})(r - 1) \\
\mathcal{N}_6(r) &= \frac{625}{768}(r + 1)(r + \frac{3}{5})(r + \frac{1}{5})(r - \frac{1}{5})(r - \frac{3}{5})
\end{aligned} \tag{5.18}$$

or,

$$\begin{aligned}
\mathcal{N}_1(r) &= -\frac{1}{768}(625r^5 - 625r^4 - 250r^3 + 250r^2 + 9r - 9) \\
\mathcal{N}_2(r) &= \frac{25}{768}(125r^5 - 75r^4 - 130r^3 + 78r^2 + 5r - 3) \\
\mathcal{N}_3(r) &= -\frac{25}{384}(125r^5 - 25r^4 - 170r^3 + 34r^2 + 45r - 9) \\
\mathcal{N}_4(r) &= \frac{25}{384}(125r^5 + 25r^4 - 170r^3 - 34r^2 + 45r + 9) \\
\mathcal{N}_5(r) &= -\frac{25}{768}(125r^5 + 75r^4 - 130r^3 - 78r^2 + 5r + 3) \\
\mathcal{N}_6(r) &= \frac{1}{768}(625r^5 + 625r^4 - 250r^3 - 250r^2 + 9r + 9)
\end{aligned} \tag{5.19}$$

with the derivatives given by

$$\begin{aligned}
\frac{\partial N_1}{\partial r} &= -\frac{1}{768}(3125r^4 - 2500r^3 - 750r^2 + 500r + 9) \\
\frac{\partial N_2}{\partial r} &= \frac{25}{768}(625r^4 - 300r^3 - 390r^2 + 156r + 5) \\
\frac{\partial N_3}{\partial r} &= -\frac{25}{384}(625r^4 - 100r^3 - 510r^2 + 68r + 45) \\
\frac{\partial N_4}{\partial r} &= \frac{25}{384}(625r^4 + 100r^3 - 510r^2 - 68r + 45) \\
\frac{\partial N_5}{\partial r} &= -\frac{25}{768}(625r^4 + 300r^3 - 390r^2 - 156r + 5) \\
\frac{\partial N_6}{\partial r} &= \frac{1}{768}(3125r^4 + 2500r^3 - 750r^2 - 500r + 9)
\end{aligned} \tag{5.20}$$



Plot of the 6 fifth-order basis functions.

These functions are used in [STONE](#) ??.

5.2.6 Sixth-order basis functions (Q_6)

The 1D basis polynomial is given by

$$f_h(r) = a + br + cr^2 + dr^3 + er^4 + fr^5 + gr^6$$

with the nodes at position -1, -2/3, -1/3, 0, +1/3, +2/3 and +1. The function $f^h(r)$ must then fulfil

$$\begin{pmatrix} 1 & -1 & 1 & -1 & 1 & -1 & 1 \\ 1 & -\frac{2}{3} & \frac{4}{9} & -\frac{8}{27} & \frac{16}{81} & -\frac{32}{243} & \frac{64}{729} \\ 1 & -\frac{1}{3} & \frac{1}{9} & -\frac{1}{27} & \frac{1}{81} & -\frac{1}{243} & \frac{1}{729} \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & \frac{1}{3} & \frac{1}{9} & \frac{1}{27} & \frac{1}{81} & \frac{1}{243} & \frac{1}{729} \\ 1 & \frac{2}{3} & \frac{4}{9} & \frac{8}{27} & \frac{16}{81} & \frac{32}{243} & \frac{64}{729} \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \\ e \\ f \\ g \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \end{pmatrix}$$

The middle line yields $a = f_4$, so that we have:

$$\begin{pmatrix} -1 & 1 & -1 & 1 & -1 & 1 \\ -\frac{2}{3} & \frac{4}{9} & -\frac{8}{27} & \frac{16}{81} & -\frac{32}{243} & \frac{64}{729} \\ -\frac{1}{3} & \frac{1}{9} & -\frac{1}{27} & \frac{1}{81} & -\frac{1}{243} & \frac{1}{729} \\ \frac{1}{3} & \frac{1}{9} & \frac{1}{27} & \frac{1}{81} & \frac{1}{243} & \frac{1}{729} \\ \frac{2}{3} & \frac{4}{9} & \frac{8}{27} & \frac{16}{81} & \frac{32}{243} & \frac{64}{729} \\ \frac{1}{3} & \frac{1}{9} & \frac{1}{27} & \frac{1}{81} & \frac{1}{243} & \frac{1}{729} \\ 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} b \\ c \\ d \\ e \\ f \\ g \end{pmatrix} = \begin{pmatrix} f_1 - f_4 \\ f_2 - f_4 \\ f_3 - f_4 \\ f_5 - f_4 \\ f_6 - f_4 \\ f_7 - f_4 \end{pmatrix}$$

Multiplying all lines by 729, we obtain:

$$\frac{1}{729} \begin{pmatrix} -729 & 729 & -729 & 729 & -729 & 729 \\ -486 & 324 & -216 & 144 & -96 & 64 \\ -243 & 81 & -27 & 9 & -3 & 1 \\ 243 & 81 & 27 & 9 & 3 & 1 \\ 486 & 324 & 216 & 144 & 96 & 64 \\ 729 & 729 & 729 & 729 & 729 & 729 \end{pmatrix} \cdot \begin{pmatrix} b \\ c \\ d \\ e \\ f \\ g \end{pmatrix} = \begin{pmatrix} f_1 - f_4 \\ f_2 - f_4 \\ f_3 - f_4 \\ f_5 - f_4 \\ f_6 - f_4 \\ f_7 - f_4 \end{pmatrix}$$

The inverse¹ of this matrix is:

$$\begin{array}{cccccc} -0.00006859 & 0.00061728 & -0.00308642 & 0.00308642 & -0.00061728 & 0.00006859 \\ 0.00006859 & -0.00092593 & 0.00925926 & 0.00925926 & -0.00092593 & 0.00006859 \\ 0.00077160 & -0.00617284 & 0.01003086 & -0.01003086 & 0.00617284 & -0.00077160 \\ -0.00077160 & 0.00925926 & -0.03009259 & 0.03009259 & 0.00925926 & -0.00077160 \\ -0.00138889 & 0.00555556 & -0.00694444 & 0.00694444 & -0.00555556 & 0.00138889 \\ 0.00138889 & -0.00833333 & 0.02083333 & 0.02083333 & -0.00833333 & 0.00138889 \end{array}$$

Obviously, this is not a very practical approach anymore. One could solve the system by hand, making sure to keep fractions but it will be cumbersome. Let us turn to another approach.

The nodes inside the reference element are as follows:

$$\begin{array}{cccccc} (1) & (2) & (3) & (4) & (5) & (6) & (7) \\ -|---|----|----+---|----|---|--- \\ -1 & -2/3 & -1/3 & 0 & 1/3 & 2/3 & 1 \end{array}$$

Basis function $\mathcal{N}_1(r)$ is a 6th order polynomial expression that should be 1 at node 1, and 0 at others, i.e. at $r = -2/3, -1/3, 0, 1/3, 2/3, 1$. It must then be of the form:

$$\mathcal{N}_1(r) = \alpha(r + \frac{2}{3})(r + \frac{1}{3})(r)(r - \frac{1}{3})(r - \frac{2}{3})(r - 1)$$

When evaluated at $r = -1$, we get

$$\mathcal{N}_1(r = -1) = \alpha(-\frac{1}{3})(-\frac{2}{3})(-1)(-\frac{4}{3})(-\frac{5}{3})(-2) = \alpha \frac{80}{81}$$

Since this quantity must be 1, we have

$$1 = \alpha \frac{80}{81} \quad \rightarrow \quad \alpha = \frac{81}{80}$$

so that

$$\begin{aligned} \mathcal{N}_1(r) &= \frac{81}{80}(r + \frac{2}{3})(r + \frac{1}{3})(r)(r - \frac{1}{3})(r - \frac{2}{3})(r - 1) \\ &= \frac{81}{80} \frac{1}{81} (3r + 2)(3r + 1)(r)(3r - 1)(3r - 2)(r - 1) \\ &= \frac{1}{80} (9r^2 - 4)(9r^2 - 1)(r^2 - r) \\ &= \frac{1}{80} (81r^4 - 45r^2 + 4)(r^2 - r) \end{aligned} \tag{5.21}$$

¹https://physandmathsolutions.com/Matrices/matrix_inverse/matrix_inverse_6x6.php

Moving to $\mathcal{N}_2(r)$, we have

$$\mathcal{N}_2(r) = \alpha(r+1)(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1)$$

which must be equal to 1 for $r = -2/3$:

$$\begin{aligned}\mathcal{N}_2(r = -2/3) &= \alpha(-\frac{2}{3}+1)(-\frac{2}{3}+\frac{1}{3})(-\frac{2}{3})(-\frac{2}{3}-\frac{1}{3})(-\frac{2}{3}-\frac{2}{3})(-\frac{2}{3}-1) \\ &= \alpha(\frac{1}{3})(-\frac{1}{3})(-\frac{2}{3})(-1)(-\frac{4}{3})(-\frac{5}{3}) \\ &= -\alpha\frac{40}{243}\end{aligned}\tag{5.22}$$

so that

$$\mathcal{N}_2(r) = -\frac{243}{40}(r+1)(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1)$$

Moving to $\mathcal{N}_3(r)$, we have

$$\mathcal{N}_3(r) = \alpha(r+1)(r+\frac{2}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1)$$

which must be equal to 1 for $r = -1/3$:

$$\begin{aligned}\mathcal{N}_3(r = -1/3) &= \alpha(-\frac{1}{3}+1)(-\frac{1}{3}+\frac{2}{3})(-\frac{1}{3})(-\frac{1}{3}-\frac{1}{3})(-\frac{1}{3}-\frac{2}{3})(-\frac{1}{3}-1) \\ &= \alpha(\frac{2}{3})(\frac{1}{3})(-\frac{1}{3})(-\frac{2}{3})(-1)(-\frac{4}{3}) \\ &= \alpha\frac{16}{243}\end{aligned}\tag{5.23}$$

so that

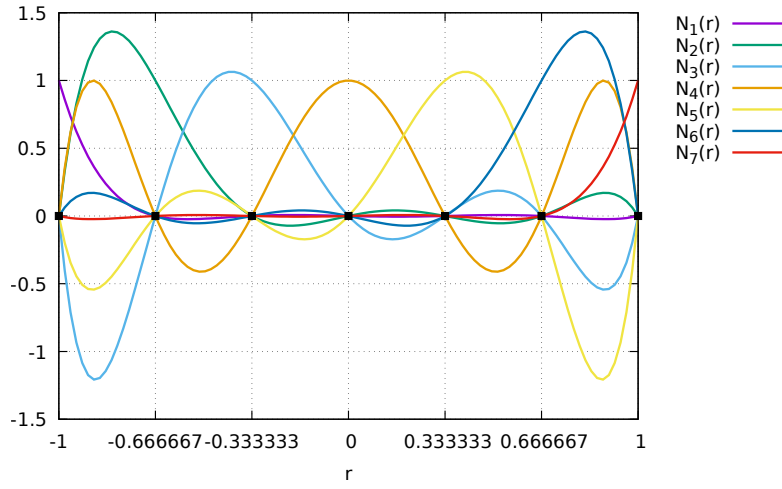
$$\mathcal{N}_3(r) = \frac{243}{16}(r+1)(r+\frac{2}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1)$$

Likewise, we arrive at the rest of the basis functions. In the end:

$$\begin{aligned}\mathcal{N}_1(r) &= \frac{81}{80}(r+\frac{2}{3})(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1) \\ \mathcal{N}_2(r) &= -\frac{243}{40}(r+1)(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1) \\ \mathcal{N}_3(r) &= \frac{243}{16}(r+1)(r+\frac{2}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})(r-1) \\ \mathcal{N}_4(r) &= -\frac{81}{4}(r+1)(r+\frac{2}{3})(r+\frac{1}{3})(r-\frac{1}{3})(r-\frac{2}{3})(r-1) \\ \mathcal{N}_5(r) &= \frac{243}{16}(r+1)(r+\frac{2}{3})(r+\frac{1}{3})(r)(r-\frac{2}{3})(r-1) \\ \mathcal{N}_6(r) &= -\frac{243}{40}(r+1)(r+\frac{2}{3})(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-1) \\ \mathcal{N}_7(r) &= \frac{81}{80}(r+1)(r+\frac{2}{3})(r+\frac{1}{3})(r)(r-\frac{1}{3})(r-\frac{2}{3})\end{aligned}\tag{5.24}$$

or

$$\begin{aligned}
\mathcal{N}_1(r) &= \frac{1}{80}(81r^6 - 81r^5 - 45r^4 + 45r^3 + 4r^2 - 4r) \\
\mathcal{N}_2(r) &= -\frac{9}{40}(27r^6 - 18r^5 - 30r^4 + 20r^3 + 3r^2 - 2r) \\
\mathcal{N}_3(r) &= \frac{9}{16}(27r^6 - 9r^5 - 39r^4 + 13r^3 + 12r^2 - 4r) \\
\mathcal{N}_4(r) &= -\frac{1}{4}(81r^6 - 126r^4 + 49r^2 - 4) \\
\mathcal{N}_5(r) &= \frac{9}{16}(27r^6 + 9r^5 - 39r^4 - 13r^3 + 12r^2 + 4r) \\
\mathcal{N}_6(r) &= -\frac{9}{40}(27r^6 + 18r^5 - 30r^4 - 20r^3 + 3r^2 + 2r) \\
\mathcal{N}_7(r) &= \frac{1}{80}(81r^6 + 81r^5 - 45r^4 - 45r^3 + 4r^2 + 4r)
\end{aligned} \tag{5.25}$$



Plot of the 7 six-order basis functions.

Using WolframAlpha², we arrive at

$$\begin{aligned}
\frac{d\mathcal{N}_1}{dr} &= \frac{1}{80}(486r^5 - 405r^4 - 180r^3 + 135r^2 + 8r - 4) \\
\frac{d\mathcal{N}_2}{dr} &= -\frac{9}{20}(81r^5 - 45r^4 - 60r^3 + 30r^2 + 3r - 1) \\
\frac{d\mathcal{N}_3}{dr} &= \frac{9}{16}(162r^5 - 45r^4 - 156r^3 + 39r^2 + 24r - 4) \\
\frac{d\mathcal{N}_4}{dr} &= \frac{1}{2}(-243r^5 + 252r^3 - 49r) \\
\frac{d\mathcal{N}_5}{dr} &= \frac{9}{16}(162r^5 + 45r^4 - 156r^3 - 39r^2 + 24r + 4) \\
\frac{d\mathcal{N}_6}{dr} &= -\frac{9}{20}(81r^5 + 45r^4 - 60r^3 - 30r^2 + 3r + 1) \\
\frac{d\mathcal{N}_7}{dr} &= \frac{1}{80}(486r^5 + 405r^4 - 180r^3 - 135r^2 + 8r + 4)
\end{aligned} \tag{5.26}$$

These functions are used in [STONE](#) ??.

²<https://www.wolframalpha.com/>

5.2.7 A generic approach to 1D basis functions

In order to define basis functions of order n each element must have $n + 1$ nodes. The $i - th$ basis function for an $n - th$ order approximation is given by:

$$\mathcal{N}_i(r) = \frac{\prod_{j=0, j \neq i}^n (r - r_j)}{\prod_{j=0, j \neq i}^n (r_i - r_j)}$$

Let us see in practice how this works and start with $n = 2$ (i.e. Q_2 basis functions). In the reference element we have $r_0 = -1$, $r_1 = 0$ and $r_2 = +1$, so that

$$\begin{aligned} \mathcal{N}_0 &= \frac{(r - r_1)(r - r_2)}{(r_0 - r_1)(r_0 - r_2)} \\ &= \frac{(r - 0)(r - 1)}{(-1 - 0)(-1 - 1)} \\ &= \frac{1}{2}r(r - 1) \\ \mathcal{N}_1 &= \frac{(r - r_0)(r - r_2)}{(r_1 - r_0)(r_1 - r_2)} \\ &= \frac{(r + 1)(r - 1)}{(0 + 1)(0 - 1)} \\ &= 1 - r^2 \\ \mathcal{N}_2 &= \frac{(r - r_0)(r - r_1)}{(r_2 - r_0)(r_2 - r_1)} \\ &= \frac{(r + 1)(r - 0)}{(1 + 1)(1 - 0)} \\ &= \frac{1}{2}r(r + 1) \end{aligned} \tag{5.27}$$

These are the basis functions obtained in Section 5.2.2.

Note that in practice it is often rarely desirable to use much higher than quadratic basis functions because higher order ones have too much oscillation.

What about derivatives?

Remark:

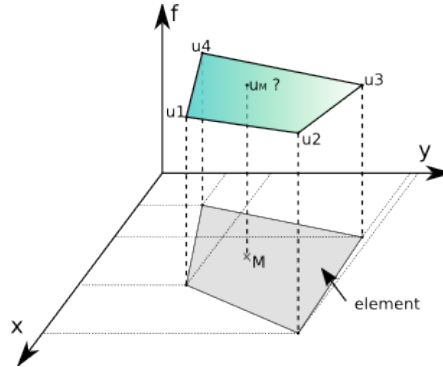
According to Guermond³: “While the choice of equidistant nodes appears somewhat natural, it is appropriate only when working with low-degree polynomials. The main difficulty comes from the oscillatory nature of the Lagrange polynomials as the number of interpolation nodes grows. This phenomenon is often referred to as the Runge phenomenon [359] (see also Meray [313])”

³https://people.tamu.edu/~guermond/M661_FALL_2015/chap2.pdf

5.3 Elements and basis functions in 2D

elements2D.tex

Let us for a moment consider a single quadrilateral element in the xy -plane, as shown on the following figure:



Let us assume that we know the values of a given field u at the four vertices. For a given point M inside the element in the plane, what is the value of the field u at this point? It makes sense to postulate that u_M will be given by

$$u_M = \phi(u_1, u_2, u_3, u_4, x_M, y_M)$$

where ϕ is a function to be determined. Although ϕ is not unique, we can decide to express the value u_M as a weighed sum of the values at the vertices u_i . One option could be to assign all four vertices the same weight, say $1/4$ so that $u_M = (u_1 + u_2 + u_3 + u_4)/4$, i.e. u_M is simply given by the arithmetic mean of the vertices values. If the function $u(x, y)$ is such that it is a constant function, say $u(x, y) = C$, then $u_M = (u_1 + u_2 + u_3 + u_4)/4 = (C + C + C + C)/4 = C$ and the result is exact. However, for any other function u the value u_M will not be as accurate. Also, this approach suffers from a major drawback as it does not use the location of point M inside the element. For instance, when $(x_M, y_M) \rightarrow (x_2, y_2)$ we expect $u_M \rightarrow u_2$ but u_M would remain equal to $(u_1 + u_2 + u_3 + u_4)/4$.

In light of this, we could now assume that the weights would depend on the position of M in a continuous fashion:

$$u(x_M, y_M) = \sum_{i=1}^4 \mathcal{N}_i(x_M, y_M) u_i = \mathcal{N}_1(x_M, y_M)u_1 + \mathcal{N}_2(x_M, y_M)u_2 + \mathcal{N}_3(x_M, y_M)u_3 + \mathcal{N}_4(x_M, y_M)u_4 \quad (5.28)$$

where the \mathcal{N}_i are continuous (and also "well behaved") functions which have the property:

$$\mathcal{N}_i(x_j, y_j) = \delta_{ij}$$

or, in other words for example:

$$\begin{aligned} \mathcal{N}_3(x_1, y_1) &= 0 \\ \mathcal{N}_3(x_2, y_2) &= 0 \\ \mathcal{N}_3(x_3, y_3) &= 1 \\ \mathcal{N}_3(x_4, y_4) &= 0 \end{aligned} \quad (5.29)$$

The functions \mathcal{N}_i are commonly called basis functions.

Omitting the M subscripts (yet stil assuming the point being inside the element), the velocity components u and v for a point inside the element are given by:

$$u^h(x, y) = \sum_{i=1}^4 \mathcal{N}_i(x, y) u_i \quad (5.30)$$

$$v^h(x, y) = \sum_{i=1}^4 \mathcal{N}_i(x, y) v_i \quad (5.31)$$

where we have added the superscript h to denote that it is an approximation of the functions of this element of diameter h (by diameter we mean here a representative scalar value of the dimension of the element).

One can now easily compute velocity gradients (and therefore the strain rate tensor) since we have assumed the basis functions to be "well behaved" (in this case first-order differentiable):

$$\dot{\varepsilon}_{xx}^h(x, y) = \frac{\partial u^h}{\partial x} = \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} u_i \quad (5.32)$$

$$\dot{\varepsilon}_{yy}^h(x, y) = \frac{\partial v^h}{\partial y} = \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} v_i \quad (5.33)$$

$$\dot{\varepsilon}_{xy}^h(x, y) = \frac{1}{2} \left(\frac{\partial u^h}{\partial y} + \frac{\partial v^h}{\partial x} \right) = \frac{1}{2} \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} u_i + \frac{1}{2} \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} v_i \quad (5.34)$$

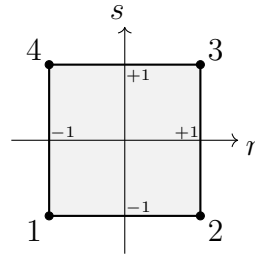
How we actually obtain the exact form of the basis functions \mathcal{N}_i is explained in the coming sections.

5.3.1 Bilinear basis functions in 2D (Q_1)

basis_Q1_2D.tex

In this section, we consider for simplicity an element which is a square defined by $-1 < r < 1$, $-1 < s < 1$ in the Cartesian coordinates system (r, s) ⁴:

(tikz_q12d.tex)



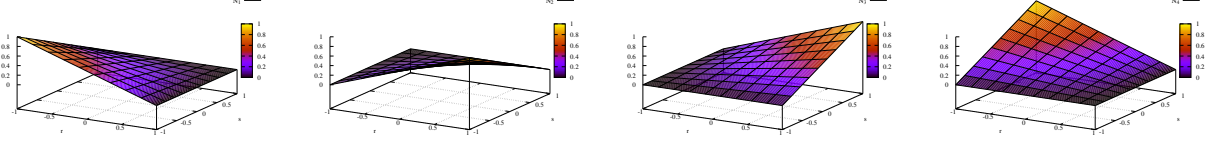
Note the counter-clockwise numbering⁵. This element is commonly called the reference element. How we go from the (x, y) coordinate system to the (r, s) once and vice versa will be dealt with later on. The basis functions in the above reference element in the reduced coordinates system (r, s) are given by:

⁴There is a reason to choose r and s as coordinates and not x and y as we will see later.

⁵Note that in many of the python codes which are part of this project the numbering starts at 0.

$$\begin{aligned}
\mathcal{N}_1(r, s) &= 0.25(1-r)(1-s) \\
\mathcal{N}_2(r, s) &= 0.25(1+r)(1-s) \\
\mathcal{N}_3(r, s) &= 0.25(1+r)(1+s) \\
\mathcal{N}_4(r, s) &= 0.25(1-r)(1+s)
\end{aligned} \tag{5.35}$$

These basis functions are the product of the linear basis functions of Section 5.2.1 in the r direction and the s direction.



Surface representation of the basis functions on the reference element. in images/basis.Q1.2D/

The partial derivatives of these functions with respect to r and s automatically follow:

$$\begin{aligned}
\frac{\partial \mathcal{N}_1}{\partial r}(r, s) &= -0.25(1-s) & \frac{\partial \mathcal{N}_1}{\partial s}(r, s) &= -0.25(1-r) \\
\frac{\partial \mathcal{N}_2}{\partial r}(r, s) &= +0.25(1-s) & \frac{\partial \mathcal{N}_2}{\partial s}(r, s) &= -0.25(1+r) \\
\frac{\partial \mathcal{N}_3}{\partial r}(r, s) &= +0.25(1+s) & \frac{\partial \mathcal{N}_3}{\partial s}(r, s) &= +0.25(1+r) \\
\frac{\partial \mathcal{N}_4}{\partial r}(r, s) &= -0.25(1+s) & \frac{\partial \mathcal{N}_4}{\partial s}(r, s) &= +0.25(1-r)
\end{aligned}$$

Let us go back to Eq. (5.31) and let us assume that the function $v(r, s) = C$ so that $v_i = C$ for $i = 1, 2, 3, 4$. It then follows that

$$v^h(r, s) = \sum_{i=1}^4 \mathcal{N}_i(r, s) v_i = C \sum_{i=1}^4 \mathcal{N}_i(r, s) = C[\mathcal{N}_1(r, s) + \mathcal{N}_2(r, s) + \mathcal{N}_3(r, s) + \mathcal{N}_4(r, s)] = C$$

This is a very important property: if the v function used to assign values at the vertices is constant, then the value of v^h *anywhere* in the element is exactly C . If we now turn to the derivatives of v with respect to r and s :

$$\frac{\partial v^h}{\partial r}(r, s) = \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial r}(r, s) v_i = C \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial r}(r, s) = C[-0.25(1-s) + 0.25(1-s) + 0.25(1+s) - 0.25(1+s)] = 0$$

$$\frac{\partial v^h}{\partial s}(r, s) = \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial s}(r, s) v_i = C \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial s}(r, s) = C[-0.25(1-r) - 0.25(1+r) + 0.25(1+r) + 0.25(1-r)] = 0$$

We reassuringly find that the derivative of a constant field anywhere in the element is exactly zero.

If we now choose $v(r, s) = ar + bs$ with a and b two constant scalars, we find:

$$\begin{aligned}
v^h(r, s) &= \sum_{i=1}^4 \mathcal{N}_i(r, s) v_i \\
&= \sum_{i=1}^4 \mathcal{N}_i(r, s)(ar_i + bs_i) \\
&= a \sum_{i=1}^4 \mathcal{N}_i(r, s)r_i + b \sum_{i=1}^4 \mathcal{N}_i(r, s)s_i \\
&= a \left[\frac{1}{4}(1-r)(1-s)(-1) + \frac{1}{4}(1+r)(1-s)(+1) + \frac{1}{4}(1+r)(1+s)(+1) + \frac{1}{4}(1-r)(1+s)(-1) \right] \\
&\quad + b \left[\frac{1}{4}(1-r)(1-s)(-1) + \frac{1}{4}(1+r)(1-s)(-1) + \frac{1}{4}(1+r)(1+s)(+1) + \frac{1}{4}(1-r)(1+s)(+1) \right] \\
&= \frac{a}{4} [-(1-r)(1-s) + (1+r)(1-s) + (1+r)(1+s) - (1-r)(1+s)] \\
&\quad + \frac{b}{4} [-(1-r)(1-s) - (1+r)(1-s) + (1+r)(1+s) + (1-r)(1+s)] \\
&= ar + bs
\end{aligned} \tag{5.36}$$

This set of bilinear basis functions is therefore capable of exactly representing a bilinear field. The derivatives are:

$$\begin{aligned}
\frac{\partial v^h}{\partial r}(r, s) &= \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial r}(r, s) v_i \\
&= a \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial r}(r, s)r_i + b \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial r}(r, s)s_i \\
&= a \left[-\frac{1}{4}(1-s)(-1) + \frac{1}{4}(1-s)(+1) + \frac{1}{4}(1+s)(+1) - \frac{1}{4}(1+s)(-1) \right] \\
&\quad + b \left[-\frac{1}{4}(1-s)(-1) + \frac{1}{4}(1-s)(-1) + \frac{1}{4}(1+s)(+1) - \frac{1}{4}(1+s)(+1) \right] \\
&= \frac{a}{4} [(1-s) + (1-s) + (1+s) + (1+s)] \\
&\quad + \frac{b}{4} [(1-s) - (1-s) + (1+s) - (1+s)] \\
&= a
\end{aligned} \tag{5.37}$$

Here again, we find that the derivative of the bilinear field inside the element is exact: $\frac{\partial v^h}{\partial r} = \frac{\partial v}{\partial r}$.

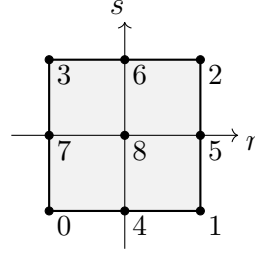
However, following the same methodology as above, one can easily prove that this is no more true for polynomials of degree strictly higher than 1. This fact has serious consequences: if the solution to the problem at hand is for instance a parabola, the Q_1 basis functions cannot represent the solution properly, but only by approximating the parabola in each element by a line. As we will see later, Q_2 basis functions can remedy this problem by containing quadratic terms.

Remark. The Q_1 basis functions are first-order polynomials. We have seen that they can be used to compute gradients. However they cannot be used to compute 2nd-order derivatives since their 2nd-order derivative is identically zero.

5.3.2 Biquadratic basis functions in 2D (Q_2)

This element is part of the so-called Lagrange family [1047]. Inside an element the local numbering of the nodes is as follows⁶:

(tikz-q22d.tex)



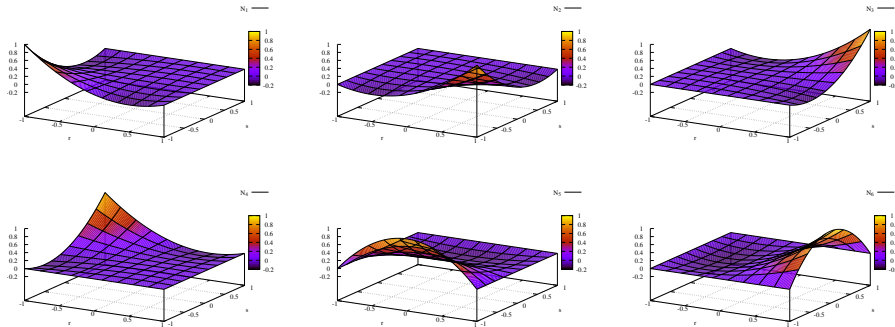
Note that this numbering is also employed in Li [779, p56]. The polynomial representation of the function ϕ over this element is then taken to be biquadratic:

$$\phi^h(r, s) = a + br + cs + drs + er^2 + fs^2 + gr^2s + hrs^2 + ir^2s^2 = \sum_{i=0}^8 \mathcal{N}_i(r, s)\phi_i$$

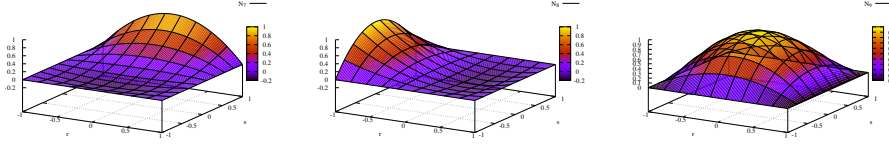
and one can show that the basis functions are:

$$\begin{aligned} \mathcal{N}_0(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s-1) \\ \mathcal{N}_1(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s-1) \\ \mathcal{N}_2(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s+1) \\ \mathcal{N}_3(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s+1) \\ \mathcal{N}_4(r, s) &= (1-r^2)\frac{1}{2}s(s-1) \\ \mathcal{N}_5(r, s) &= \frac{1}{2}r(r+1)(1-s^2) \\ \mathcal{N}_6(r, s) &= (1-r^2)\frac{1}{2}s(s+1) \\ \mathcal{N}_7(r, s) &= \frac{1}{2}r(r-1)(1-s^2) \\ \mathcal{N}_8(r, s) &= (1-r^2)(1-s^2) \end{aligned}$$

Note that we have $\mathcal{N}_i(r_j, s_j) = \delta_{ij}$ and then obviously $\mathcal{N}_i(r_i, s_i) = 1$.



⁶I have adopted here a numbering scheme starting at zero! Also, it is a numbering among many other possible choices!



Surface representation of the basis functions on the reference element. in images/basis_Q2_2D/

Their derivatives are given by:

$$\begin{aligned}
 \frac{\partial \mathcal{N}_0}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s-1) & \frac{\partial \mathcal{N}_0}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s-1) \\
 \frac{\partial \mathcal{N}_1}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s-1) & \frac{\partial \mathcal{N}_1}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s-1) \\
 \frac{\partial \mathcal{N}_2}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s+1) & \frac{\partial \mathcal{N}_2}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s+1) \\
 \frac{\partial \mathcal{N}_3}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s+1) & \frac{\partial \mathcal{N}_3}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s+1) \\
 \frac{\partial \mathcal{N}_4}{\partial r} &= (-2r)\frac{1}{2}s(s-1) & \frac{\partial \mathcal{N}_4}{\partial s} &= (1-r^2)\frac{1}{2}(2s-1) \\
 \frac{\partial \mathcal{N}_5}{\partial r} &= \frac{1}{2}(2r+1)(1-s^2) & \frac{\partial \mathcal{N}_5}{\partial s} &= \frac{1}{2}r(r+1)(-2s) \\
 \frac{\partial \mathcal{N}_6}{\partial r} &= (-2r)\frac{1}{2}s(s+1) & \frac{\partial \mathcal{N}_6}{\partial s} &= (1-r^2)\frac{1}{2}(2s+1) \\
 \frac{\partial \mathcal{N}_7}{\partial r} &= \frac{1}{2}(2r-1)(1-s^2) & \frac{\partial \mathcal{N}_7}{\partial s} &= \frac{1}{2}r(r-1)(-2s) \\
 \frac{\partial \mathcal{N}_8}{\partial r} &= (-2r)(1-s^2) & \frac{\partial \mathcal{N}_8}{\partial s} &= (1-r^2)(-2s)
 \end{aligned}$$

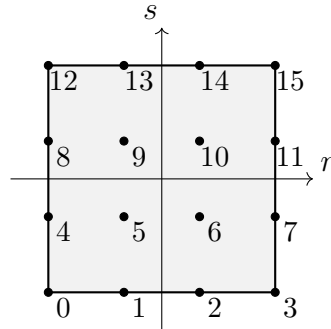
These basis functions are used for example in [STONE 18](#).

5.3.3 Bicubic basis functions in 2D (Q_3)

basis_Q3_2D.tex

Inside an element a possible local numbering of the nodes is as follows:

(tikz-q32d.tex)

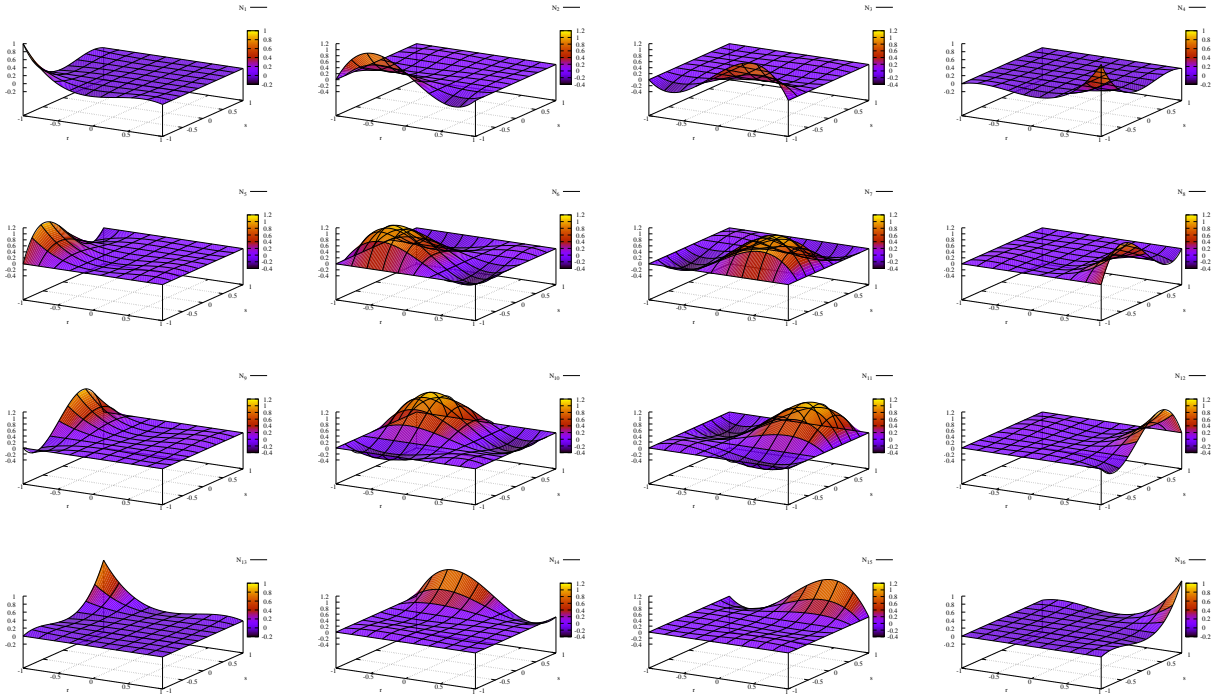


As shown in Section 5.2.3 the 1D cubic basis functions are given by:

$$\begin{aligned}
 \mathcal{N}_1(r) &= (-1 + r + 9r^2 - 9r^3)/16 & \mathcal{N}_1(s) &= (-1 + s + 9s^2 - 9s^3)/16 \\
 \mathcal{N}_2(r) &= (+9 - 27r - 9r^2 + 27r^3)/16 & \mathcal{N}_2(s) &= (+9 - 27s - 9s^2 + 27s^3)/16 \\
 \mathcal{N}_3(r) &= (+9 + 27r - 9r^2 - 27r^3)/16 & \mathcal{N}_3(s) &= (+9 + 27s - 9s^2 - 27s^3)/16 \\
 \mathcal{N}_4(r) &= (-1 - r + 9r^2 + 9r^3)/16 & \mathcal{N}_4(s) &= (-1 - s + 9s^2 + 9s^3)/16
 \end{aligned}$$

and the resulting 2D basis functions are simply the tensor product of the above 1D ones:

$$\begin{aligned}
\mathcal{N}_{01}(r, s) &= \mathcal{N}_1(r)\mathcal{N}_1(s) = (-1 + r + 9r^2 - 9r^3)/16 \cdot (-1 + s + 9s^2 - 9s^3)/16 \\
\mathcal{N}_{02}(r, s) &= \mathcal{N}_2(r)\mathcal{N}_1(s) = (+9 - 27r - 9r^2 + 27r^3)/16 \cdot (-1 + s + 9s^2 - 9s^3)/16 \\
\mathcal{N}_{03}(r, s) &= \mathcal{N}_3(r)\mathcal{N}_1(s) = (+9 + 27r - 9r^2 - 27r^3)/16 \cdot (-1 + s + 9s^2 - 9s^3)/16 \\
\mathcal{N}_{04}(r, s) &= \mathcal{N}_4(r)\mathcal{N}_1(s) = (-1 - r + 9r^2 + 9r^3)/16 \cdot (-1 + s + 9s^2 - 9s^3)/16 \\
\mathcal{N}_{05}(r, s) &= \mathcal{N}_1(r)\mathcal{N}_2(s) = (-1 + r + 9r^2 - 9r^3)/16 \cdot (9 - 27s - 9s^2 + 27s^3)/16 \\
\mathcal{N}_{06}(r, s) &= \mathcal{N}_2(r)\mathcal{N}_2(s) = (+9 - 27r - 9r^2 + 27r^3)/16 \cdot (9 - 27s - 9s^2 + 27s^3)/16 \\
\mathcal{N}_{07}(r, s) &= \mathcal{N}_3(r)\mathcal{N}_2(s) = (+9 + 27r - 9r^2 - 27r^3)/16 \cdot (9 - 27s - 9s^2 + 27s^3)/16 \\
\mathcal{N}_{08}(r, s) &= \mathcal{N}_4(r)\mathcal{N}_2(s) = (-1 - r + 9r^2 + 9r^3)/16 \cdot (9 - 27s - 9s^2 + 27s^3)/16 \\
\mathcal{N}_{09}(r, s) &= \mathcal{N}_1(r)\mathcal{N}_3(s) = (-1 + r + 9r^2 - 9r^3)/16 \cdot (9 + 27s - 9s^2 - 27s^3)/16 \\
\mathcal{N}_{10}(r, s) &= \mathcal{N}_2(r)\mathcal{N}_3(s) = (+9 - 27r - 9r^2 + 27r^3)/16 \cdot (9 + 27s - 9s^2 - 27s^3)/16 \\
\mathcal{N}_{11}(r, s) &= \mathcal{N}_3(r)\mathcal{N}_3(s) = (+9 + 27r - 9r^2 - 27r^3)/16 \cdot (9 + 27s - 9s^2 - 27s^3)/16 \\
\mathcal{N}_{12}(r, s) &= \mathcal{N}_4(r)\mathcal{N}_3(s) = (-1 - r + 9r^2 + 9r^3)/16 \cdot (9 + 27s - 9s^2 - 27s^3)/16 \\
\mathcal{N}_{13}(r, s) &= \mathcal{N}_1(r)\mathcal{N}_4(s) = (-1 + r + 9r^2 - 9r^3)/16 \cdot (-1 - s + 9s^2 + 9s^3)/16 \\
\mathcal{N}_{14}(r, s) &= \mathcal{N}_2(r)\mathcal{N}_4(s) = (+9 - 27r - 9r^2 + 27r^3)/16 \cdot (-1 - s + 9s^2 + 9s^3)/16 \\
\mathcal{N}_{15}(r, s) &= \mathcal{N}_3(r)\mathcal{N}_4(s) = (+9 + 27r - 9r^2 - 27r^3)/16 \cdot (-1 - s + 9s^2 + 9s^3)/16 \\
\mathcal{N}_{16}(r, s) &= \mathcal{N}_4(r)\mathcal{N}_4(s) = (-1 - r + 9r^2 + 9r^3)/16 \cdot (-1 - s + 9s^2 + 9s^3)/16 \quad (5.38)
\end{aligned}$$



Surface representation of the basis functions on the reference element. in `images/basis.Q3.2D/`

The derivatives are trivial to obtain from the derivatives of the 1D basis functions, e.g.

$$\frac{\partial \mathcal{N}_{13}}{\partial r} = \frac{\partial \mathcal{N}_1}{\partial r} \mathcal{N}_3(s)$$

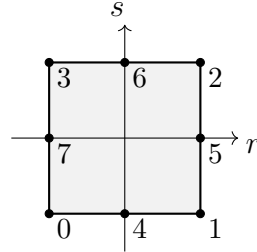
These basis functions are used in [STONE 19](#).

5.3.4 Eight node serendipity basis functions in 2D ($Q_2^{(8)}$)

basis_Q28_2D.tex

The serendipity elements are those rectangular elements which have no interior nodes (See for example Reddy [1051, p65]). Inside an element a possible local numbering of the nodes is as follows:

(tikz_serendipity2d.tex)



The main difference with the Q_2 element resides in the fact that there is no node in the middle of the element. The polynomial representation of the function ϕ over the element is then

$$\phi_h(r, s) = a + br + cs + drs + er^2 + fs^2 + gr^2s + hrs^2$$

Note that absence of the r^2s^2 term which was previously associated to the center node. We find that

$$\mathcal{N}_0(r, s) = \frac{1}{4}(1-r)(1-s)(-r-s-1) \quad (5.39)$$

$$\mathcal{N}_1(r, s) = \frac{1}{4}(1+r)(1-s)(r-s-1) \quad (5.40)$$

$$\mathcal{N}_2(r, s) = \frac{1}{4}(1+r)(1+s)(r+s-1) \quad (5.41)$$

$$\mathcal{N}_3(r, s) = \frac{1}{4}(1-r)(1+s)(-r+s-1) \quad (5.42)$$

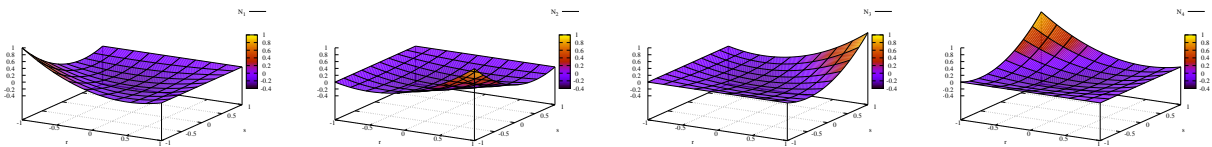
$$\mathcal{N}_4(r, s) = \frac{1}{2}(1-r^2)(1-s) \quad (5.43)$$

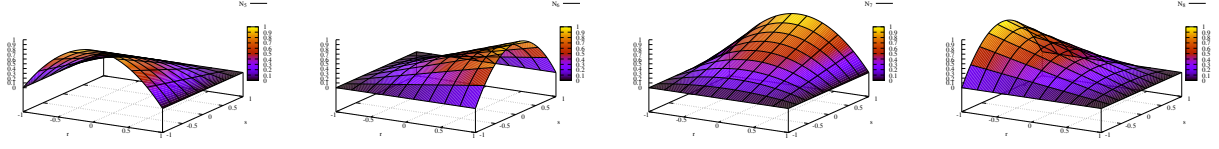
$$\mathcal{N}_5(r, s) = \frac{1}{2}(1+r)(1-s^2) \quad (5.44)$$

$$\mathcal{N}_6(r, s) = \frac{1}{2}(1-r^2)(1+s) \quad (5.45)$$

$$\mathcal{N}_7(r, s) = \frac{1}{2}(1-r)(1-s^2) \quad (5.46)$$

The basis functions at the mid side nodes are products of a second order polynomial parallel to side and a linear function perpendicular to the side while basis functions for corner nodes are modifications of the bilinear quadrilateral element.





Surface representation of the basis functions on the reference element. in `images/basis_Q28_2D/`

The first-order derivatives are given by:

$$\frac{\partial \mathcal{N}_0}{\partial r}(r, s) = -\frac{1}{4}(s-1)(2r+s) \quad (5.47)$$

$$\frac{\partial \mathcal{N}_1}{\partial r}(r, s) = -\frac{1}{4}(s-1)(2r-s) \quad (5.48)$$

$$\frac{\partial \mathcal{N}_2}{\partial r}(r, s) = \frac{1}{4}(s+1)(2r+s) \quad (5.49)$$

$$\frac{\partial \mathcal{N}_3}{\partial r}(r, s) = \frac{1}{4}(s+1)(2r-s) \quad (5.50)$$

$$\frac{\partial \mathcal{N}_4}{\partial r}(r, s) = r(s-1) \quad (5.51)$$

$$\frac{\partial \mathcal{N}_5}{\partial r}(r, s) = \frac{1}{2}(1-s^2) \quad (5.52)$$

$$\frac{\partial \mathcal{N}_6}{\partial r}(r, s) = -r(s+1) \quad (5.53)$$

$$\frac{\partial \mathcal{N}_7}{\partial r}(r, s) = -\frac{1}{2}(1-s^2) \quad (5.54)$$

$$\frac{\partial \mathcal{N}_0}{\partial s}(r, s) = -\frac{1}{4}(r-1)(r+2s) \quad (5.55)$$

$$\frac{\partial \mathcal{N}_1}{\partial s}(r, s) = -\frac{1}{4}(r+1)(r-2s) \quad (5.56)$$

$$\frac{\partial \mathcal{N}_2}{\partial s}(r, s) = \frac{1}{4}(r+1)(r+2s) \quad (5.57)$$

$$\frac{\partial \mathcal{N}_3}{\partial s}(r, s) = \frac{1}{4}(r-1)(r-2s) \quad (5.58)$$

$$\frac{\partial \mathcal{N}_4}{\partial s}(r, s) = -\frac{1}{2}(1-r^2) \quad (5.59)$$

$$\frac{\partial \mathcal{N}_5}{\partial s}(r, s) = -(r+1)s \quad (5.60)$$

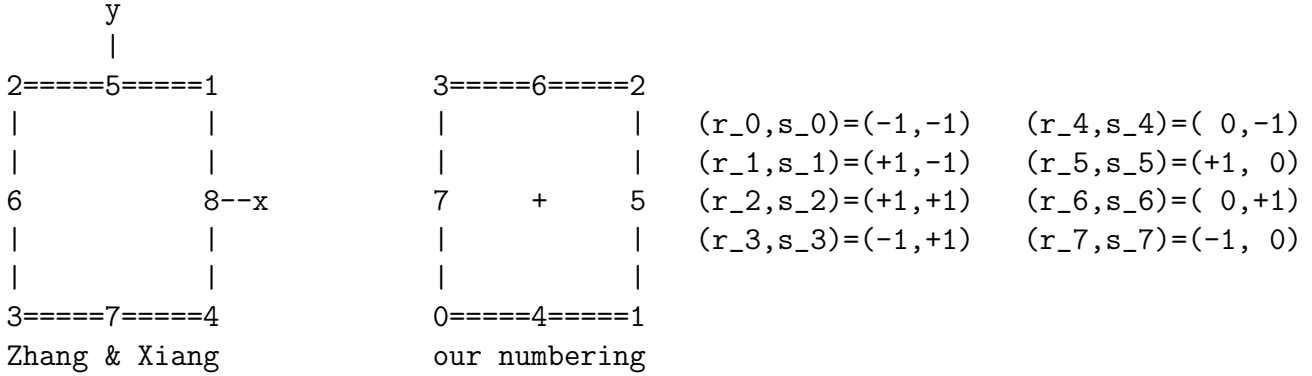
$$\frac{\partial \mathcal{N}_6}{\partial s}(r, s) = \frac{1}{2}(1-r^2) \quad (5.61)$$

$$\frac{\partial \mathcal{N}_7}{\partial s}(r, s) = (r-1)s \quad (5.62)$$

These basis functions are used in [STONE 52](#).

5.3.5 Eight node serendipity basis functions in 2D ($QH8 - C1$)

This element is proposed in Zhang & Xiang (2020) [1404]. Two remarks must be made: 1) Eq. (29) of their publication which is the definition of the basis functions contains an error⁷. 2) The authors use a rather uncommon and annoying rotated numbering:



For each element they define (their numbering):

$$\begin{aligned}
 A &= \frac{1}{2}[(x_1 - x_3)(y_2 - y_4) - (x_2 - x_4)(y_1 - y_3)] \\
 m_x &= (x_1 - x_4)(y_2 - y_3) - (x_2 - x_3)(y_1 - y_4) \\
 m_y &= (x_3 - x_4)(y_1 - y_2) - (x_1 - x_2)(y_3 - y_4)
 \end{aligned}$$

Note that A is the area of the element, and that in the case when the element is a rectangle then $m_x = m_y = 0$.

$$\begin{aligned}
 \mathcal{N}_1(r, s) &= n_1(r, s) + (m_x^2 - m_x m_y + m_y^2) \frac{E(r, s)}{D} \\
 \mathcal{N}_2(r, s) &= n_2(r, s) + (m_x^2 + m_x m_y + m_y^2) \frac{E(r, s)}{D} \\
 \mathcal{N}_3(r, s) &= n_3(r, s) + (m_x^2 - m_x m_y + m_y^2) \frac{E(r, s)}{D} \\
 \mathcal{N}_4(r, s) &= n_4(r, s) + (m_x^2 + m_x m_y + m_y^2) \frac{E(r, s)}{D} \\
 \mathcal{N}_5(r, s) &= n_5(r, s) - m_x(2Am_x + m_y^2) \frac{E(r, s)}{AD} \\
 \mathcal{N}_6(r, s) &= n_6(r, s) - m_y(2Am_y + m_x^2) \frac{E(r, s)}{AD} \\
 \mathcal{N}_7(r, s) &= n_7(r, s) + m_x(-2Am_x + m_y^2) \frac{E(r, s)}{AD} \\
 \mathcal{N}_8(r, s) &= n_8(r, s) + m_y(-2Am_y + m_x^2) \frac{E(r, s)}{AD}
 \end{aligned}$$

with

$$E(r, s) = (1 - r^2)(1 - s^2) \quad D = 4(4A^2 + m_x^2 + m_y^2)$$

and where the n_i functions are the basis functions of the 'regular' 8-node element (see Section 5.3.4).

This is implemented in [STONE](#) 52.

not finished. SHOW CONSISTENCY !! like in paper email sent to author about mistake.

⁷Answer from the author: "N5 to N8 is missing an A in the denominator and the calculation program does not have this problem"

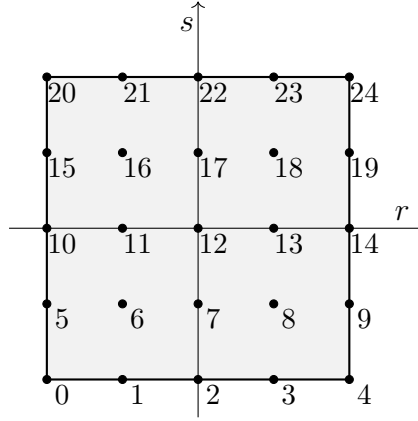
Let us verify consistency:

$$\begin{aligned}
\sum_{i=1}^8 \mathcal{N}_i(r, s) &= \underbrace{\sum_{i=1}^8 n_i(r, s)}_{=0} + \frac{E(r, s)}{D} \left[(m_x^2 - m_x m_y + m_y^2) + (m_x^2 + m_x m_y + m_y^2) + (m_x^2 - m_x m_y + m_y^2) + (m_x^2 + m_x m_y + m_y^2) \right] \\
&\quad - m_x(2Am_x + m_y^2) \frac{1}{A} - m_y(2Am_y + m_x^2) \frac{1}{A} + m_x(-2Am_x + m_y^2) \frac{1}{A} + m_y(-2Am_y + m_x^2) \frac{1}{A} \\
&= \frac{E(r, s)}{D} \left[(4m_x^2 + 4m_y^2) + \frac{1}{A}(-4Am_x^2 - 4Am_y^2) \right] \\
&= 0
\end{aligned}$$

5.3.6 Biquartic basis functions in 2D (Q_4)

Inside an element the local numbering of the nodes is as follows:

(tikz-q42d.tex)

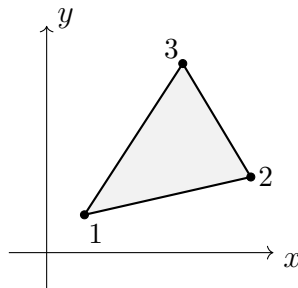


5.3.7 Linear basis functions for triangles in 2D (P_1)

basis_P1_2D.tex

Here we do not start from a reference element but consider instead a generic triangle:

(tikz_P1.tex)



This is the simplest 2D element, which is also called linear triangular element. Velocities (or displacements) (u^h, v^h) in the element are interpolated from nodal velocities (u_i, v_i) using basis functions

\mathcal{N}_i as follows,

$$\vec{v}^h = \begin{pmatrix} u^h(x, y) \\ v^h(x, y) \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^3 \mathcal{N}_i(x, y) u_i \\ \sum_{i=1}^3 \mathcal{N}_i(x, y) v_i \end{pmatrix} = \begin{pmatrix} \mathcal{N}_1(x, y) & 0 & \mathcal{N}_2(x, y) & 0 & \mathcal{N}_3(x, y) & 0 \\ 0 & \mathcal{N}_1(x, y) & 0 & \mathcal{N}_2(x, y) & 0 & \mathcal{N}_3(x, y) \end{pmatrix}$$

or simply

For this element, we have three nodes at the vertices of the triangle, which are numbered around the element in the counterclockwise direction. Each node has two degrees of freedom (i.e. it can move in the x and y directions). The velocities u^h and v^h are assumed to be linear functions within the element, that is,

$$\begin{aligned} u^h(x, y) &= b_1 + b_2x + b_3y \\ v^h(x, y) &= b_4 + b_5x + b_6y \end{aligned} \tag{5.65}$$

where b_i are constants to be determined and which depend on the triangle shape. Note that the strain rate components are then given by

$$\begin{aligned} \dot{\epsilon}_{xx}(\vec{v}) &= b_2 \\ \dot{\epsilon}_{yy}(\vec{v}) &= b_6 \\ \dot{\epsilon}_{xy}(\vec{v}) &= (b_3 + b_5)/2 \end{aligned}$$

and are constant throughout the element.

The velocities should satisfy the following six equations (when it is evaluated at a node we should recover the nodal velocity):

$$\begin{aligned} u_1 &= u^h(x_1, y_1) = b_1 + b_2x_1 + b_3y_1 \\ u_2 &= u^h(x_2, y_2) = b_1 + b_2x_2 + b_3y_2 \\ u_3 &= u^h(x_3, y_3) = b_1 + b_2x_3 + b_3y_3 \\ v_1 &= v^h(x_1, y_1) = b_4 + b_5x_1 + b_6y_1 \\ v_2 &= v^h(x_2, y_2) = b_4 + b_5x_2 + b_6y_2 \\ v_3 &= v^h(x_3, y_3) = b_4 + b_5x_3 + b_6y_3 \end{aligned}$$

Let us focus on the three equations with the u component of the velocity. These can be re-written:

$$\begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ 1 & x_3 & y_3 \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

In order to obtain b_1, b_2, b_3 we need to solve this system, or simply to compute the inverse of the 3×3 \mathbf{M} matrix, as explained in Appendix D.0.2. We define $D = \det(\mathbf{M})$ and we get

$$\begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = \frac{1}{D} \tilde{\mathbf{M}} \cdot \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix}$$

The matrix $\tilde{\mathbf{M}}$ is given by:

$$\tilde{\mathbf{M}} = \begin{pmatrix} x_2y_3 - x_3y_2 & x_3y_1 - x_1y_3 & x_1y_2 - x_2y_1 \\ y_2 - y_3 & y_3 - y_1 & y_1 - y_2 \\ x_3 - x_2 & x_1 - x_3 & x_2 - x_1 \end{pmatrix}$$

so that

$$\begin{aligned}
b_1 &= \frac{1}{D}[(x_2y_3 - x_3y_2)u_1 + (x_3y_1 - x_1y_3)u_2 + (x_1y_2 - x_2y_1)u_3] \\
b_2 &= \frac{1}{D}[(y_2 - y_3)u_1 + (y_3 - y_1)u_2 + (y_1 - y_2)u_3] \\
b_3 &= \frac{1}{D}[(x_3 - x_2)u_1 + (x_1 - x_3)u_2 + (x_2 - x_1)u_3]
\end{aligned} \tag{5.66}$$

We then have

$$\begin{aligned}
u^h(x, y) &= b_1 + b_2x + b_3y \\
&= \frac{1}{D}[(x_2y_3 - x_3y_2)u_1 + (x_3y_1 - x_1y_3)u_2 + (x_1y_2 - x_2y_1)u_3] \\
&\quad + \frac{1}{D}[(y_2 - y_3)u_1 + (y_3 - y_1)u_2 + (y_1 - y_2)u_3]x \\
&\quad + \frac{1}{D}[(x_3 - x_2)u_1 + (x_1 - x_3)u_2 + (x_2 - x_1)u_3]y \\
&= \frac{1}{D}[(x_2y_3 - x_3y_2) + (y_2 - y_3)x + (x_3 - x_2)y]u_1 \\
&\quad + \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x + (x_1 - x_3)y]u_2 \\
&\quad + \frac{1}{D}[(x_1y_2 - x_2y_1) + (y_1 - y_2)x + (x_2 - x_1)y]u_3 \\
&= \mathcal{N}_1(x, y)u_1 + \mathcal{N}_2(x, y)u_2 + \mathcal{N}_3(x, y)u_3
\end{aligned} \tag{5.67}$$

with the linear basis functions are given by:

$$\begin{aligned}
\mathcal{N}_1(x, y) &= \frac{1}{D}[(x_2y_3 - x_3y_2) + (y_2 - y_3)x + (x_3 - x_2)y] \\
\mathcal{N}_2(x, y) &= \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x + (x_1 - x_3)y] \\
\mathcal{N}_3(x, y) &= \frac{1}{D}[(x_1y_2 - x_2y_1) + (y_1 - y_2)x + (x_2 - x_1)y]
\end{aligned}$$

We can then easily verify that for example

$$\mathcal{N}_2(x_1, y_1) = \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x_1 + (x_1 - x_3)y_1] = 0 \tag{5.68}$$

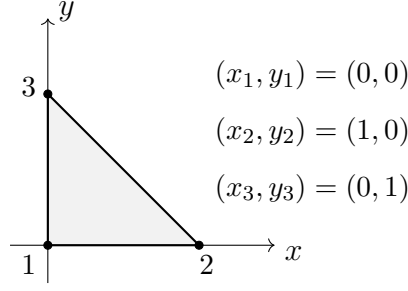
$$\mathcal{N}_2(x_2, y_2) = \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x_2 + (x_1 - x_3)y_2] = 1 \tag{5.69}$$

$$\mathcal{N}_2(x_3, y_3) = \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x_3 + (x_1 - x_3)y_3] = 0 \tag{5.70}$$

Note that the area A of the triangle is given by:

$$A = \frac{1}{2}D = \frac{1}{2} \begin{vmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ 1 & x_3 & y_3 \end{vmatrix}$$

If we now consider the reference element in the reduced coordinates space (r, s) :



The basis polynomial is then

$$f(r, s) = a + br + cs$$

and the basis functions:

$$\mathcal{N}_0(r, s) = 1 - r - s \quad (5.71)$$

$$\mathcal{N}_1(r, s) = r \quad (5.72)$$

$$\mathcal{N}_2(r, s) = s \quad (5.73)$$

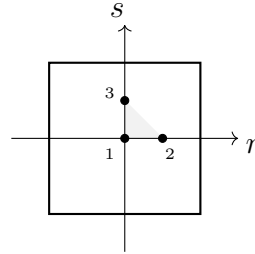
Once again we can verify that $\mathcal{N}_i(x_j, y_j) = \delta_{ij}$ and $\sum_i \mathcal{N}_i(r, s) = 1$.

5.3.8 Linear basis functions for quadrilaterals in 2D (P_1)

basis_Pm1_2D.tex

On the reference element $\Omega = [-1, 1] \times [-1, 1]$ we have three nodes placed as follows:

(tikz_pm1_2D.tex)



Let us assume that the function $f(r, s)$ is to be approximated on $[-1, 1] \times [-1, 1]$ by

$$f^h(r, s) = a + br + cs$$

Note that this is a linear function, not a bilinear one (a direct consequence of this is the fact that this function cannot be continuous from one element to another). The function f^h then must fulfill:

$$f^h(r_1, s_1) = a = f_1$$

$$f^h(r_2, s_2) = a + \frac{b}{2} = f_2$$

$$f^h(r_3, s_3) = a + \frac{c}{2} = f_3$$

This leads to :

$$a = f_1 \quad b = 2(f_2 - f_1) \quad c = 2(f_3 - f_1)$$

Then

$$f(r, s) = f_1 + 2(f_2 - f_1)r + 2(f_3 - f_1)s$$

or,

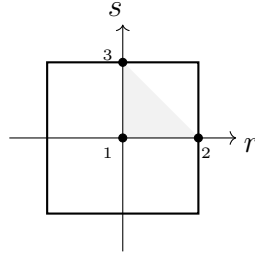
$$f(r) = \sum_{i=1}^3 \mathcal{N}_i(r, s) f_i$$

with

$$\begin{aligned} \mathcal{N}_1(r) &= 1 - 2(r + s) \\ \mathcal{N}_2(r) &= 2r \\ \mathcal{N}_3(r) &= 2s \end{aligned} \tag{5.74}$$

Note that we could also have placed the nodes at a different location:

(tikz_pm1_2D_bis.tex)



and we would then have

$$\begin{aligned} \mathcal{N}_1(r) &= 1 - r - s \\ \mathcal{N}_2(r) &= r \\ \mathcal{N}_3(r) &= s \end{aligned} \tag{5.75}$$

as obtained in Section 5.3.7.

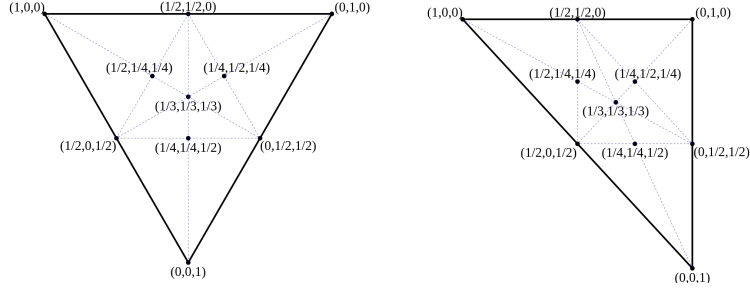
5.3.9 Enriched linear basis functions in triangles (P_1^+)

basis_p1p_2D.tex

As we will see in Section 7.3.14 the above P_1 can be enriched with a so-called bubble function. The bubble function of the MINI element is described in Arnold *et al.* (1984) [26] as being $\lambda_1 \lambda_2 \lambda_3$ where λ_i are the so-called barycentric coordinates⁸.

$$\begin{aligned} \lambda_1 &= \frac{(y_2 - y_3)(x - x_3) + (x_3 - x_2)(y - y_3)}{(y_2 - y_3)(x_1 - x_3) + (x_3 - x_2)(y_1 - y_3)} \\ \lambda_2 &= \frac{(y_3 - y_1)(x - x_3) + (x_1 - x_3)(y - y_3)}{(y_2 - y_3)(x_1 - x_3) + (x_3 - x_2)(y_1 - y_3)} \\ \lambda_3 &= 1 - \lambda_1 - \lambda_2 \end{aligned}$$

⁸https://en.wikipedia.org/wiki/Barycentric_coordinate_system



Barycentric coordinates $(\lambda_1, \lambda_2, \lambda_3)$ on an equilateral triangle and on a right triangle.

In the reference triangle, the barycentric coordinates write

$$\begin{aligned}\lambda_1 &= \frac{(s_2 - s_3)(r - r_3) + (r_3 - r_2)(s - s_3)}{(s_2 - s_3)(r_1 - r_3) + (r_3 - r_2)(s_1 - s_3)} = \frac{(-1)(r) + (-1)(s - 1)}{(-1)(0) + (-1)(-1)} = -r - s + 1 \\ \lambda_2 &= \frac{(s_3 - s_1)(r - r_3) + (r_1 - r_3)(s - s_3)}{(s_2 - s_3)(r_1 - r_3) + (r_3 - r_2)(s_1 - s_3)} = \frac{(1)(r) + (0)(s - 1)}{(-1)(0) + (-1)(-1)} = r \\ \lambda_3 &= 1 - \lambda_1 - \lambda_2 = 1 - (-r - s + 1) - r = s\end{aligned}$$

As we have seen before the bubble function is given by $\lambda_1 \lambda_2 \lambda_3 = (1 - r - s)rs$ and the polynomial form for the basis functions is given by:

$$f(r, s) = a + br + cs + d(1 - r - s)rs$$

Setting the location of the bubble at $r = s = 1/3$, i.e. $\lambda_1 \lambda_2 \lambda_3 = 1/3$, we then have

$$\begin{aligned}f(r_1, s_1) &= f_1 = a + br_1 + cs_1 + d(1 - r_1 - s_1)r_1s_1 = a \\ f(r_2, s_2) &= f_2 = a + br_2 + cs_2 + d(1 - r_2 - s_2)r_2s_2 = a + b \\ f(r_3, s_3) &= f_3 = a + br_3 + cs_3 + d(1 - r_3 - s_3)r_3s_3 = a + c \\ f(r_4, s_4) &= f_4 = a + br_4 + cs_4 + d(1 - r_4 - s_4)r_4s_4 = a + \frac{b}{3} + \frac{c}{3} + \frac{1}{27}\end{aligned}$$

where point 4 is the location of the bubble. This yields

$$a = f_1 \quad b = f_2 - a = f_2 - f_1 \quad c = f_3 - a = f_3 - f_1$$

and

$$d = 27 \left(f_4 - a - \frac{b}{3} - \frac{c}{3} \right) = 27 \left(f_4 - f_1 - \frac{f_2 - f_1}{3} - \frac{f_3 - f_1}{3} \right) = 27 \left(f_4 - \frac{f_1}{3} - \frac{f_2}{3} - \frac{f_3}{3} \right)$$

Finally

$$\begin{aligned}f(r, s) &= a + br + cs + d(1 - r - s)rs \\ &= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + 27 \left(f_4 - \frac{f_1}{3} - \frac{f_2}{3} - \frac{f_3}{3} \right) (1 - r - s)rs \\ &= [1 - r - s - 9(1 - r - s)rs]f_1 + [r - 9(1 - r - s)rs]f_2 + [s - 9(1 - r - s)rs]f_3 + [27(1 - r - s)rs]f_4\end{aligned}$$

so that

$$f(r, s) = \sum_{i=1}^4 \mathcal{N}_i(r, s) f_i$$

with

$$\begin{aligned}\mathcal{N}_1(r, s) &= 1 - r - s - 9(1 - r - s)rs \\ \mathcal{N}_2(r, s) &= r - 9(1 - r - s)rs \\ \mathcal{N}_3(r, s) &= s - 9(1 - r - s)rs \\ \mathcal{N}_4(r, s) &= 27(1 - r - s)rs\end{aligned}$$

It is trivial to verify that $\sum_i \mathcal{N}_i = 1$ for all values of r, s and the gradients of the basis functions are:

$$\frac{\partial \mathcal{N}_1}{\partial r}(r, s) = -1 - 9(1 - 2r - s)s \quad (5.76)$$

$$\frac{\partial \mathcal{N}_2}{\partial r}(r, s) = +1 - 9(1 - 2r - s)s \quad (5.77)$$

$$\frac{\partial \mathcal{N}_3}{\partial r}(r, s) = -9(1 - 2r - s)s \quad (5.78)$$

$$\frac{\partial \mathcal{N}_4}{\partial r}(r, s) = 27(1 - 2r - s)s \quad (5.79)$$

$$\frac{\partial \mathcal{N}_1}{\partial s}(r, s) = -1 - 9(1 - r - 2s)r \quad (5.80)$$

$$\frac{\partial \mathcal{N}_2}{\partial s}(r, s) = -9(1 - r - 2s)r \quad (5.81)$$

$$\frac{\partial \mathcal{N}_3}{\partial s}(r, s) = +1 - 9(1 - r - 2s)r \quad (5.82)$$

$$\frac{\partial \mathcal{N}_4}{\partial s}(r, s) = 27(1 - r - 2s)r \quad (5.83)$$

$$\frac{\partial \mathcal{N}_1}{\partial s}(r, s) = -1 - 9(1 - r - 2s)r \quad (5.84)$$

We have two coordinate systems for the element: the global Cartesian coordinates (x, y) and the natural/reduced coordinates (r, s) . Inside the element, the relation between the two is given by

$$\begin{aligned} x &= N_1x_1 + N_2x_2 + N_3x_3 + N_4x_4 = \sum_i \mathcal{N}_i(r, s)x_i \\ y &= N_1y_1 + N_2y_2 + N_3y_3 + N_4y_4 = \sum_i \mathcal{N}_i(r, s)y_i \end{aligned} \quad (5.85)$$

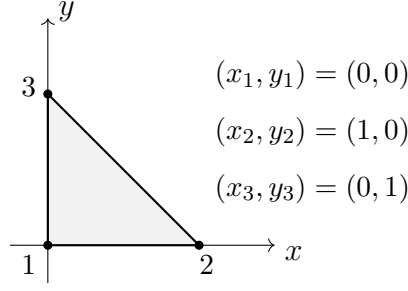
or,

$$\begin{aligned} x &= [1 - r - s - 9(1 - r - s)rs]x_1 + [r - 9(1 - r - s)rs]x_2 + [s - 9(1 - r - s)rs]x_3 + [27(1 - r - s)rs]x_4 \\ &= x_1 - r(x_1 - x_2) - s(x_1 - x_3) + (1 - r - s)rs(-9x_1 - 9x_2 - 9x_3 + 27x_4) \\ &= x_1 - r(x_1 - x_2) - s(x_1 - x_3) + (1 - r - s)rs(-9x_1 - 9x_2 - 9x_3 + 27(x_1 + x_2 + x_3)/3) \\ &= x_1 - r(x_1 - x_2) - s(x_1 - x_3) \\ &= x_1 - rx_{12} - sx_{13} \\ y &= [1 - r - s - 9(1 - r - s)rs]y_1 + [r - 9(1 - r - s)rs]y_2 + [s - 9(1 - r - s)rs]y_3 + [27(1 - r - s)rs]y_4 \\ &= y_1 - r(y_1 - y_2) - s(y_1 - y_3) + (1 - r - s)rs(-9y_1 - 9y_2 - 9y_3 + 27y_4) \\ &= y_1 - r(y_1 - y_2) - s(y_1 - y_3) + (1 - r - s)rs(-9y_1 - 9y_2 - 9y_3 + 27(y_1 + y_2 + y_3)/3) \\ &= y_1 - r(y_1 - y_2) - s(y_1 - y_3) \\ &= y_1 - ry_{12} - sy_{13} \end{aligned}$$

5.3.10 Quadratic basis functions for triangles in 2D (P_2)

basis_p2_2D.tex

(tikz_P2.tex)



The basis polynomial is then

$$f(r, s) = c_1 + c_2 r + c_3 s + c_4 r^2 + c_5 r s + c_6 s^2$$

We have

$$\begin{aligned} f_1 = f(r_1, s_1) &= c_1 \\ f_2 = f(r_2, s_2) &= c_1 + c_2 + c_4 \\ f_3 = f(r_3, s_3) &= c_1 + c_3 + c_6 \\ f_4 = f(r_4, s_4) &= c_1 + c_2/2 + c_4/4 \\ f_5 = f(r_5, s_5) &= c_1 + c_2/2 + c_3/2 + c_4/4 + c_5/4 + c_6/4 \\ f_6 = f(r_6, s_6) &= c_1 + c_3/2 + c_6/4 \end{aligned}$$

This can be cast as $\vec{f} = \mathbf{A} \cdot \vec{c}$ where \mathbf{A} is a 6×6 matrix:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1/2 & 0 & 1/4 & 0 & 0 \\ 1 & 1/2 & 1/2 & 1/4 & 1/4 & 1/4 \\ 1 & 0 & 1/2 & 0 & 0 & 1/4 \end{pmatrix}$$

As it turns out it is rather trivial to compute the inverse of this matrix:

$$\mathbf{A}^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ -3 & -1 & 0 & 4 & 0 & 0 \\ -3 & 0 & -1 & 0 & 0 & 4 \\ 2 & 2 & 0 & -4 & 0 & 0 \\ 4 & 0 & 0 & -4 & 4 & -4 \\ 2 & 0 & 2 & 0 & 0 & -4 \end{pmatrix}$$

Using $\vec{c} = \mathbf{A}^{-1} \cdot \vec{f}$ one then obtains:

$$\begin{aligned} c_1 &= f_1 \\ c_2 &= -3f_1 - f_2 + 4f_4 \\ c_3 &= -3f_1 - f_3 + 4f_6 \\ c_4 &= 2f_1 + 2f_2 - 4f_4 \\ c_5 &= 4f_1 - 4f_4 + 4f_5 - 4f_6 \\ c_6 &= (2f_1 + 2f_3 - 4f_6) \end{aligned}$$

and then

$$\begin{aligned} f(r, s) &= f_1 + (-3f_1 - f_2 + 4f_4)r + (-3f_1 - f_3 + 4f_6)s \\ &\quad + (2f_1 + 2f_2 - 4f_4)r^2 + (4f_1 - 4f_4 + 4f_5 - 4f_6)rs + (2f_1 + 2f_3 - 4f_6)s^2 \\ &= \sum_{i=1}^6 \mathcal{N}_i(r, s) f_i \end{aligned} \tag{5.86}$$

with

$$\begin{aligned}
\mathcal{N}_1(r, s) &= 1 - 3r - 3s + 2r^2 + 4rs + 2s^2 \\
\mathcal{N}_2(r, s) &= -r + 2r^2 \\
\mathcal{N}_3(r, s) &= -s + 2s^2 \\
\mathcal{N}_4(r, s) &= 4r - 4r^2 - 4rs \\
\mathcal{N}_5(r, s) &= 4rs \\
\mathcal{N}_6(r, s) &= 4s - 4rs - 4s^2
\end{aligned}$$

The derivatives are as follows:

$$\begin{aligned}
\frac{\partial \mathcal{N}_1}{\partial r}(r, s) &= -3 + 4r + 4s \\
\frac{\partial \mathcal{N}_2}{\partial r}(r, s) &= -1 + 4r \\
\frac{\partial \mathcal{N}_3}{\partial r}(r, s) &= 0 \\
\frac{\partial \mathcal{N}_4}{\partial r}(r, s) &= 4 - 8r - 4s \\
\frac{\partial \mathcal{N}_5}{\partial r}(r, s) &= 4s \\
\frac{\partial \mathcal{N}_6}{\partial r}(r, s) &= -4s
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{N}_1}{\partial s}(r, s) &= -3 + 4r + 4s \\
\frac{\partial \mathcal{N}_2}{\partial s}(r, s) &= 0 \\
\frac{\partial \mathcal{N}_3}{\partial s}(r, s) &= -1 + 4s \\
\frac{\partial \mathcal{N}_4}{\partial s}(r, s) &= -4r \\
\frac{\partial \mathcal{N}_5}{\partial s}(r, s) &= 4r \\
\frac{\partial \mathcal{N}_6}{\partial s}(r, s) &= 4 - 4r - 8s
\end{aligned}$$

5.3.11 Enriched quadratic basis functions in triangles (P_2^+)

basis_p2p_2D.tex

This is used by the Crouzeix-Raviart element, see Section 7.3.16.

TIKZ!

```

03      (r_1,s_1)=(0,0)
||\\    (r_2,s_2)=(1,0)
||  \\  (r_3,s_3)=(0,1)
||   \\ (r_4,s_4)=(1/2,0)
06  05  (r_5,s_5)=(1/2,1/2)

```

```

|| 07 \\      (r_6,s_6)=(0,1/2)
||      \\      (r_7,s_7)=(1/3,1/3)
01==04==02

```

The basis functions are given by:

[find reference](#)

$$\mathcal{N}_1(r, s) = (1 - r - s)(1 - 2r - 2s + 3rs) \quad (5.87)$$

$$\mathcal{N}_2(r, s) = r(2r - 1 + 3s - 3rs - 3s^2) \quad (5.88)$$

$$\mathcal{N}_3(r, s) = s(2s - 1 + 3r - 3r^2 - 3rs) \quad (5.89)$$

$$\mathcal{N}_4(r, s) = 4(1 - r - s)r(1 - 3s) \quad (5.90)$$

$$\mathcal{N}_5(r, s) = 4rs[-2 + 3r + 3s] \quad (5.91)$$

$$\mathcal{N}_6(r, s) = 4(1 - r - s)s(1 - 3r) \quad (5.92)$$

$$\mathcal{N}_7(r, s) = 27(1 - r - s)rs \quad (5.93)$$

It is then easy to verify that for all basis functions we have $\mathcal{N}_i(r_j, s_j) = \delta_{ij}$ where j denotes one of the seven nodes. The derivatives are as follows:

$$\frac{\partial \mathcal{N}_1}{\partial r}(r, s) = r(4 - 6s) - 3s^2 + 7s - 3 \quad (5.94)$$

$$\frac{\partial \mathcal{N}_2}{\partial r}(r, s) = r(4 - 6s) - 3s^2 + 3s - 1 \quad (5.95)$$

$$\frac{\partial \mathcal{N}_3}{\partial r}(r, s) = -3s(2r + s - 1) \quad (5.96)$$

$$\frac{\partial \mathcal{N}_4}{\partial r}(r, s) = 4(3s - 1)(2r + s - 1) \quad (5.97)$$

$$\frac{\partial \mathcal{N}_5}{\partial r}(r, s) = 4s(6r + 3s - 2) \quad (5.98)$$

$$\frac{\partial \mathcal{N}_6}{\partial r}(r, s) = 4s(6r + 3s - 4) \quad (5.99)$$

$$\frac{\partial \mathcal{N}_7}{\partial r}(r, s) = -27s(2r + s - 1) \quad (5.100)$$

$$\frac{\partial \mathcal{N}_1}{\partial s}(r, s) = -3r^2 + r(7 - 6s) + 4s - 3 \quad (5.101)$$

$$\frac{\partial \mathcal{N}_2}{\partial s}(r, s) = -3r(r + 2s - 1) \quad (5.102)$$

$$\frac{\partial \mathcal{N}_3}{\partial s}(r, s) = -3r^2 + r(3 - 6s) + 4s - 1 \quad (5.103)$$

$$\frac{\partial \mathcal{N}_4}{\partial s}(r, s) = 4r(3r + 6s - 4) \quad (5.104)$$

$$\frac{\partial \mathcal{N}_5}{\partial s}(r, s) = 4r(3r + 6s - 2) \quad (5.105)$$

$$\frac{\partial \mathcal{N}_6}{\partial s}(r, s) = 4(3r - 1)(r + 2s - 1) \quad (5.106)$$

$$\frac{\partial \mathcal{N}_7}{\partial s}(r, s) = -27r(r + 2s - 1) \quad (5.107)$$

Note that the basis functions can also be expressed as a function of the barycentric coordinates, as in the MILAMIN code [299] or in Cuvelier *et al.* (1986) [298]⁹

TIKZ!

```

03
||\
|| \
||  \
05    04
|| 07 \
||    \
01==06==02

```

$$N_1(\lambda_1, \lambda_2, \lambda_3) = \eta_1(2\eta_1 - 1) + 3\eta_1\eta_2\eta_3 \quad (5.108)$$

$$N_2(\lambda_1, \lambda_2, \lambda_3) = \eta_2(2\eta_2 - 1) + 3\eta_1\eta_2\eta_3 \quad (5.109)$$

$$N_3(\lambda_1, \lambda_2, \lambda_3) = \eta_3(2\eta_3 - 1) + 3\eta_1\eta_2\eta_3 \quad (5.110)$$

$$N_4(\lambda_1, \lambda_2, \lambda_3) = 4\eta_2\eta_3 - 12\eta_1\eta_2\eta_3 \quad (5.111)$$

$$N_5(\lambda_1, \lambda_2, \lambda_3) = 4\eta_1\eta_3 - 12\eta_1\eta_2\eta_3 \quad (5.112)$$

$$N_6(\lambda_1, \lambda_2, \lambda_3) = 4\eta_1\eta_2 - 12\eta_1\eta_2\eta_3 \quad (5.113)$$

$$N_7(\lambda_1, \lambda_2, \lambda_3) = 27\eta_1\eta_2\eta_3 \quad (5.114)$$

VERIFY that when $\eta_1 = 1 - r - s$, $\eta_2 = r$ and $\eta_3 = s$ we find the above r, s basis functions

5.3.12 Cubic basis functions for triangles (P_3)

basis_p3_2D.tex

TIKZ!

```

9
|\      (r_0,s_0)=(0,0)    (r_5,s_5)=(1/3,1/3)
| \     (r_1,s_1)=(1/3,0)  (r_6,s_6)=(2/3,1/3)
7   8    (r_2,s_2)=(2/3,0) (r_7,s_7)=(0,2/3)
|   \    (r_3,s_3)=(1,0)   (r_8,s_8)=(1/3,2/3)
4   5   6 (r_4,s_4)=(0,1/3) (r_9,s_9)=(0,1)
|       \
0==1==2==3

```

The basis polynomial is then

$$f(r, s) = c_0 + c_1r + c_2s + c_3r^2 + c_4rs + c_5s^2 + c_6r^3 + c_7r^2s + c_8rs^2 + c_9s^3$$

with the support nodes being given by

⁹Note that the numbering of the nodes in the book is different with respect to the one above.

$$(r_0, s_0) = (0, 0) \quad (5.115)$$

$$(r_1, s_1) = (1/3, 0) \quad (5.116)$$

$$(r_2, s_2) = (2/3, 0) \quad (5.117)$$

$$(r_3, s_3) = (1, 0) \quad (5.118)$$

$$(r_4, s_4) = (0, 1/3) \quad (5.119)$$

$$(r_5, s_5) = (1/3, 1/3) \quad (5.120)$$

$$(r_6, s_6) = (2/3, 1/3) \quad (5.121)$$

$$(r_7, s_7) = (0, 2/3) \quad (5.122)$$

$$(r_8, s_8) = (1/3, 2/3) \quad (5.123)$$

$$(r_9, s_9) = (0, 1) \quad (5.124)$$

$$f_0 = f(r_0, s_0) = c_0$$

$$f_1 = f(r_1, s_1) = c_0 + c_1 + \frac{1}{9}c_3 + \frac{1}{27}c_6$$

$$f_2 = f(r_2, s_2) = c_0 + \frac{2}{3}c_1 + \frac{4}{9}c_3 + \frac{8}{27}c_6$$

$$f_3 = f(r_3, s_3) = c_0 + c_1 + c_3 + c_6$$

$$f_4 = f(r_4, s_4) = c_0 + \frac{1}{3}c_2 + \frac{1}{9}c_5 + \frac{1}{27}c_9$$

$$f_5 = f(r_5, s_5) = c_0 + \frac{1}{3}c_1 + \frac{1}{3}c_2 + \frac{1}{9}c_3 + \frac{1}{9}c_4 + \frac{1}{9}c_5 + \frac{1}{27}c_6r^3 + \frac{1}{27}c_7 + \frac{1}{27}c_8 + \frac{1}{27}c_9$$

$$f_6 = f(r_6, s_6) = c_0 + \frac{2}{3}c_1 + \frac{1}{3}c_2 + \frac{4}{9}c_3 + \frac{2}{9}c_4 + \frac{1}{9}c_5 + \frac{8}{27}c_6 + \frac{4}{27}c_7 + \frac{2}{27}c_8 + \frac{1}{9}c_9$$

$$f_7 = f(r_7, s_7) = c_0 + \frac{2}{3}c_2 + \frac{4}{9}c_5 + \frac{8}{27}c_9$$

$$f_8 = f(r_8, s_8) = c_0 + \frac{1}{3}c_1 + \frac{2}{3}c_2 + \frac{1}{9}c_3 + \frac{2}{9}c_4 + \frac{4}{9}c_5 + \frac{1}{27}c_6 + \frac{2}{27}c_7 + \frac{4}{27}c_8 + \frac{8}{27}c_9s^3$$

$$f_9 = f(r_9, s_9) = c_0 + c_2 + c_5 + c_9$$

or,

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & \frac{1}{3} & 0 & \frac{1}{9} & 0 & 0 & \frac{1}{27} & 0 & 0 & 0 \\ 1 & \frac{2}{3} & 0 & \frac{4}{9} & 0 & 0 & \frac{8}{27} & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & \frac{1}{3} & 0 & 0 & \frac{1}{9} & 0 & 0 & 0 & \frac{1}{27} \\ 1 & \frac{1}{3} & \frac{1}{3} & \frac{1}{9} & \frac{1}{9} & \frac{1}{9} & \frac{1}{27} & \frac{1}{27} & \frac{1}{27} & \frac{1}{27} \\ 1 & \frac{2}{3} & \frac{1}{3} & \frac{4}{9} & \frac{2}{9} & \frac{1}{9} & \frac{8}{27} & \frac{4}{27} & \frac{2}{27} & \frac{1}{27} \\ 1 & 0 & \frac{2}{3} & 0 & 0 & \frac{4}{9} & 0 & 0 & 0 & \frac{8}{27} \\ 1 & \frac{1}{3} & \frac{2}{3} & \frac{1}{9} & \frac{2}{9} & \frac{4}{9} & \frac{1}{27} & \frac{2}{27} & \frac{4}{27} & \frac{8}{27} \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \\ c_8 \\ c_9 \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \end{pmatrix}$$

or,

$$\frac{1}{27} \begin{pmatrix} 27 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 27 & 9 & 0 & 3 & 0 & 0 & 1 & 0 & 0 & 0 \\ 27 & 18 & 0 & 12 & 0 & 0 & 8 & 0 & 0 & 0 \\ 27 & 27 & 0 & 27 & 0 & 0 & 27 & 0 & 0 & 0 \\ 27 & 0 & 9 & 0 & 0 & 3 & 0 & 0 & 0 & 1 \\ 27 & 9 & 9 & 3 & 3 & 3 & 1 & 1 & 1 & 1 \\ 27 & 18 & 9 & 12 & 6 & 3 & 8 & 4 & 2 & 1 \\ 27 & 0 & 18 & 0 & 0 & 12 & 0 & 0 & 0 & 8 \\ 27 & 9 & 18 & 3 & 6 & 12 & 1 & 2 & 4 & 8 \\ 27 & 0 & 27 & 0 & 0 & 27 & 0 & 0 & 0 & 27 \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \\ c_8 \\ c_9 \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \end{pmatrix}$$

The inverse of the matrix is

$$\frac{1}{2} \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -11 & 18 & -9 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ -11 & 0 & 0 & 0 & 18 & 0 & 0 & -9 & 0 & 2 \\ 18 & -45 & 36 & -9 & 0 & 0 & 0 & 0 & 0 & 0 \\ 36 & -45 & 9 & 0 & -45 & 54 & -9 & 9 & -9 & 0 \\ 18 & 0 & 0 & 0 & -45 & 0 & 0 & 36 & 0 & -9 \\ -9 & 27 & -27 & 9 & 0 & 0 & 0 & 0 & 0 & 0 \\ -27 & 54 & -27 & 0 & 27 & -54 & 27 & 0 & 0 & 0 \\ -27 & 27 & 0 & 0 & 54 & -54 & 0 & -27 & 27 & 0 \\ -9 & 0 & 0 & 0 & 27 & 0 & 0 & -27 & 0 & 9 \end{pmatrix}$$

so that the solution of the system $\mathbf{A} \cdot \vec{c} = \vec{f}$ is $\vec{c} = \mathbf{A}^{-1} \cdot \vec{f}$, or:

$$\begin{aligned} c_0 &= 1 \\ c_1 &= \frac{1}{2}(-11f_0 + 18f_1 - 9f_2 + 2f_3) \\ c_2 &= \frac{1}{2}(-11f_0 + 18f_4 - 9f_7 + 2f_9) \\ c_3 &= \text{etc ...} \end{aligned} \tag{5.125}$$

which we insert in

$$f(r, s) = c_0 + c_1 r + c_2 s + c_3 r^2 + c_4 r s + c_5 s^2 + c_6 r^3 + c_7 r^2 s + c_8 r s^2 + c_9 s^3$$

and we then obtain

$$\begin{aligned}
f(r, s) &= \frac{1}{2} (2 - 11r - 11s + 18r^2 + 36rs + 18s^2 - 9r^3 - 27r^2s - 27rs^2 - 9s^3) f_0 \\
&+ \frac{1}{2} (18r - 45r^2 - 45rs + 27r^3 + 54r^2s + 27rs^2) f_1 \\
&+ \frac{1}{2} (-9r + 36r^2 + 9rs - 27r^3 - 27r^2s) f_2 \\
&+ \frac{1}{2} (2r - 9r^2 + 9r^3) f_3 \\
&+ \frac{1}{2} (18s - 45rs - 45s^2 + 27r^2s + 54rs^2 + 27s^3) f_4 \\
&+ \frac{1}{2} (54rs - 54r^2s - 54rs^2) f_5 \\
&+ \frac{1}{2} (-9rs + 27r^2s) f_6 \\
&+ \frac{1}{2} (-9s + 9rs + 36s^2 - 27rs^2 - 27s^3) f_7 \\
&+ \frac{1}{2} (-9rs + 27rs^2) f_8 \\
&+ \frac{1}{2} (2s - 9s^2 + 9s^3) f_9 \\
&= \sum_{i=0}^9 \mathcal{N}_i(r, s) f_i
\end{aligned}$$

$$\begin{aligned}
\mathcal{N}_0(r, s) &= \frac{1}{2} (2 - 11r - 11s + 18r^2 + 36rs + 18s^2 - 9r^3 - 27r^2s - 27rs^2 - 9s^3) \\
\mathcal{N}_1(r, s) &= \frac{1}{2} (18r - 45r^2 - 45rs + 27r^3 + 54r^2s + 27rs^2) \\
\mathcal{N}_2(r, s) &= \frac{1}{2} (-9r + 36r^2 + 9rs - 27r^3 - 27r^2s) \\
\mathcal{N}_3(r, s) &= \frac{1}{2} (2r - 9r^2 + 9r^3) \\
\mathcal{N}_4(r, s) &= \frac{1}{2} (18s - 45rs - 45s^2 + 27r^2s + 54rs^2 + 27s^3) \\
\mathcal{N}_5(r, s) &= \frac{1}{2} (54rs - 54r^2s - 54rs^2) \\
\mathcal{N}_6(r, s) &= \frac{1}{2} (-9rs + 27r^2s) \\
\mathcal{N}_7(r, s) &= \frac{1}{2} (-9s + 9rs + 36s^2 - 27rs^2 - 27s^3) \\
\mathcal{N}_8(r, s) &= \frac{1}{2} (-9rs + 27rs^2) \\
\mathcal{N}_9(r, s) &= \frac{1}{2} (2s - 9s^2 + 9s^3)
\end{aligned} \tag{5.126}$$

and then

$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial r}(r, s) &= \frac{1}{2}(-11 + 36r + 36s - 27r^2 - 54rs - 27s^2) \\
\frac{\partial \mathcal{N}_1}{\partial r}(r, s) &= \frac{1}{2}(18 - 90r - 45s + 81r^2 + 108rs + 27s^2) \\
\frac{\partial \mathcal{N}_2}{\partial r}(r, s) &= \frac{1}{2}(-9 + 72r + 9s - 81r^2 - 54rs) \\
\frac{\partial \mathcal{N}_3}{\partial r}(r, s) &= \frac{1}{2}(2 - 18r + 27r^2) \\
\frac{\partial \mathcal{N}_4}{\partial r}(r, s) &= \frac{1}{2}(-45s + 54rs + 54s^2) \\
\frac{\partial \mathcal{N}_5}{\partial r}(r, s) &= \frac{1}{2}(54s - 108rs - 54s^2) \\
\frac{\partial \mathcal{N}_6}{\partial r}(r, s) &= \frac{1}{2}(-9s + 54rs) \\
\frac{\partial \mathcal{N}_7}{\partial r}(r, s) &= \frac{1}{2}(9s - 27s^2) \\
\frac{\partial \mathcal{N}_8}{\partial r}(r, s) &= \frac{1}{2}(-9s + 27s^2) \\
\frac{\partial \mathcal{N}_9}{\partial r}(r, s) &= 0
\end{aligned}$$

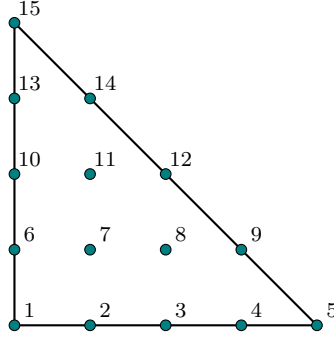
$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial s}(r, s) &= \frac{1}{2}(-11 + 36r + 36s - 27r^2 - 54rs - 27s^2) \\
\frac{\partial \mathcal{N}_1}{\partial s}(r, s) &= \frac{1}{2}(-45r + 54r^2 + 54rs) \\
\frac{\partial \mathcal{N}_2}{\partial s}(r, s) &= \frac{1}{2}(9r - 27r^2) \\
\frac{\partial \mathcal{N}_3}{\partial s}(r, s) &= 0 \\
\frac{\partial \mathcal{N}_4}{\partial s}(r, s) &= \frac{1}{2}(18 - 45r - 90s + 27r^2 + 108rs + 81s^2) \\
\frac{\partial \mathcal{N}_5}{\partial s}(r, s) &= \frac{1}{2}(54r - 54r^2 - 108rs) \\
\frac{\partial \mathcal{N}_6}{\partial s}(r, s) &= \frac{1}{2}(-9r + 27r^2) \\
\frac{\partial \mathcal{N}_7}{\partial s}(r, s) &= \frac{1}{2}(-9 + 9r + 72s - 54rs - 81s^2) \\
\frac{\partial \mathcal{N}_8}{\partial s}(r, s) &= \frac{1}{2}(-9r + 54rs) \\
\frac{\partial \mathcal{N}_9}{\partial s}(r, s) &= \frac{1}{2}(2 - 18s + 27s^2)
\end{aligned}$$

It is implemented in [STONE](#) 120. See also python code in `images/basis.P3` which I wrote to test these basis functions.

5.3.13 Quartic basis functions for triangles (P_4)

basis-p4-2D.tex

(tikz-p4.tex)



The support nodes coordinates are as follows:

| $i \rightarrow$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|-----------------|---|---------------|---------------|---------------|---|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|----|
| r_i | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{3}{4}$ | 1 | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{3}{4}$ | 0 | $\frac{1}{4}$ | $\frac{1}{2}$ | 0 | $\frac{1}{4}$ | 0 |
| s_i | 0 | 0 | 0 | 0 | 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{3}{4}$ | $\frac{3}{4}$ | 1 |

Inside the element a field f is represented by a 4-th order polynomial:

$$\begin{aligned}
 f^h(r, s) = & c_0 + c_1 r + c_2 s \\
 & + c_3 r^2 + c_4 r s + c_5 s^2 \\
 & + c_6 r^3 + c_7 r^2 s + c_8 r s^2 + c_9 s^3 \\
 & + c_{10} r^4 + c_{11} r^3 s + c_{12} r^2 s^2 + c_{13} r s^3 + c_{14} s^4
 \end{aligned} \tag{5.127}$$

At each node the function takes a value f_i , $i \in [1, 15]$ so that we have:

$$\begin{pmatrix}
 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 1 & \frac{1}{4} & 0 & \frac{1}{16} & 0 & 0 & \frac{1}{64} & 0 & 0 & 0 & \frac{1}{256} & 0 & 0 & 0 & 0 \\
 1 & \frac{1}{2} & 0 & \frac{1}{4} & 0 & 0 & \frac{1}{8} & 0 & 0 & 0 & \frac{1}{16} & 0 & 0 & 0 & 0 \\
 1 & \frac{3}{4} & 0 & \frac{9}{16} & 0 & 0 & \frac{27}{64} & 0 & 0 & 0 & \frac{81}{256} & 0 & 0 & 0 & 0 \\
 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
 \\
 1 & 0 & \frac{1}{4} & 0 & 0 & \frac{1}{16} & 0 & 0 & 0 & \frac{1}{64} & 0 & 0 & 0 & 0 & \frac{1}{256} \\
 1 & \frac{1}{4} & \frac{1}{4} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} & \frac{1}{64} & \frac{1}{64} & \frac{1}{64} & \frac{1}{64} & \frac{1}{256} & \frac{1}{256} & \frac{1}{256} & \frac{1}{256} & \frac{1}{256} \\
 1 & \frac{1}{2} & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} \\
 1 & \frac{3}{4} & \frac{3}{4} & \frac{9}{16} & \frac{9}{16} & \frac{9}{16} & \frac{27}{64} & \frac{27}{64} & \frac{27}{64} & \frac{27}{64} & \frac{81}{256} & \frac{81}{256} & \frac{81}{256} & \frac{81}{256} & \frac{81}{256} \\
 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 \\
 1 & 0 & \frac{1}{2} & 0 & 0 & \frac{1}{4} & 0 & 0 & 0 & \frac{1}{8} & 0 & 0 & 0 & 0 & \frac{1}{16} \\
 1 & \frac{1}{4} & \frac{1}{2} & \frac{1}{16} & \frac{1}{8} & \frac{1}{4} & \frac{1}{64} & \frac{1}{32} & \frac{1}{16} & \frac{1}{8} & \frac{1}{256} & \frac{1}{128} & \frac{1}{64} & \frac{1}{32} & \frac{1}{16} \\
 1 & \frac{1}{2} & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} & \frac{1}{16} \\
 \\
 1 & 0 & \frac{3}{4} & 0 & 0 & \frac{9}{16} & 0 & 0 & 0 & \frac{27}{64} & 0 & 0 & 0 & 0 & \frac{81}{256} \\
 1 & \frac{1}{4} & \frac{3}{4} & \frac{1}{16} & \frac{3}{16} & \frac{9}{16} & \frac{1}{64} & \frac{3}{64} & \frac{9}{64} & \frac{27}{64} & \frac{1}{256} & \frac{3}{256} & \frac{9}{256} & \frac{27}{256} & \frac{81}{256} \\
 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1
 \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \\ c_8 \\ c_9 \\ c_{10} \\ c_{11} \\ c_{12} \\ c_{13} \\ c_{14} \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ f_2 \\ \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \\ f_{10} \\ f_{11} \\ f_{12} \\ f_{13} \\ f_{14} \end{pmatrix}$$

or,

$$\frac{1}{256} \begin{pmatrix} 256 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 256 & 64 & 0 & 16 & 0 & 0 & 4 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 256 & 128 & 0 & 64 & 0 & 0 & 32 & 0 & 0 & 0 & 16 & 0 & 0 & 0 & 0 \\ 256 & 192 & 0 & 144 & 0 & 0 & 108 & 0 & 0 & 0 & 81 & 0 & 0 & 0 & 0 \\ 256 & 256 & 0 & 256 & 0 & 0 & 256 & 0 & 0 & 0 & 256 & 0 & 0 & 0 & 0 \\ \\ 256 & 0 & 64 & 0 & 0 & 16 & 0 & 0 & 0 & 4 & 0 & 0 & 0 & 0 & 1 \\ 256 & 64 & 64 & 16 & 16 & 16 & 4 & 4 & 4 & 4 & 1 & 1 & 1 & 1 & 1 \\ 256 & 128 & 64 & 64 & 32 & 16 & 32 & 16 & 8 & 4 & 16 & 8 & 4 & 2 & 1 \\ 256 & 192 & 64 & 144 & 48 & 16 & 108 & 36 & 12 & 4 & 81 & 27 & 9 & 3 & 1 \\ \\ 256 & 0 & 128 & 0 & 0 & 64 & 0 & 0 & 0 & 32 & 0 & 0 & 0 & 0 & 16 \\ 256 & 64 & 128 & 16 & 32 & 64 & 4 & 8 & 16 & 32 & 1 & 2 & 4 & 8 & 16 \\ 256 & 128 & 128 & 64 & 64 & 64 & 32 & 32 & 32 & 32 & 16 & 16 & 16 & 16 & 16 \\ \\ 256 & 0 & 192 & 0 & 0 & 144 & 0 & 0 & 0 & 108 & 0 & 0 & 0 & 0 & 81 \\ 256 & 64 & 192 & 16 & 48 & 144 & 4 & 12 & 36 & 108 & 1 & 3 & 9 & 27 & 81 \\ \\ 256 & 0 & 256 & 0 & 0 & 256 & 0 & 0 & 0 & 256 & 0 & 0 & 0 & 0 & 256 \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ \\ c_3 \\ c_4 \\ c_5 \\ c_6 \\ c_7 \\ c_8 \\ c_9 \\ \\ c_{10} \\ c_{11} \\ c_{12} \\ c_{13} \\ c_{14} \end{pmatrix} = \begin{pmatrix} f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \\ f \end{pmatrix}$$

The inverse of the matrix is:

$$\frac{1}{3} \begin{pmatrix} 3 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -25 & 48 & -36 & 16 & -3 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -25 & 0 & 0 & 0 & 0 & 48 & 0 & 0 & 0 & -36 & 0 & 0 & 16 & 0 & -3 \\ \\ 70 & -208 & 228 & -112 & 22 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 140 & -208 & 84 & -16 & 0 & -208 & 288 & -96 & 16 & 84 & -96 & 12 & -16 & 16 & 0 \\ 70 & 0 & 0 & 0 & 0 & -208 & 0 & 0 & 0 & 228 & 0 & 0 & -112 & 0 & 22 \\ \\ -80 & 288 & -384 & 224 & -48 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -240 & 576 & -432 & 96 & 0 & 288 & -672 & 480 & -96 & -48 & 96 & -48 & 0 & 0 & 0 \\ -240 & 288 & -48 & 0 & 0 & 576 & -672 & 96 & 0 & -432 & 480 & -48 & 96 & -96 & 0 \\ -80 & 0 & 0 & 0 & 0 & 288 & 0 & 0 & 0 & -384 & 0 & 0 & 224 & 0 & -48 \\ \\ 32 & -128 & 192 & -128 & 32 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 128 & -384 & 384 & -128 & 0 & -128 & 384 & -384 & 128 & 0 & 0 & 0 & 0 & 0 & 0 \\ 192 & -384 & 192 & 0 & 0 & -384 & 768 & -384 & 0 & 192 & -384 & 192 & 0 & 0 & 0 \\ 128 & -128 & 0 & 0 & 0 & -384 & 384 & 0 & 0 & 384 & -384 & 0 & -128 & 128 & 0 \\ 32 & 0 & 0 & 0 & 0 & -128 & 0 & 0 & 0 & 192 & 0 & 0 & -128 & 0 & 32 \end{pmatrix}$$

so that we obtain

$$\begin{aligned}
f(r, s) &= \frac{1}{3}(3 - 25r - 25s + 70r^2 + 140rs + 70s^2 - 80r^3 - 240r^2s - 240rs^2 - 80s^3 + 32r^4 + 128r^3s + 192r^2s^2 + 128rs^3 + 32s^4)f_0 \\
&+ \frac{1}{3}(48r - 208r^2 - 208rs + 288r^3 + 576r^2s + 288rs^2 - 128r^4 - 384r^3s - 384r^2s^2 - 128rs^3)f_1 \\
&+ \frac{1}{3}(-36r + 228r^2 + 84rs - 384r^3 - 432r^2s - 48rs^2 + 192r^4 + 384r^3s + 192r^2s^2)f_2 \\
&+ \frac{1}{3}(16r - 112r^2 - 16rs + 224r^3 + 96r^2s - 128r^4 - 128r^3s)f_3 \\
&+ \frac{1}{3}(-3r + 22r^2 - 48r^3 + 32r^4)f_4 \\
&+ \frac{1}{3}(48s - 208rs - 208s^2 + 288r^2s + 576rs^2 + 288s^3 - 128r^3s - 384r^2s^2 - 384rs^3 - 128s^4)f_5 \\
&+ \frac{1}{3}(288rs - 672r^2s - 672rs^2 + 384r^3s + 768r^2s^2 + 384rs^3)f_6 \\
&+ \frac{1}{3}(-96rs + 480r^2s + 96rs^2 - 384r^3s - 384r^2s^2)f_7 \\
&+ \frac{1}{3}(16rs - 96r^2s + 128r^3s)f_8 \\
&+ \frac{1}{3}(-36s + 84rs + 228s^2 - 48r^2s - 432rs^2 - 384s^3 + 192r^2s^2 + 384rs^3 + 192s^4)f_9 \\
&+ \frac{1}{3}(-96rs + 96r^2s + 480rs^2 - 384r^2s^2 - 384rs^3)f_{10} \\
&+ \frac{1}{3}(12rs - 48r^2s - 48rs^2 + 192r^2s^2)f_{11} \\
&+ \frac{1}{3}(16s - 16rs - 112s^2 + 96rs^2 + 224s^3 - 128rs^3 - 128s^4)f_{12} \\
&+ \frac{1}{3}(16rs - 96rs^2 + 128rs^3)f_{13} \\
&+ \frac{1}{3}(-3s + 22s^2 - 48s^3 + 32s^4)f_{14} \\
&= \sum_{i=0}^{14} \mathcal{N}_i(r, s)f_i
\end{aligned} \tag{5.128}$$

and then

$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial r}(r, s) &= \frac{1}{3}(-25 + 140r + 140s - 240r^2 - 480rs - 240s^2 + 128r^3 + 384r^2s + 384rs^2 + 128s^3) \\
\frac{\partial \mathcal{N}_1}{\partial r}(r, s) &= \frac{1}{3}(48 - 416r - 208s + 864r^2 + 1152rs + 288s^2 - 512r^3 - 1152r^2s - 768rs^2 - 128s^3) \\
\frac{\partial \mathcal{N}_2}{\partial r}(r, s) &= \frac{1}{3}(-36 + 456r + 84s - 1152r^2 - 864rs - 48s^2 + 768r^3 + 1152r^2s + 384rs^2) \\
\frac{\partial \mathcal{N}_3}{\partial r}(r, s) &= \frac{1}{3}(16 - 224r - 16s + 672r^2 + 192rs - 512r^3 - 384r^2s) \\
\frac{\partial \mathcal{N}_4}{\partial r}(r, s) &= \frac{1}{3}(-3 + 44r - 144r^2 + 128r^3) \\
\frac{\partial \mathcal{N}_5}{\partial r}(r, s) &= \frac{1}{3}(-208s + 576rs + 576s^2 - 384r^2s - 768rs^2 - 384s^3) \\
\frac{\partial \mathcal{N}_6}{\partial r}(r, s) &= \frac{1}{3}(288s - 1344rs - 672s^2 + 1152r^2s + 1536rs^2 + 384s^3) \\
\frac{\partial \mathcal{N}_7}{\partial r}(r, s) &= \frac{1}{3}(-96s + 960rs + 96s^2 - 1152r^2s - 768rs^2) \\
\frac{\partial \mathcal{N}_8}{\partial r}(r, s) &= \frac{1}{3}(16s - 192rs + 384r^2s) \\
\frac{\partial \mathcal{N}_9}{\partial r}(r, s) &= \frac{1}{3}(84s - 96rs - 432s^2 + 384rs^2 + 384s^3) \\
\frac{\partial \mathcal{N}_{10}}{\partial r}(r, s) &= \frac{1}{3}(-96s + 192rs + 480s^2 - 768rs^2 - 384s^3) \\
\frac{\partial \mathcal{N}_{11}}{\partial r}(r, s) &= \frac{1}{3}(12s - 96rs - 48s^2 + 384rs^2) \\
\frac{\partial \mathcal{N}_{12}}{\partial r}(r, s) &= \frac{1}{3}(-16s + 96s^2 - 128s^3) \\
\frac{\partial \mathcal{N}_{13}}{\partial r}(r, s) &= \frac{1}{3}(16s - 96s^2 + 128s^3) \\
\frac{\partial \mathcal{N}_{14}}{\partial r}(r, s) &= 0
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial s}(r, s) &= \frac{1}{3}(-25 + 140r + 140s - 240r^2 - 480rs - 240s^2 + 128r^3 + 384r^2s + 384rs^2 + 128s^3) \\
\frac{\partial \mathcal{N}_1}{\partial s}(r, s) &= \frac{1}{3}(-208r + 576r^2 + 576rs - 384r^3 - 768r^2s - 384rs^2) \\
\frac{\partial \mathcal{N}_2}{\partial s}(r, s) &= \frac{1}{3}(84r - 432r^2 - 96rs + 384r^3 + 384r^2s) \\
\frac{\partial \mathcal{N}_3}{\partial s}(r, s) &= \frac{1}{3}(-16r + 96r^2 - 128r^3) \\
\frac{\partial \mathcal{N}_4}{\partial s}(r, s) &= 0 \\
\frac{\partial \mathcal{N}_5}{\partial s}(r, s) &= \frac{1}{3}(48 - 208r - 416s + 288r^2 + 1152rs + 864s^2 - 128r^3 - 768r^2s - 1152rs^2 - 512s^3) \\
\frac{\partial \mathcal{N}_6}{\partial s}(r, s) &= \frac{1}{3}(288r - 672r^2 - 1344rs + 384r^3 + 1536r^2s + 1152rs^2) \\
\frac{\partial \mathcal{N}_7}{\partial s}(r, s) &= \frac{1}{3}(-96r + 480r^2 + 192rs - 384r^3 - 768r^2s) \\
\frac{\partial \mathcal{N}_8}{\partial s}(r, s) &= \frac{1}{3}(16r - 96r^2 + 128r^3) \\
\frac{\partial \mathcal{N}_9}{\partial s}(r, s) &= \frac{1}{3}(-36 + 84r + 456s - 48r^2 - 864rs - 1152s^2 + 384r^2s + 1152rs^2 + 768s^3) \\
\frac{\partial \mathcal{N}_{10}}{\partial s}(r, s) &= \frac{1}{3}(-96r + 96r^2 + 960rs - 768r^2s - 1152rs^2) \\
\frac{\partial \mathcal{N}_{11}}{\partial s}(r, s) &= \frac{1}{3}(12r - 48r^2 - 96rs + 384r^2s) \\
\frac{\partial \mathcal{N}_{12}}{\partial s}(r, s) &= \frac{1}{3}(16 - 16r - 224s + 192rs + 672s^2 - 384rs^2 - 512s^3) \\
\frac{\partial \mathcal{N}_{13}}{\partial s}(r, s) &= \frac{1}{3}(16r - 192rs + 384rs^2) \\
\frac{\partial \mathcal{N}_{14}}{\partial s}(r, s) &= \frac{1}{3}(-3 + 44s - 144s^2 + 128s^3)
\end{aligned}$$

It is implemented in [STONE](#) 120. See also python code in `images/basis_P4` which I wrote to test these basis functions.

5.3.14 Enriched linear basis functions in quadrilaterals (Q_1^+) -WIP

basis_q1p_2D.tex

```

4=====3
|           |   (r_1,s_1)=(-1,-1)
|           |   (r_2,s_2)=(1,-1)
|    5      |   (r_3,s_3)=(1,1)
|           |   (r_4,s_4)=(-1,1)
|           |   (r_5,s_5)=(0,0)
1=====2

```

- In Bai [38] (1997): "It is well known that the equal-order bilinear velocity-bilinear continuous pressure element - the $Q_1 \times Q_1$, element - exhibits a certain spurious pressure mode. In the paper we propose a new stabilized $Q_1 \times Q_1$ combination for the velocity and pressure with three internal degrees of freedom added to the velocity space, that is, one degree of freedom

for each component of the velocity and one degree of freedom shared by both components of the velocity.”

Two versions are proposed, if I understand it correctly. The first one is given in Eq. (7) (three extra dofs: u_5, v_5, w):

$$\begin{aligned} u^h(r, s) &= \sum_{i=1}^4 N_i(r, s) u_i + \left[u_5 - \frac{w}{4}(1-s) \right] (1-r^2)(1-s^2) \\ v^h(r, s) &= \sum_{i=1}^4 N_i(r, s) v_i + \left[v_5 - \frac{w}{4}(1-r) \right] (1-r^2)(1-s^2) \end{aligned} \quad (5.129)$$

The second one in Eq. (23) (four extra dofs: u_5, v_5, u_6, v_6):

$$\begin{aligned} u^h(r, s) &= \sum_{i=1}^4 N_i(r, s) u_i + [u_5 + u_6(r+s)] (1-r^2)(1-s^2) \\ v^h(r, s) &= \sum_{i=1}^4 N_i(r, s) v_i + [v_5 + v_6(r+s)] (1-r^2)(1-s^2) \end{aligned} \quad (5.130)$$

- In Franca, Oliveira, and Sarkis [411] (2007): ”Stabilized finite element method for Stokes equations with piecewise continuous bilinear approximations for both velocity and pressure variables. The velocity field is enriched with piecewise polynomial bubble functions with null average at element edges.”

It looks like they are proposing (see their Eq. (2.6)):

$$\begin{aligned} u^h(r, s) &= \sum_{i=1}^4 N_i(r, s) u_i + (\alpha + \gamma s) \frac{1}{2} (r^2 + s^2 - \frac{4}{3}) \\ v^h(r, s) &= \sum_{i=1}^4 N_i(r, s) v_i + (\beta + \gamma r) \frac{1}{2} (r^2 + s^2 - \frac{4}{3}) \end{aligned} \quad (5.131)$$

- In Kwon and Park [738] (2014): ”We introduce a new stable MINI-element pair for incompressible Stokes equations on quadrilateral meshes, which uses the smallest number of bubbles for the velocity. The pressure is discretized with the P_1 -midpoint-edge-continuous elements and each component of the velocity field is done with the standard Q_1 -conforming elements enriched by one bubble a quadrilateral.”
- In Lamichhane [741] (2017): ”We consider a quadrilateral MINI finite element for approximating the solution of Stokes equations using a quadrilateral mesh. We use the standard bilinear finite element space enriched with element-wise defined bubble functions for the velocity and the standard bilinear finite element space for the pressure space. With a simple modification of the standard bubble function we show that a single bubble function is sufficient to ensure the inf-sup condition. This is a refinement of Bai (1997) [38] where the author enriches the velocity space with more than a single vector bubble function per element. In this article we show that with a small modification of the standard bubble function we can get the stability just by using a single vector bubble function per element.”

The two bubble functions are defined on the reference element $[-1, 1] \times [-1, 1]$:

$$b^{(1)}(r, s) = (1 - r)(1 - s) \cdot (1 - r^2)(1 - s^2) \quad (5.132)$$

$$b^{(2)}(r, s) = \left(1 + \frac{r + s}{4}\right) \cdot (1 - r^2)(1 - s^2) \quad (5.133)$$

Both bubble functions are exactly one in the middle of the element and exactly zero on the edges of the element as expected from basis functions.

We then have

$$\begin{aligned} \frac{\partial b^{(1)}}{\partial r}(r, s) &= (1 - s)^2(1 + s)[-2(1 - r)(1 + r) + (1 - r)^2] \\ &= (1 - s)^2(1 + s)[-2 + 2r^2 + 1 - 2r + r^2] \\ &= (1 - s)^2(1 + s)[-1 - 2r + 3r^2] \end{aligned} \quad (5.134)$$

$$\frac{\partial b^{(1)}}{\partial s}(r, s) = (1 - r)^2(1 + r)[-1 - 2s + 3s^2] \quad (5.135)$$

$$\begin{aligned} \frac{\partial b^{(2)}}{\partial r}(r, s) &= \frac{1}{4}(1 - s^2)(1 - r^2 + (4 + r + s)(-2r)) \\ &= \frac{1}{4}(1 - s^2)(1 - 8r - 3r^2 - 2rs) \end{aligned} \quad (5.136)$$

$$\begin{aligned} \frac{\partial b^{(2)}}{\partial s}(r, s) &= \frac{1}{4}(1 - r^2)(1 - s^2 + (4 + r + s)(-2s)) \\ &= \frac{1}{4}(1 - r^2)(1 - 8s - 3s^2 - 2rs) \end{aligned} \quad (5.137)$$

We postulate that a function f has the following representation in the element:

$$f^h(r, s) = a + br + cs + drs + e b(r, s)$$

where $b(r, s)$ stands for the bubble function which is of the form $b(r, s) = (1 - r^2)(1 - s^2)\phi(r, s)$ and ϕ is a (bi)-linear function of r, s .

We need

$$f^h(r_1, s_1) = a - b - c + d = f_1 \quad (5.138)$$

$$f^h(r_2, s_2) = a + b - c - d = f_2 \quad (5.139)$$

$$f^h(r_3, s_3) = a + b + c + d = f_3 \quad (5.140)$$

$$f^h(r_4, s_4) = a - b + c - d = f_4 \quad (5.141)$$

$$f^h(r_5, s_5) = a + e = f_5 \quad (5.142)$$

This can be written as a linear system:

$$\begin{pmatrix} 1 & -1 & -1 & 1 & 0 \\ 1 & 1 & -1 & -1 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & -1 & 1 & -1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \\ e \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \end{pmatrix}$$

and the solution is then:

$$\begin{pmatrix} a \\ b \\ c \\ d \\ e \end{pmatrix} = \frac{1}{4} \begin{pmatrix} 1 & 1 & 1 & 1 & 0 \\ -1 & 1 & 1 & -1 & 0 \\ -1 & -1 & 1 & 1 & 0 \\ 1 & -1 & 1 & -1 & 0 \\ -1 & -1 & -1 & -1 & 4 \end{pmatrix} \cdot \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \end{pmatrix}$$

or,

$$\begin{aligned}
a &= \frac{1}{4}(f_1 + f_2 + f_3 + f_4) \\
b &= \frac{1}{4}(-f_1 + f_2 + f_3 - f_4) \\
c &= \frac{1}{4}(-f_1 - f_2 + f_3 + f_4) \\
d &= \frac{1}{4}(f_1 - f_2 + f_3 - f_4) \\
e &= \frac{1}{4}(-f_1 - f_2 - f_3 - f_4 + 4f_5)
\end{aligned} \tag{5.143}$$

Then

$$\begin{aligned}
4f^h(r, s) &= 4[a + br + cs + drs + e(1 - r^2)(1 - s^2)\phi(r, s)] \\
&= (f_1 + f_2 + f_3 + f_4) \\
&\quad + (-f_1 + f_2 + f_3 - f_4)r \\
&\quad + (-f_1 - f_2 + f_3 + f_4)s \\
&\quad + (f_1 - f_2 + f_3 - f_4)rs \\
&\quad + (-f_1 - f_2 - f_3 - f_4 + 4f_5)(1 - r^2)(1 - s^2)\phi(r, s) \\
&= (1 - r - s + rs - b(r, s))f_1 \\
&\quad + (1 + r - s - rs - b(r, s))f_2 \\
&\quad + (1 + r + s + rs - b(r, s))f_3 \\
&\quad + (1 - r + s - rs - b(r, s))f_4 \\
&\quad + 4b(r, s)f_5
\end{aligned} \tag{5.144}$$

or,

$$\begin{aligned}
f^h(r, s) &= \underbrace{\left(\frac{1}{4}(1 - r)(1 - s) - \frac{1}{4}b(r, s)\right)}_{\mathcal{N}_1} f_1 + \underbrace{\left(\frac{1}{4}(1 + r)(1 - s) - \frac{1}{4}b(r, s)\right)}_{\mathcal{N}_2} f_2 \\
&\quad + \underbrace{\left(\frac{1}{4}(1 + r)(1 + s) - \frac{1}{4}b(r, s)\right)}_{\mathcal{N}_3} f_3 + \underbrace{\left(\frac{1}{4}(1 - r)(1 + s) - \frac{1}{4}b(r, s)\right)}_{\mathcal{N}_4} f_4 \\
&\quad + \underbrace{b(r, s)}_{\mathcal{N}_5} f_5
\end{aligned} \tag{5.145}$$

As in the P_1^+ case the resulting basis functions are a combination of the regular Q_1 basis functions and the bubble.

– Zeroth-order consistency check $f(r, s) = C$:

$$f^h(r, s) = \sum_{i=1}^5 \mathcal{N}_i(r, s) f_i = C \sum_{i=1}^5 \mathcal{N}_i(r, s) = C \tag{5.146}$$

- First-order consistency check $f(r, s) = r$ (or $f(r, s) = s$):

$$\begin{aligned}
f^h(r, s) &= \sum_{i=1}^5 \mathcal{N}_i(r, s) f_i \\
&= \mathcal{N}_1(r, s)(-1) + \mathcal{N}_2(r, s)(+1) + \mathcal{N}_3(r, s)(+1) + \mathcal{N}_4(r, s)(-1) + \mathcal{N}_5(r, s)(0) \\
&= -\mathcal{N}_1(r, s) + \mathcal{N}_2(r, s) + \mathcal{N}_3(r, s) - \mathcal{N}_4(r, s) \\
&= r
\end{aligned} \tag{5.147}$$

- Second-order consistency check $f(r, s) = rs$ ($f_1 = (-1)(-1) = 1$, $f_2 = (+1)(-1) = -1$, etc ...)

$$\begin{aligned}
f^h(r, s) &= \sum_{i=1}^5 \mathcal{N}_i(r, s) f_i \\
&= \mathcal{N}_1(r, s)(+1) + \mathcal{N}_2(r, s)(-1) + \mathcal{N}_3(r, s)(+1) + \mathcal{N}_4(r, s)(-1) + \mathcal{N}_5(r, s)(0) \\
&= \mathcal{N}_1 - \mathcal{N}_2 + \mathcal{N}_3 - \mathcal{N}_4 \\
&= \left(\frac{1}{4}(1-r)(1-s) - \frac{1}{4}b(r, s) \right) - \left(\frac{1}{4}(1+r)(1-s) - \frac{1}{4}b(r, s) \right) \\
&+ \left(\frac{1}{4}(1+r)(1+s) - \frac{1}{4}b(r, s) \right) - \left(\frac{1}{4}(1-r)(1+s) - \frac{1}{4}b(r, s) \right) \\
&= \frac{1}{4}(1-r)(1-s) - \frac{1}{4}(1+r)(1-s) + \frac{1}{4}(1+r)(1+s) - \frac{1}{4}(1-r)(1+s) \\
&= \frac{1}{2}(-r)(1-s) + \frac{1}{2}(+r)(1+s) \\
&= rs
\end{aligned} \tag{5.148}$$

We find that the basis functions can represent a bilinear field exactly.

Consistency check for quadratic terms, i.e. $f(r, s) = r^2$ (or $f(r, s) = s^2$):

$$\begin{aligned}
f^h(r, s) &= \sum_{i=1}^5 \mathcal{N}_i(r, s) f_i \\
&= \mathcal{N}_1(r, s) \cdot (+1) + \mathcal{N}_2(r, s) \cdot (+1) + \mathcal{N}_3(r, s) \cdot (+1) + \mathcal{N}_4(r, s) \cdot (+1) + \mathcal{N}_5(r, s) \cdot (0) \\
&= \left(\frac{1}{4}(1-r)(1-s) - \frac{1}{4}b(r, s) \right) + \left(\frac{1}{4}(1+r)(1-s) - \frac{1}{4}b(r, s) \right) \\
&+ \left(\frac{1}{4}(1+r)(1+s) - \frac{1}{4}b(r, s) \right) + \left(\frac{1}{4}(1-r)(1+s) - \frac{1}{4}b(r, s) \right) \\
&= \frac{1}{2}(1-s) + \frac{1}{2}(1+s) - b(r, s) \\
&= 1 - b(r, s)
\end{aligned} \tag{5.149}$$

We have

$$\begin{aligned}
\int_{-1}^{+1} \int_{-1}^{+1} (1 - b_1(r, s)) dr ds &= \int_{-1}^{+1} \int_{-1}^{+1} [1 - (1-r^2)(1-s^2)(1-r)(1-s)] dr ds = 20/9 \simeq 2.22 \\
\int_{-1}^{+1} \int_{-1}^{+1} (1 - b_2(r, s, \beta)) dr ds &= \int_{-1}^{+1} \int_{-1}^{+1} [1 - (1-r^2)(1-s^2)(1+\beta(r+s))] dr ds = 20/9 \quad \forall \beta
\end{aligned}$$

Both bubbles yield the same average. This is not helpful.

Let us now look at the (root) mean square:

$$\begin{aligned} \int_{-1}^{+1} \int_{-1}^{+1} (1 - b_1(r, s))^2 dr ds &= 21284/11025 \simeq 1.93052 \\ \int_{-1}^{+1} \int_{-1}^{+1} (1 - b_2(r, s, \beta))^2 dr ds &= \frac{4}{1575} (128\beta^2 + 623) \end{aligned} \quad (5.150)$$

The problem is that the minimum is reached for $\beta = 0$ which is not allowed so we cannot choose β so as to minimise the error. For $\beta = 0.25$ as used in the paper:

$$\int_{-1}^{+1} \int_{-1}^{+1} (1 - b_2(r, s))^2 dr ds = 2524/1575 \simeq 1.60254$$

On the other hand, this means that using the second bubble function does a better job at representing square terms (r^2 , s^2) than using the first one.

One can also revisit the second bubble function: in Lamichhane (2017) [741] it is postulated to be defined by

$$b^{(2)}(r, s) = (a + br + cs)(1 - r^2)(1 - s^2) \quad abc \neq 0 \quad (5.151)$$

on the reference element $[-1, 1] \times [-1, 1]$. Then the author states that 'for simplicity we choose':

$$b^{(2)}(r, s) = \frac{1}{4}(4 + r + s)(1 - r^2)(1 - s^2) \quad (5.152)$$

and that 'the factor 1/4 is used to force the value of the bubble function at the centroid of the square to be 1'.

Looking closer, we see that forcing the bubble to be 1 in $(r, s) = (0, 0)$ does impose $a = 1$ but leaves b, c free, i.e. the bubble is then:

$$b^{(2)}(r, s) = (1 + br + cs)(1 - r^2)(1 - s^2) \quad bc \neq 0 \quad (5.153)$$

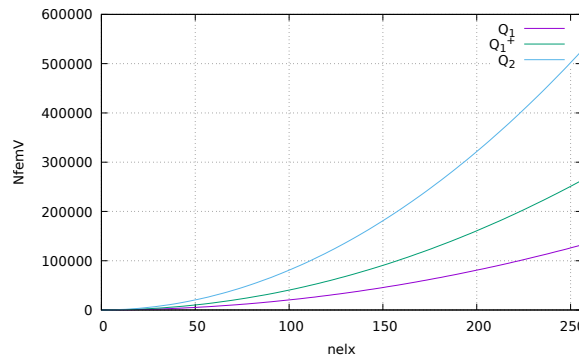
For symmetry reasons I would be tempted to indeed take $b = c$ but I am then left with

$$b^{(2)}(r, s) = [1 + b(r + s)](1 - r^2)(1 - s^2) \quad b \neq 0 \quad (5.154)$$

which means that Lamichhane sets $b = c = 1/4$ in his paper.

Question: We know that $b = 0$ is not allowed, but could it not be possible to design an analytical or numerical test or a theory to choose an 'optimal' value (in some sense) for b ?

Let us consider a square mesh with $nelx^2$ elements for simplicity. The number of V dofs for a Q_1 space would be $(nelx + 1)^2 = nelx^2 + 2nelx + 1$. The number of V dofs for a Q_1^+ space would be $(nelx + 1)^2 + nelx^2 = 2nelx^2 + 2nelx + 1$. The number of V dofs for a Q_2 space would be $(2 * nelx + 1)^2 = 4nelx^2 + 4nelx + 1$. Asymptotically, for large values of $nelx$, we find that a Q_1^+ space requires twice as many dofs as Q_1 while Q_2 requires 4 times as many.



Relevant Literature:

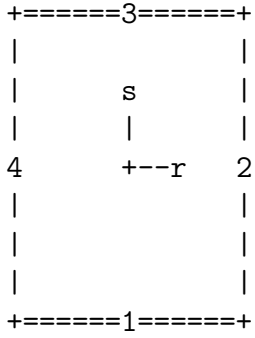
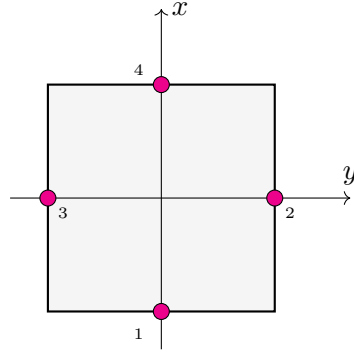
- Mons & Roge (1992) [892],
- Li *et al.* (2009) [781],
- Knobloch & Tobiska (2000) [714],
- Franca *et al.* (1993) [409],
- Idelsohn *et al.* (1995) [619].

5.3.15 The rotated Q_1 (Rannacher-Turek element)

basis_q1rc_2D.tex

The nodes are not on the corners of the element but in the middle of the element edges:

(tikz_RTQ1P0.tex)



There are two types of basis functions: the Middle Point (MP) variant such that $\mathcal{N}_i(\vec{r}_j) = \delta_{ij}$ and the Mid Value (MV) variant such that $\frac{1}{|\Gamma_i|} \int_{\Gamma_i} \mathcal{N}_j d\Gamma = \delta_{ij}$.

The Middle Point (MP) variant . We have $\tilde{Q}_1 = \text{span}\{1, r, s, r^2 - s^2\}$ so a function $f \in \tilde{Q}_1$ is such that

$$f(r, s) = a + br + cs + d(r^2 - s^2) \quad (5.155)$$

This function must be so that

$$f_1 = f(r = 0, s = -1) = a - c - d \quad (5.156)$$

$$f_2 = f(r = +1, s = 0) = a + b + d \quad (5.157)$$

$$f_3 = f(r = 0, s = +1) = a + c - d \quad (5.158)$$

$$f_4 = f(r = -1, s = 0) = a - b + d \quad (5.159)$$

and then

$$\begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{pmatrix} = \begin{pmatrix} 1 & 0 & -1 & -1 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & -1 \\ 1 & -1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix}$$

This system can easily be solved, and a, b, c, d are then replaced in Eq. (5.155), which yields

$$f(r, s) = \mathcal{N}_1(r, s)f_1 + \mathcal{N}_2(r, s)f_2 + \mathcal{N}_3(r, s)f_3 + \mathcal{N}_4(r, s)f_4 \quad (5.160)$$

inside the element with

$$\begin{aligned} \mathcal{N}_1(r, s) &= \frac{1}{4}(1 - 2s - (r^2 - s^2)) \\ \mathcal{N}_2(r, s) &= \frac{1}{4}(1 + 2r + (r^2 - s^2)) \\ \mathcal{N}_3(r, s) &= \frac{1}{4}(1 + 2s - (r^2 - s^2)) \\ \mathcal{N}_4(r, s) &= \frac{1}{4}(1 - 2r + (r^2 - s^2)) \end{aligned}$$

We of course recover the partition of unity property, i.e. $\sum \mathcal{N}_i(r, s) = 1$ for any coordinate r, s inside the reference element.

Remark. *These basis functions have been independently proposed by Donea, Giuliani, Morgan, and Quartapelle [340] (1981). The authors prove herein that this element is checkerboard-free (although they do not show any example of simulation carried out with this element).*

$$\frac{\partial \mathcal{N}_1}{\partial r} = \frac{1}{2}(-r) \quad (5.161)$$

$$\frac{\partial \mathcal{N}_2}{\partial r} = \frac{1}{2}(1 + r) \quad (5.162)$$

$$\frac{\partial \mathcal{N}_3}{\partial r} = \frac{1}{2}(-r) \quad (5.163)$$

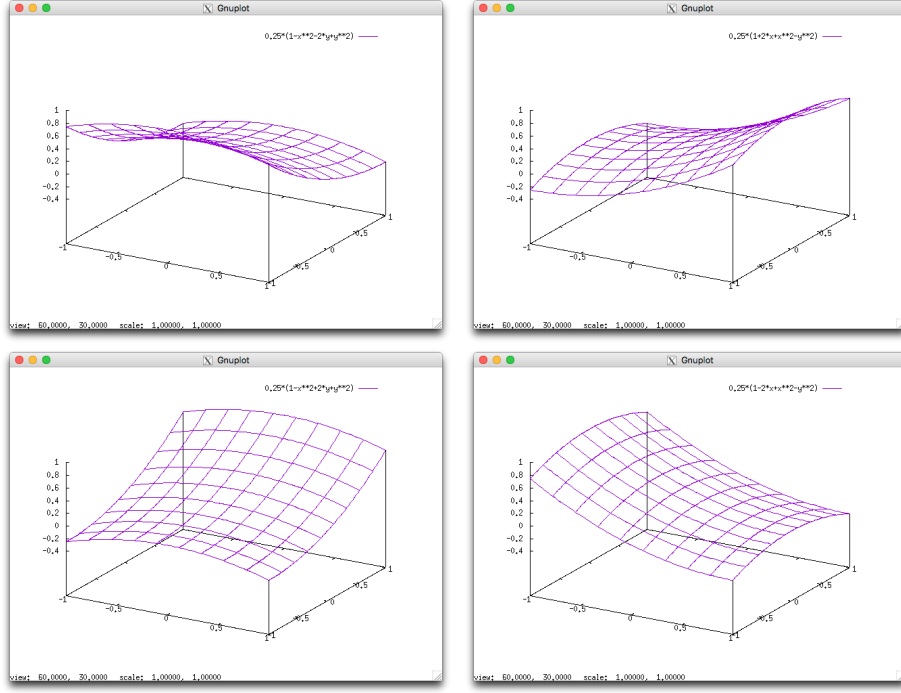
$$\frac{\partial \mathcal{N}_4}{\partial r} = \frac{1}{2}(-1 + r) \quad (5.164)$$

$$\frac{\partial \mathcal{N}_1}{\partial s} = \frac{1}{2}(-1 + s) \quad (5.165)$$

$$\frac{\partial \mathcal{N}_2}{\partial s} = \frac{1}{2}(-s) \quad (5.166)$$

$$\frac{\partial \mathcal{N}_3}{\partial s} = \frac{1}{2}(1 + s) \quad (5.167)$$

$$\frac{\partial \mathcal{N}_4}{\partial s} = \frac{1}{2}(-s) \quad (5.168)$$



Graphical representation of the \bar{Q}_1 basis functions

The Mid Value (MV) variant .

These basis functions are implemented in deal.II ¹⁰ for $x \in [0, 1]$ and $y \in [0, 1]$:

$$\mathcal{N}_1(x, y) = 0.75 + 1.5x - 2.5y - 1.5(x^2 - y^2) \quad \text{bottom} \quad (5.169)$$

$$\mathcal{N}_2(x, y) = -0.25 - 0.5x + 1.5y + 1.5(x^2 - y^2) \quad \text{right} \quad (5.170)$$

$$\mathcal{N}_3(x, y) = -0.25 + 1.5x - 0.5y - 1.5(x^2 - y^2) \quad \text{top} \quad (5.171)$$

$$\mathcal{N}_4(x, y) = 0.75 - 2.5x + 1.5y + 1.5(x^2 - y^2) \quad \text{left} \quad (5.172)$$

We then proceed to rewrite these for $r \in [-1, 1]$ and $t \in [-1, 1]$:

$$\mathcal{N}_1(r, s) = \frac{1}{4} - \frac{1}{2}s - \frac{3}{8}(r^2 - s^2) \quad \text{bottom} \quad (5.173)$$

$$\mathcal{N}_2(r, s) = \frac{1}{4} + \frac{1}{2}r + \frac{3}{8}(r^2 - s^2) \quad \text{right} \quad (5.174)$$

$$\mathcal{N}_3(r, s) = \frac{1}{4} + \frac{1}{2}s - \frac{3}{8}(r^2 - s^2) \quad \text{top} \quad (5.175)$$

$$\mathcal{N}_4(r, s) = \frac{1}{4} - \frac{1}{2}r + \frac{3}{8}(r^2 - s^2) \quad \text{left} \quad (5.176)$$

It is easy to verify that these functions verify the property

$$\frac{1}{|\Gamma_i|} \int_{\Gamma_i} N_j d\Gamma = \delta_{ij}$$

These basis functions are used in Shipeng & Zhongci (2006) [1162] and mentioned in John [650,

¹⁰https://www.dealii.org/8.5.0/doxygen/deal.II/polynomials_rannacher_turek_8cc_source.html

p.722].

$$\begin{aligned}\frac{\partial N_1}{\partial r} &= -\frac{3}{4}r \\ \frac{\partial N_2}{\partial r} &= \frac{1}{2} + \frac{3}{4}r \\ \frac{\partial N_3}{\partial r} &= -\frac{3}{4}r \\ \frac{\partial N_4}{\partial r} &= -\frac{1}{2} + \frac{3}{4}r\end{aligned}$$

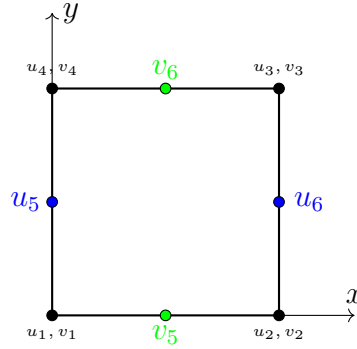
$$\begin{aligned}\frac{\partial N_1}{\partial t} &= -\frac{1}{2} + \frac{3}{4}t \\ \frac{\partial N_2}{\partial t} &= -\frac{3}{4}t \\ \frac{\partial N_3}{\partial t} &= \frac{1}{2} + \frac{3}{4}t \\ \frac{\partial N_4}{\partial t} &= -\frac{3}{4}t\end{aligned}$$

5.3.16 The 2D enriched $Q_1^+ \times P_0$ of Fortin

basis_q1fortin_2D.tex

We here consider the enriched $Q_1 \times P_0$ element introduced first by Fortin (1981) [401]. The layout of the degrees of freedom is as follows:

(tikz_q1pp02D.tex)



The approximation of the velocity components u and v inside the element is

$$u^h(r, s) = a^u \mathcal{N}_1(r, s) + b^u \mathcal{N}_2(r, s) + c^u \mathcal{N}_3(r, s) + d^u \mathcal{N}_4(r, s) + d b_5^u(r, s) + e b_6^u(r, s)$$

$$v^h(r, s) = a^v \mathcal{N}_1(r, s) + b^v \mathcal{N}_2(r, s) + c^v \mathcal{N}_3(r, s) + d^v \mathcal{N}_4(r, s) + d^v b_5^v(r, s) + e^v b_6^v(r, s)$$

where $\mathcal{N}_{1,2,3,4}$ are the standard Q_1 basis functions in 2D and with

$$b_5^u(r, s) = \frac{1}{2}(1-r)(1-s^2) \quad b_6^u(r, s) = \frac{1}{2}(1+r)(1-s^2)$$

and

$$b_5^v(r, s) = \frac{1}{2}(1-r^2)(1-s) \quad b_6^v(r, s) = \frac{1}{2}(1-r^2)(1+s)$$

In the end one arrives at

$$\begin{aligned}
\mathcal{N}_1^u(r, s) &= \mathcal{N}_1(r, s) - \frac{1}{2}b_5^u(r, s) \\
\mathcal{N}_2^u(r, s) &= \mathcal{N}_2(r, s) - \frac{1}{2}b_6^u(r, s) \\
\mathcal{N}_3^u(r, s) &= \mathcal{N}_3(r, s) - \frac{1}{2}b_6^u(r, s) \\
\mathcal{N}_4^u(r, s) &= \mathcal{N}_4(r, s) - \frac{1}{2}b_5^u(r, s) \\
\mathcal{N}_5^u(r, s) &= b_5^u(r, s) \\
\mathcal{N}_6^u(r, s) &= b_6^u(r, s) \\
\\
\mathcal{N}_1^v(r, s) &= \mathcal{N}_1(r, s) - \frac{1}{2}b_5^v(r, s) \\
\mathcal{N}_2^v(r, s) &= \mathcal{N}_2(r, s) - \frac{1}{2}b_5^v(r, s) \\
\mathcal{N}_3^v(r, s) &= \mathcal{N}_3(r, s) - \frac{1}{2}b_6^v(r, s) \\
\mathcal{N}_4^v(r, s) &= \mathcal{N}_4(r, s) - \frac{1}{2}b_6^v(r, s) \\
\mathcal{N}_5^v(r, s) &= b_5^v(r, s) \\
\mathcal{N}_6^v(r, s) &= b_6^v(r, s)
\end{aligned} \tag{5.177}$$

We can check for the zero-th order consistency: Let $u(r, s) = C$, then

$$u^h(r, s) = \sum_{i=1}^6 \mathcal{N}_i^u(r, s) u_i = C \sum_{i=1}^6 \mathcal{N}_i^u(r, s) = C \sum_{i=1}^4 \mathcal{N}_i(r, s) = C \tag{5.178}$$

5.3.17 The P_1^{NC} space

p1nc.tex

| P1 | P1NC | |
|---------|---------|-------------------|
| 2 | . | (r0,s0)=(1/2,0) |
| \ | \ | |
| \ | 2 1 | (r1,s1)=(1/2,1/2) |
| \ | \ | |
| 0-----1 | .--0--. | (r2,s2)=(0,1/2) |

The basis functions are $1 - 2\lambda_i$, where λ_i are the barycentric coordinates, so we arrive at

$$\begin{aligned}
\mathcal{N}_0(r, s) &= 1 - 2\lambda_3 = 1 - 2s \\
\mathcal{N}_1(r, s) &= 1 - 2\lambda_1 = 1 - 2(1 - r - s) = -1 + 2r + 2s \\
\mathcal{N}_2(r, s) &= 1 - 2\lambda_2 = 1 - 2r
\end{aligned}$$

with of course

$$\mathcal{N}_0(r, s) + \mathcal{N}_1(r, s) + \mathcal{N}_2(r, s) = (1 - 2s) + (-1 + 2r + 2s) + (1 - 2r) = 1$$

We have

$$\begin{aligned}
\mathcal{N}_0(r_0, s_0) &= 1 - 2 \cdot 0 = 1 \\
\mathcal{N}_0(r_1, s_1) &= 1 - 2 \cdot 1/2 = 0 \\
\mathcal{N}_0(r_2, s_2) &= 1 - 2 \cdot 1/2 = 0 \\
\mathcal{N}_1(r_0, s_0) &= 1 - 2(1 - 1/2 - 0) = 0 \\
\mathcal{N}_1(r_1, s_1) &= 1 - 2(1 - 1/2 - 1/2) = 1 \\
\mathcal{N}_1(r_2, s_2) &= 1 - 2(1 - 0 - 1/2) = 0 \\
\mathcal{N}_2(r_0, s_0) &= 1 - 2 \cdot 1/2 = 0 \\
\mathcal{N}_2(r_1, s_1) &= 1 - 2 \cdot 1/2 = 0 \\
\mathcal{N}_2(r_2, s_2) &= 1 - 2 \cdot 0 = 1
\end{aligned}$$

Automatically,

$$\begin{aligned}
\partial_r \mathcal{N}_0(r, s) &= 0 \\
\partial_r \mathcal{N}_1(r, s) &= 2 \\
\partial_r \mathcal{N}_2(r, s) &= -2 \\
\partial_s \mathcal{N}_0(r, s) &= -2 \\
\partial_s \mathcal{N}_1(r, s) &= 2 \\
\partial_s \mathcal{N}_2(r, s) &= 0
\end{aligned}$$

Another way to obtain the basis functions is as follows.

$$f^h(r, s) = a + br + cs$$

$$\begin{aligned}
f_0 &= f^h(r_0, s_0) = a + b/2 \\
f_1 &= f^h(r_1, s_1) = a + b/2 + c/2 \\
f_2 &= f^h(r_2, s_2) = a + c/2
\end{aligned}$$

or,

$$\begin{pmatrix} 1 & 1/2 & 0 \\ 1 & 1/2 & 1/2 \\ 1 & 0 & 1/2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ f_2 \end{pmatrix}$$

i.e.

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1 & -1 & 1 \\ 0 & 2 & -2 \\ -2 & 2 & 0 \end{pmatrix} \cdot \begin{pmatrix} f_0 \\ f_1 \\ f_2 \end{pmatrix}$$

so

$$\begin{aligned}
f^h(r, s) &= a + br + cs \\
&= (f_0 - f_1 + f_2) + 2(f_1 - f_2)r + 2(-f_0 + f_1)s \\
&= (1 - 2s)f_0 + (-1 + 2r + 2s)f_1 + (1 - 2r)f_2 \\
&= \sum_{i=0}^2 \mathcal{N}_i(r, s) f_i
\end{aligned}$$

5.4 Elements and basis functions in 3D

elements3D.tex

5.4.1 Linear basis functions in tetrahedra (P_1)

basis_p1p_3D.tex

The location of the four nodes are:

$$\begin{aligned}(\mathbf{r}_0, \mathbf{s}_0) &= (0, 0, 0) \\(\mathbf{r}_1, \mathbf{s}_1) &= (1, 0, 0) \\(\mathbf{r}_2, \mathbf{s}_2) &= (0, 1, 0) \\(\mathbf{r}_3, \mathbf{s}_3) &= (0, 0, 1)\end{aligned}$$

The basis polynomial is given by

$$f(r, s, t) = c_0 + c_1 r + c_2 s + c_3 t$$

and it needs to satisfy these four conditions:

$$f_1 = f(r_1, s_1, t_1) = c_0 \quad (5.179)$$

$$f_2 = f(r_2, s_2, t_2) = c_0 + c_1 \quad (5.180)$$

$$f_3 = f(r_3, s_3, t_3) = c_0 + c_2 \quad (5.181)$$

$$f_4 = f(r_4, s_4, t_4) = c_0 + c_3 \quad (5.182)$$

which yields:

$$c_0 = f_1 \quad c_1 = f_2 - f_1 \quad c_2 = f_3 - f_1 \quad c_3 = f_4 - f_1$$

Finally,

$$\begin{aligned}f(r, s, t) &= c_0 + c_1 r + c_2 s + c_3 t \\&= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + (f_4 - f_1)t \\&= f_1(1 - r - s - t) + f_2 r + f_3 s + f_4 t \\&= \sum_i \mathcal{N}_i(r, s, t) f_i\end{aligned}$$

Finally,

$$\begin{aligned}\mathcal{N}_1(r, s, t) &= 1 - r - s - t \\ \mathcal{N}_2(r, s, t) &= r \\ \mathcal{N}_3(r, s, t) &= s \\ \mathcal{N}_4(r, s, t) &= t\end{aligned}$$

Derivatives are trivial to obtain.

5.4.2 Enriched linear in tetrahedra(P_1^+)

basis_p1p_3D.tex

These basis functions would be used in the MINI element, see Section 7.3.14.

In 3D the bubble function looks like $rst(1 - r - s - t)$ so that

$$f(r, s, t) = a + b r + c s + d t + e rst(1 - r - s - t)$$

We have node 1 at location $(r, s, t) = (0, 0, 0)$, node 2 at $(r, s, t) = (1, 0, 0)$, node 3 at $(r, s, t) = (0, 1, 0)$, node 4 at $(r, s, t) = (0, 0, 1)$ and we set the location of the bubble (node 5) at $r = s = t = 1/4$ so that

$$\begin{aligned} f(r_1, s_1, t_1) &= f_1 = a + b r_1 + c s_1 + d t_1 + e r_1 s_1 t_1 (1 - r_1 - s_1 - t_1) \\ f(r_2, s_2, t_2) &= f_2 = a + b r_2 + c s_2 + d t_2 + e r_2 s_2 t_2 (1 - r_2 - s_2 - t_2) \\ f(r_3, s_3, t_3) &= f_3 = a + b r_3 + c s_3 + d t_3 + e r_3 s_3 t_3 (1 - r_3 - s_3 - t_3) \\ f(r_4, s_4, t_4) &= f_4 = a + b r_4 + c s_4 + d t_4 + e r_4 s_4 t_4 (1 - r_4 - s_4 - t_4) \\ f(r_5, s_5, t_5) &= f_5 = a + b r_5 + c s_5 + d t_5 + e r_5 s_5 t_5 (1 - r_5 - s_5 - t_5) \end{aligned} \quad (5.183)$$

i.e.,

$$\begin{aligned} f_1 &= a \\ f_2 &= a + b \\ f_3 &= a + c \\ f_4 &= a + d \\ f_5 &= a + b/4 + c/4 + d/4 + e/64(1 - 1/4 - 1/4 - 1/4) \\ &= a + b/4 + c/4 + d/4 + e/256 \end{aligned}$$

Then

$$\begin{aligned} a &= f_1 \\ b &= f_2 - f_1 \\ c &= f_3 - f_1 \\ d &= f_4 - f_1 \\ e &= 256(f_5 - a - b/4 - c/4 - d/4) \\ &= 256(f_5 - f_1 - (f_2 - f_1)/4 - (f_3 - f_1)/4 - (f_4 - f_1)/4) \\ &= 256(-f_1/4 - f_2/4 - f_3/4 - f_4/4 + f_5) \\ &= 64(-f_1 - f_2 - f_3 - f_4 + 4f_5) \end{aligned} \quad (5.184)$$

Finally:

$$\begin{aligned} f(r, s, t) &= a + br + cs + dt + erst(1 - r - s - t) \\ &= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + (f_4 - f_1)t + 64(-f_1 - f_2 - f_3 - f_4 + 4f_5)rst(1 - r - s - t) \\ &= f_1[1 - r - s - t - 64rst(1 - r - s - t)] \\ &+ f_2[r - 64rst(1 - r - s - t)] \\ &+ f_3[s - 64rst(1 - r - s - t)] \\ &+ f_4[t - 64rst(1 - r - s - t)] \\ &+ f_5[256rst(1 - r - s - t)] \\ &= \sum_{i=1}^5 \mathcal{N}_i(r, s, t) f_i \end{aligned} \quad (5.185)$$

with

$$\mathcal{N}_1(r, s, t) = 1 - r - s - t - 64rst(1 - r - s - t) \quad (5.186)$$

$$\mathcal{N}_2(r, s, t) = r - 64rst(1 - r - s - t) \quad (5.187)$$

$$\mathcal{N}_3(r, s, t) = s - 64rst(1 - r - s - t) \quad (5.188)$$

$$\mathcal{N}_4(r, s, t) = t - 64rst(1 - r - s - t) \quad (5.189)$$

$$\mathcal{N}_5(r, s, t) = +256rst(1 - r - s - t) \quad (5.190)$$

The derivatives are given by:

$$\frac{\partial \mathcal{N}_1}{\partial r}(r, s, t) = -1 - 64st(1 - 2r - s - t)$$

$$\frac{\partial \mathcal{N}_2}{\partial r}(r, s, t) = +1 - 64st(1 - 2r - s - t)$$

$$\frac{\partial \mathcal{N}_3}{\partial r}(r, s, t) = -64st(1 - 2r - s - t)$$

$$\frac{\partial \mathcal{N}_4}{\partial r}(r, s, t) = -64st(1 - 2r - s - t)$$

$$\frac{\partial \mathcal{N}_5}{\partial r}(r, s, t) = 256st(1 - 2r - s - t)$$

$$\frac{\partial \mathcal{N}_1}{\partial s}(r, s, t) = -1 - 64rt(1 - r - 2s - t)$$

$$\frac{\partial \mathcal{N}_2}{\partial s}(r, s, t) = -64rt(1 - r - 2s - t)$$

$$\frac{\partial \mathcal{N}_3}{\partial s}(r, s, t) = +1 - 64rt(1 - r - 2s - t)$$

$$\frac{\partial \mathcal{N}_4}{\partial s}(r, s, t) = -64rt(1 - r - 2s - t)$$

$$\frac{\partial \mathcal{N}_5}{\partial s}(r, s, t) = 256rt(1 - r - 2s - t)$$

$$\frac{\partial \mathcal{N}_1}{\partial t}(r, s, t) = -1 - 64rs(1 - r - s - 2t)$$

$$\frac{\partial \mathcal{N}_2}{\partial t}(r, s, t) = -64rs(1 - r - s - 2t)$$

$$\frac{\partial \mathcal{N}_3}{\partial t}(r, s, t) = -64rs(1 - r - s - 2t)$$

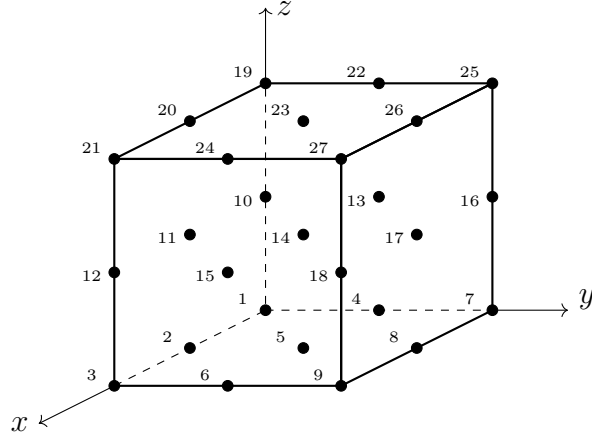
$$\frac{\partial \mathcal{N}_4}{\partial t}(r, s, t) = +1 - 64rs(1 - r - s - 2t)$$

$$\frac{\partial \mathcal{N}_5}{\partial t}(r, s, t) = 256rs(1 - r - s - 2t)$$

5.4.3 Triquadratic basis functions in 3D (Q_2)

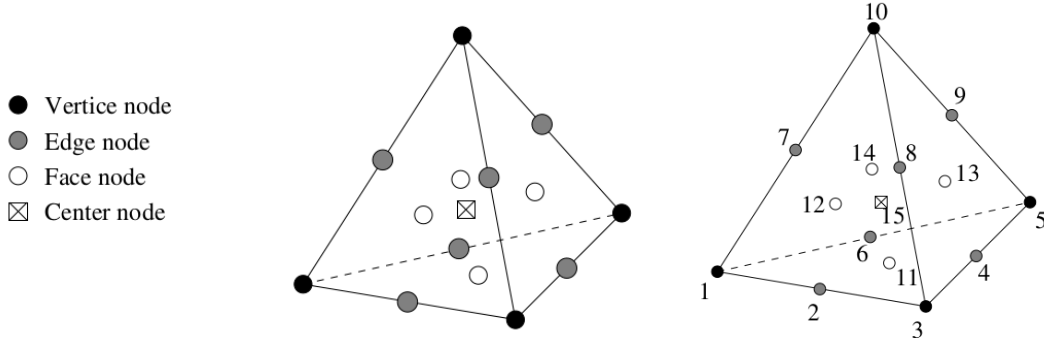
basis.q2_3D.tex

(tikz.q2.tex)



$$\begin{aligned}
\mathcal{N}_1 &= 0.5r(r-1) \cdot 0.5s(s-1) \cdot 0.5t(t-1) \\
\mathcal{N}_2 &= (1-r^2) \cdot 0.5s(s-1) \cdot 0.5t(t-1) \\
\mathcal{N}_3 &= 0.5r(r+1) \cdot 0.5s(s-1) \cdot 0.5t(t-1) \\
\mathcal{N}_4 &= 0.5r(r-1) \cdot (1-s^2) \cdot 0.5t(t-1) \\
\mathcal{N}_5 &= (1-r^2) \cdot (1-s^2) \cdot 0.5t(t-1) \\
\mathcal{N}_6 &= 0.5r(r+1) \cdot (1-s^2) \cdot 0.5t(t-1) \\
\mathcal{N}_7 &= 0.5r(r-1) \cdot 0.5s(s+1) \cdot 0.5t(t-1) \\
\mathcal{N}_8 &= (1-r^2) \cdot 0.5s(s+1) \cdot 0.5t(t-1) \\
\mathcal{N}_9 &= 0.5r(r+1) \cdot 0.5s(s+1) \cdot 0.5t(t-1) \\
\mathcal{N}_{10} &= 0.5r(r-1) \cdot 0.5s(s-1) \cdot (1-t^2) \\
\mathcal{N}_{11} &= (1-r^2) \cdot 0.5s(s-1) \cdot (1-t^2) \\
\mathcal{N}_{12} &= 0.5r(r+1) \cdot 0.5s(s-1) \cdot (1-t^2) \\
\mathcal{N}_{13} &= 0.5r(r-1) \cdot (1-s^2) \cdot (1-t^2) \\
\mathcal{N}_{14} &= (1-r^2) \cdot (1-s^2) \cdot (1-t^2) \\
\mathcal{N}_{15} &= 0.5r(r+1) \cdot (1-s^2) \cdot (1-t^2) \\
\mathcal{N}_{16} &= 0.5r(r-1) \cdot 0.5s(s+1) \cdot (1-t^2) \\
\mathcal{N}_{17} &= (1-r^2) \cdot 0.5s(s+1) \cdot (1-t^2) \\
\mathcal{N}_{18} &= 0.5r(r+1) \cdot 0.5s(s+1) \cdot (1-t^2) \\
\mathcal{N}_{19} &= 0.5r(r-1) \cdot 0.5s(s-1) \cdot 0.5t(t+1) \\
\mathcal{N}_{20} &= (1-r^2) \cdot 0.5s(s-1) \cdot 0.5t(t+1) \\
\mathcal{N}_{21} &= 0.5r(r+1) \cdot 0.5s(s-1) \cdot 0.5t(t+1) \\
\mathcal{N}_{22} &= 0.5r(r-1) \cdot (1-s^2) \cdot 0.5t(t+1) \\
\mathcal{N}_{23} &= (1-r^2) \cdot (1-s^2) \cdot 0.5t(t+1) \\
\mathcal{N}_{24} &= 0.5r(r+1) \cdot (1-s^2) \cdot 0.5t(t+1) \\
\mathcal{N}_{25} &= 0.5r(r-1) \cdot 0.5s(s+1) \cdot 0.5t(t+1) \\
\mathcal{N}_{26} &= (1-r^2) \cdot 0.5s(s+1) \cdot 0.5t(t+1) \\
\mathcal{N}_{27} &= 0.5r(r+1) \cdot 0.5s(s+1) \cdot 0.5t(t+1)
\end{aligned}$$

5.4.4 Enriched quadratic basis functions in tetrahedra (P_2^+)



The velocity basis functions are:

$$\phi_i = \lambda_i(2\lambda_i - 1) + 3(\lambda_i\lambda_j\lambda_k + \lambda_i\lambda_j\lambda_l + \lambda_i\lambda_k\lambda_l) - 4\lambda_i\lambda_j\lambda_k\lambda_l \quad (5.191)$$

$$\phi_{ij} = 4\lambda_i\lambda_j - 12(\lambda_i\lambda_j\lambda_k + \lambda_i\lambda_j\lambda_l) + 32\lambda_i\lambda_j\lambda_k\lambda_l \quad (5.192)$$

$$\phi_{ijk} = 27\lambda_i\lambda_j\lambda_k - 108\lambda_i\lambda_j\lambda_k\lambda_l \quad (5.193)$$

$$\phi_c = 256\lambda_i\lambda_j\lambda_k\lambda_l \quad (5.194)$$

REFS ??? better definition of functions !

5.4.5 Linear basis functions for hexahedra (P_1)

This is the $\mathbf{Q}_2 \times P_{-1}$ element. I choose the reduced coordinates of the pressure nodes to be :

| point | r | s | t |
|-------|------|------|------|
| 1 | 1/2 | -1/2 | -1/2 |
| 2 | -1/2 | 1/2 | -1/2 |
| 3 | -1/2 | -1/2 | 1/2 |
| 4 | 1/2 | 1/2 | 1/2 |

Inside the element the pressure is given as a linear function of the reduced coordinates r, s, t :

$$p(r, s, t) = a + br + cs + dt$$

This expression must exactly interpolate the pressure at all four pressure nodes:

$$p_1 = p(r_1, s_1, t_1) = a + br_1 + cs_1 + dt_1 = a + b/2 - c/2 - d/2$$

$$p_2 = p(r_2, s_2, t_2) = a + br_2 + cs_2 + dt_2 = a - b/2 + c/2 - d/2$$

$$p_3 = p(r_3, s_3, t_3) = a + br_3 + cs_3 + dt_3 = a - b/2 - c/2 + d/2$$

$$p_4 = p(r_4, s_4, t_4) = a + br_4 + cs_4 + dt_4 = a + b/2 + c/2 + d/2$$

or,

$$\begin{pmatrix} 1 & 1/2 & -1/2 & -1/2 \\ 1 & -1/2 & +1/2 & -1/2 \\ 1 & -1/2 & -1/2 & +1/2 \\ 1 & 1/2 & +1/2 & +1/2 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$$

The matrix is invertible and we get:

$$\begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} 1/4 & 1/4 & 1/4 & 1/4 \\ 1/2 & -1/2 & -1/2 & 1/2 \\ -1/2 & 1/2 & -1/2 & 1/2 \\ -1/2 & -1/2 & 1/2 & 1/2 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$$

so

$$\begin{aligned}
p(r, s, t) &= a + br + cs + dt \\
&= \frac{1}{4}(p_1 + p_2 + p_3 + p_4) + \frac{1}{2}(p_1 - p_2 - p_3 + p_4)r + \frac{1}{2}(-p_1 + p_2 - p_3 + p_4)s + \frac{1}{2}(-p_1 - p_2 + p_3 + p_4)t \\
&= \frac{1}{4}(1 + 2r - 2s - 2t)p_1 + \frac{1}{4}(1 - 2r + 2s - 2t)p_2 + \frac{1}{4}(1 - 2r - 2s + 2t)p_3 + \frac{1}{4}(1 + 2r + 2s + 2t)p_4 \\
&= \sum_{i=1}^4 N_i(r, s, t)p_i
\end{aligned} \tag{5.195}$$

with

$$\begin{aligned}
N_1(r, s, t) &= \frac{1}{4}(1 + 2r - 2s - 2t) \\
N_2(r, s, t) &= \frac{1}{4}(1 - 2r + 2s - 2t) \\
N_3(r, s, t) &= \frac{1}{4}(1 - 2r - 2s + 2t) \\
N_4(r, s, t) &= \frac{1}{4}(1 + 2r + 2s + 2t)
\end{aligned}$$

I could also have chosen

| point | r | s | t |
|-------|---|---|---|
| 1 | 0 | 0 | 0 |
| 2 | 1 | 0 | 0 |
| 3 | 0 | 1 | 0 |
| 4 | 0 | 0 | 1 |

This expression must exactly interpolate the pressure at all four pressure nodes:

$$\begin{aligned}
p_1 &= p(r_1, s_1, t_1) = a + br_1 + cs_1 + dt_1 = a \\
p_2 &= p(r_2, s_2, t_2) = a + br_2 + cs_2 + dt_2 = a + b \\
p_3 &= p(r_3, s_3, t_3) = a + br_3 + cs_3 + dt_3 = a + c \\
p_4 &= p(r_4, s_4, t_4) = a + br_4 + cs_4 + dt_4 = a + d
\end{aligned}$$

i.e.

$$a = p_1 \quad b = p_2 - p_1 \quad c = p_3 - p_1 \quad d = p_4 - p_1$$

or,

$$p^h(r, s) = a + br + cs + dt = p_1 + (p_2 - p_1)r + (p_3 - p_1)s + (p_4 - p_1)t = p_1(1 - r - s - t) + rp_2 + sp_3 + tp_4$$

so

$$N_1(r, s, t) = 1 - r - s - t \tag{5.196}$$

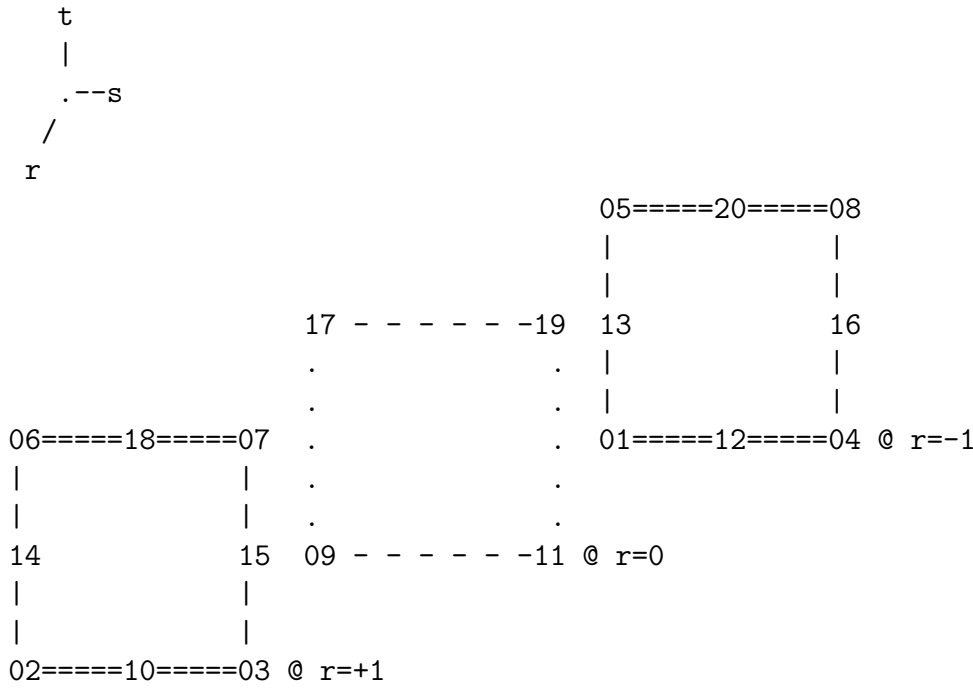
$$N_2(r, s, t) = r \tag{5.197}$$

$$N_3(r, s, t) = s \tag{5.198}$$

$$N_4(r, s, t) = t \tag{5.199}$$

5.4.6 20-node serendipity basis functions in 3D ($Q_2^{(20)}$)

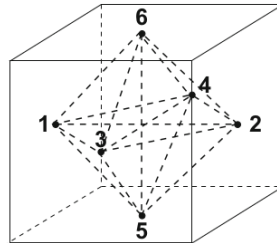
The serendipity elements are those rectangular elements which have no interior nodes [1051, p91].



find/build basis functions!

5.4.7 The rotated Q_1

The nodes are not on the corners of the element but in the middle of the element faces:



Node numbering and connectivity pattern of the reference element. Taken from [445]

We have $\tilde{Q}_1 = \text{span}\{1, r, s, t, r^2 - s^2, s^2 - t^2\}$.

The Middle Point (MP) variant

The basis functions are given by (see Georgiev *et al.* (2008) [445]):

$$N_1(r, s, t) = \frac{1}{6}(1 - 3r + 2r^2 - s^2 - t^2) \quad (5.200)$$

$$N_2(r, s, t) = \frac{1}{6}(1 + 3r + 2r^2 - s^2 - t^2) \quad (5.201)$$

$$N_3(r, s, t) = \frac{1}{6}(1 - r^2 - 3s + 2s^2 - t^2) \quad (5.202)$$

$$N_4(r, s, t) = \frac{1}{6}(1 - r^2 + 3s + 2s^2 - t^2) \quad (5.203)$$

$$N_5(r, s, t) = \frac{1}{6}(1 - r^2 - s^2 - 3t + 2t^2) \quad (5.204)$$

$$N_6(r, s, t) = \frac{1}{6}(1 - r^2 - s^2 + 3t + 2t^2) \quad (5.205)$$

$$\frac{\partial N_1}{\partial r} = \frac{1}{6}(-3 + 4r) \quad (5.206)$$

$$\frac{\partial N_2}{\partial r} = \frac{1}{6}(3 + 4r) \quad (5.207)$$

$$\frac{\partial N_3}{\partial r} = \frac{1}{6}(-2r) \quad (5.208)$$

$$\frac{\partial N_4}{\partial r} = \frac{1}{6}(-2r) \quad (5.209)$$

$$\frac{\partial N_5}{\partial r} = \frac{1}{6}(-2r) \quad (5.210)$$

$$\frac{\partial N_6}{\partial r} = \frac{1}{6}(-2r) \quad (5.211)$$

$$\frac{\partial N_1}{\partial s} = \frac{1}{6}(-2s) \quad (5.212)$$

$$\frac{\partial N_2}{\partial s} = \frac{1}{6}(-2s) \quad (5.213)$$

$$\frac{\partial N_3}{\partial s} = \frac{1}{6}(-3 + 4s) \quad (5.214)$$

$$\frac{\partial N_4}{\partial s} = \frac{1}{6}(3 + 4s) \quad (5.215)$$

$$\frac{\partial N_5}{\partial s} = \frac{1}{6}(-2s) \quad (5.216)$$

$$\frac{\partial N_6}{\partial s} = \frac{1}{6}(-2s) \quad (5.217)$$

$$\frac{\partial N_1}{\partial t} = \frac{1}{6}(-2t) \quad (5.218)$$

$$\frac{\partial N_2}{\partial t} = \frac{1}{6}(-2t) \quad (5.219)$$

$$\frac{\partial N_3}{\partial t} = \frac{1}{6}(-2t) \quad (5.220)$$

$$\frac{\partial N_4}{\partial t} = \frac{1}{6}(-2t) \quad (5.221)$$

$$\frac{\partial N_5}{\partial t} = \frac{1}{6}(-3 + 4t) \quad (5.222)$$

$$\frac{\partial N_6}{\partial t} = \frac{1}{6}(3 + 4t) \quad (5.223)$$

The Mid Value (MV) variant .

$$N_1(r, s, t) = \frac{1}{12}(2 - 6r + 6r^2 - 3s^2 - 3t^2) \quad (5.224)$$

$$N_2(r, s, t) = \frac{1}{12}(2 + 6r + 6r^2 - 3s^2 - 3t^2) \quad (5.225)$$

$$N_3(r, s, t) = \frac{1}{12}(2 - 3r^2 - 6s + 6s^2 - 3t^2) \quad (5.226)$$

$$N_4(r, s, t) = \frac{1}{12}(2 - 3r^2 + 6s + 6s^2 - 3t^2) \quad (5.227)$$

$$N_5(r, s, t) = \frac{1}{12}(2 - 3r^2 - 3s^2 - 6t + 6t^2) \quad (5.228)$$

$$N_6(r, s, t) = \frac{1}{12}(2 - 3r^2 - 3s^2 + 6t + 6t^2) \quad (5.229)$$

$$\frac{\partial N_1}{\partial r} = \frac{1}{12}(-6 + 12r) = \frac{1}{2}(-1 + 2r) \quad (5.230)$$

$$\frac{\partial N_2}{\partial r} = \frac{1}{12}(6 + 12r) = \frac{1}{2}(1 + 2r) \quad (5.231)$$

$$\frac{\partial N_3}{\partial r} = \frac{1}{12}(-6r) = -\frac{1}{2}r \quad (5.232)$$

$$\frac{\partial N_4}{\partial r} = \frac{1}{12}(-6r) = -\frac{1}{2}r \quad (5.233)$$

$$\frac{\partial N_5}{\partial r} = \frac{1}{12}(-6r) = -\frac{1}{2}r \quad (5.234)$$

$$\frac{\partial N_6}{\partial r} = \frac{1}{12}(-6r) = -\frac{1}{2}r \quad (5.235)$$

$$\frac{\partial N_1}{\partial s} = \frac{1}{12}(-6s) = -\frac{1}{2}s \quad (5.236)$$

$$\frac{\partial N_2}{\partial s} = \frac{1}{12}(-6s) = -\frac{1}{2}s \quad (5.237)$$

$$\frac{\partial N_3}{\partial s} = \frac{1}{12}(-6 + 12s) = \frac{1}{2}(-1 + 2s) \quad (5.238)$$

$$\frac{\partial N_4}{\partial s} = \frac{1}{12}(6 + 12s) = \frac{1}{2}(1 + 2s) \quad (5.239)$$

$$\frac{\partial N_5}{\partial s} = \frac{1}{12}(-6s) = -\frac{1}{2}s \quad (5.240)$$

$$\frac{\partial N_6}{\partial s} = \frac{1}{12}(-6s) = -\frac{1}{2}s \quad (5.241)$$

$$\frac{\partial N_1}{\partial t} = \frac{1}{12}(-6t) = -\frac{1}{2}t \quad (5.242)$$

$$\frac{\partial N_2}{\partial t} = \frac{1}{12}(-6t) = -\frac{1}{2}t \quad (5.243)$$

$$\frac{\partial N_3}{\partial t} = \frac{1}{12}(-6t) = -\frac{1}{2}t \quad (5.244)$$

$$\frac{\partial N_4}{\partial t} = \frac{1}{12}(-6t) = -\frac{1}{2}t \quad (5.245)$$

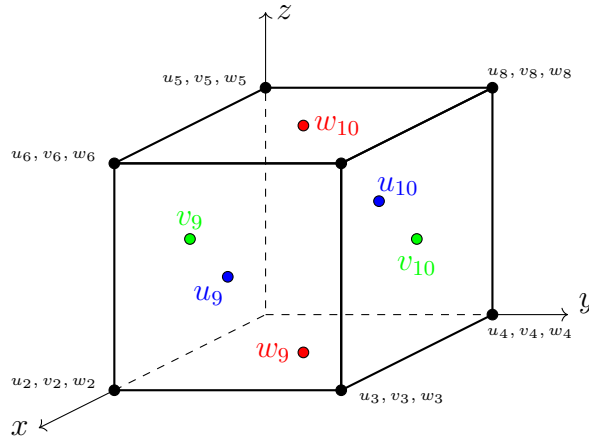
$$\frac{\partial N_5}{\partial t} = \frac{1}{12}(-6 + 12t) = \frac{1}{2}(-1 + 2t) \quad (5.246)$$

$$\frac{\partial N_6}{\partial t} = \frac{1}{12}(6 + 12t) = \frac{1}{2}(1 + 2t) \quad (5.247)$$

5.4.8 The 3D enriched $Q_1^+ \times P_0$ of Fortin

This element is mentioned on p249 of Cuvelier, Segal & van Steenhoven [298]: "The enriched trilinear velocity-constant pressure element is probably the simplest admissible 3D element." Fortin [401] designed a simple LBB-stable Q_1 element to which mid-face nodes are added, i.e. a 'bubble' Q_2 function is added on each face. However, only $\vec{v} \cdot \vec{n}$ is present on these mid-face nodes:

(tikz-q1pp0.tex)



Fortin states: "this element satisfies the B.B. condition and is probably the simplest 3-D element to do So. This unfortunately does not mean that it is more accurate (at least on regular meshes)." and "the element satisfies the B.B. condition. It can therefore be used in a non-regular mesh without fear. The number of degrees of freedom is approximately double with respect to the $Q_1 \times P_0$ element and this is reflected by an increased number of vortices and a reduction of their size. However, there seems to be a qualitative deficiency of these vortices since they do not easily assemble into complex flows. Only numerical experiments can give the final answer." This element is mentioned/used in [1080, 84, 1303].

Considering a single element, we have

- Q_1 : $2 \times 2 \times 2 \times 3 = 24$ velocity dofs
- Q_1^+ : $2 \times 2 \times 2 \times 3 + 6 = 30$ velocity dofs:

$$\vec{V}^T = \underbrace{(u_1, v_1, w_1, \dots, u_8, v_8, w_8)}_{Q_1 \text{ dofs}} \underbrace{(u_9, v_9, w_9, u_{10}, v_{10}, w_{10})}_{\text{bubble dofs}}$$

The big difference with all other elements so far is the fact that the dofs u_9, v_9, w_9 are not colocated (same for the other three). u_9 lives in the middle of the $r = -1$ face, v_9 lives in the middle of the $s = -1$ face and w_9 lives in the middle of the $t = -1$ face.

- Q_2 : $3 \times 3 \times 3 \times 3 = 81$ velocity dofs

Considering a 3D mesh composed of $nel = nelx \times nely \times nelz$ elements:

- Q_1 : the total number of Velocity dofs is $N_{femV} = (nelx + 1) \times (nely + 1) \times (nelz + 1) \times 3$
- Q_1^+ : the total number of nodes is

$$N_{femV} = (nelx+1) \times (nely+1) \times (nelz+1) \times 3 + (nelx+1) \times nely \times nelz + nelx \times (nely+1) \times nelz + nelx \times nely \times (nelz+1)$$

- Q_2 : the total number of Velocity dofs is $N_{femV} = (2nelx + 1) \times (2nely + 1) \times (2nelz + 1) \times 3$

When $nelx = nely = nelz = n \gg 1$ then the numbers above converge to $3n^3$, $6n^3$ and $24n^3$ respectively. This means that for large meshes the enriched Q_1 uses twice as many dofs as the standard Q_1 while the Q_2 element uses 8 times more.

x -component of velocity The polynomial representation of the velocity in the element is given by

$$u^h(r, s, t) = a + br + cs + dt + ers + frt + gst + hrst + kb_9(r, s, t) + lb_{10}(r, s, t)$$

where the two bubble functions are:

$$b_9^u(r, s, t) = \frac{1}{2}(1-r)(1-s^2)(1-t^2) \quad b_{10}^u(r, s, t) = \frac{1}{2}(1+r)(1-s^2)(1-t^2)$$

The coordinates of the u_9 dof is $(-1, 0, 0)$ and the coordinate of the u_{10} dof is $(1, 0, 0)$. We see that the bubble functions are 1 at their nodes and zero at all other nodes. We can actually use a different basis for $1, r, s, t, rs, rt, st, rst$ and we instead choose the standard Q_1 functions so that u^h becomes:

$$u^h(r, s, t) = aN_1 + bN_2 + cN_3 + dN_4 + eN_5 + fN_6 + gN_7 + hN_8 + kb_9(r, s, t) + lb_{10}(r, s, t)$$

We then must find the set of coefficients $\{a \dots l\}$ and we will do so by requiring that $u^h(r_i, s_i, t_i) = u_i$ for $i = 1, 10$.

The coordinates of all 10 nodes and the values of basis functions at these locations are:

| node # | r | s | t | N_1 | N_2 | N_3 | N_4 | N_5 | N_6 | N_7 | N_8 | b_9^u | b_{10}^u |
|--------|-----|-----|-----|-------|-------|-------|-------|-------|-------|-------|-------|---------|------------|
| 1 | -1 | -1 | -1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | +1 | -1 | -1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | +1 | +1 | -1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | -1 | +1 | -1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | -1 | -1 | +1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 6 | +1 | -1 | +1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 7 | +1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 8 | -1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 9 | -1 | 0 | 0 | 1/4 | 0 | 0 | 1/4 | 1/4 | 0 | 0 | 1/4 | 1 | 0 |
| 10 | +1 | 0 | 0 | 0 | 1/4 | 1/4 | 0 | 0 | 1/4 | 1/4 | 0 | 0 | 1 |

We then have the following ten equations:

$$\begin{aligned}
u_1 &= u^h(r_1, s_1, t_1) &= a \\
u_2 &= u^h(r_1, s_1, t_1) &= b \\
u_3 &= u^h(r_1, s_1, t_1) &= c \\
u_4 &= u^h(r_1, s_1, t_1) &= d \\
u_5 &= u^h(r_1, s_1, t_1) &= e \\
u_6 &= u^h(r_1, s_1, t_1) &= f \\
u_7 &= u^h(r_1, s_1, t_1) &= g \\
u_8 &= u^h(r_1, s_1, t_1) &= h \\
u_9 &= u^h(r_9, s_9, t_9) &= \frac{1}{4}(a + d + e + h) + k \\
u_{10} &= u^h(r_{10}, s_{10}, t_{10}) &= \frac{1}{4}(b + c + f + g) + l
\end{aligned}$$

or,

$$\begin{pmatrix}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\
1/4 & 0 & 0 & 1/4 & 1/4 & 0 & 0 & 1/4 & 1 & 0 \\
0 & 1/4 & 1/4 & 0 & 0 & 1/4 & 1/4 & 0 & 0 & 1
\end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \\ e \\ f \\ g \\ h \\ k \\ l \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \\ u_6 \\ u_7 \\ u_8 \\ u_9 \\ u_{10} \end{pmatrix}$$

This yields

$$\begin{aligned}
a &= u_1 \\
b &= u_2 \\
c &= u_3 \\
d &= u_4 \\
e &= u_5 \\
f &= u_6 \\
g &= u_7 \\
h &= u_8 \\
k &= u_9 - \frac{1}{4}(u_1 + u_4 + u_5 + u_8) \\
l &= u_{10} - \frac{1}{4}(u_2 + u_3 + u_6 + u_7)
\end{aligned}$$

and then

$$\begin{aligned}
u^h(r, s, t) &= aN_1 + bN_2 + cN_3 + dN_4 + eN_5 + fN_6 + gN_7 + hN_8 + kb_9^u(r, s, t) + lb_{10}^u(r, s, t) \\
&= u_1N_1 + u_2N_2 + u_3N_3 + u_4N_4 + u_5N_5 + u_6N_6 + u_7N_7 + u_8N_8 \\
&\quad + \left[u_9 - \frac{1}{4}(u_1 + u_4 + u_5 + u_8) \right] b_9^u(r, s, t) + \left[u_{10} - \frac{1}{4}(u_2 + u_3 + u_6 + u_7) \right] b_{10}^u(r, s, t) \\
&= \left(u_1 - \frac{1}{4}b_9 \right) N_1 + \left(u_2 - \frac{1}{4}b_{10} \right) N_2 + \left(u_3 - \frac{1}{4}b_{10} \right) N_3 + \left(u_4 - \frac{1}{4}b_9 \right) N_4 + \\
&\quad \left(u_5 - \frac{1}{4}b_9 \right) N_5 + \left(u_6 - \frac{1}{4}b_{10} \right) N_6 + \left(u_7 - \frac{1}{4}b_{10} \right) N_7 + \left(u_8 - \frac{1}{4}b_9 \right) N_8 + \\
&\quad b_9^u(r, s, t)u_9 + b_{10}^u(r, s, t)u_{10}
\end{aligned}$$

Finally, we can write the basis functions for the u field:

$$\begin{aligned}
N_1^u(r, s, t) &= N_1(r, s, t) - \frac{1}{4}b_9^u(r, s, t) \\
N_2^u(r, s, t) &= N_2(r, s, t) - \frac{1}{4}b_{10}^u(r, s, t) \\
N_3^u(r, s, t) &= N_3(r, s, t) - \frac{1}{4}b_{10}^u(r, s, t) \\
N_4^u(r, s, t) &= N_4(r, s, t) - \frac{1}{4}b_9^u(r, s, t) \\
N_5^u(r, s, t) &= N_5(r, s, t) - \frac{1}{4}b_9^u(r, s, t) \\
N_6^u(r, s, t) &= N_6(r, s, t) - \frac{1}{4}b_{10}^u(r, s, t) \\
N_7^u(r, s, t) &= N_7(r, s, t) - \frac{1}{4}b_{10}^u(r, s, t) \\
N_8^u(r, s, t) &= N_8(r, s, t) - \frac{1}{4}b_9^u(r, s, t) \\
N_9^u(r, s, t) &= b_9^u(r, s, t) \\
N_{10}^u(r, s, t) &= b_{10}^u(r, s, t)
\end{aligned}$$

And it is easy to verify that

$$\sum_{i=1}^{10} N_i^u(r, s, t) = 1 \quad \forall r, s, t$$

During the implementation phase we will need the derivatives of the basis functions, which are trivial for the standard Q_1 basis functions N_i . Remain then

$$\begin{aligned}
\partial_r b_9^u(r, s, t) &= \frac{\partial}{\partial r} \left(\frac{1}{2}(1-r)(1-s^2)(1-t^2) \right) = -\frac{1}{2}(1-s^2)(1-t^2) \\
\partial_s b_9^u(r, s, t) &= \frac{\partial}{\partial s} \left(\frac{1}{2}(1-r)(1-s^2)(1-t^2) \right) = -(1-r)s(1-t^2) \\
\partial_t b_9^u(r, s, t) &= \frac{\partial}{\partial t} \left(\frac{1}{2}(1-r)(1-s^2)(1-t^2) \right) = -(1-r)(1-s^2)t \\
\partial_r b_{10}^u(r, s, t) &= \frac{\partial}{\partial r} \left(\frac{1}{2}(1+r)(1-s^2)(1-t^2) \right) = \frac{1}{2}(1-s^2)(1-t^2) \\
\partial_s b_{10}^u(r, s, t) &= \frac{\partial}{\partial s} \left(\frac{1}{2}(1+r)(1-s^2)(1-t^2) \right) = -(1+r)s(1-t^2) \\
\partial_t b_{10}^u(r, s, t) &= \frac{\partial}{\partial t} \left(\frac{1}{2}(1+r)(1-s^2)(1-t^2) \right) = -(1+r)(1-s^2)t
\end{aligned}$$

y-component of velocity The bubbles are given by

$$b_9^v(r, s, t) = \frac{1}{2}(1-r^2)(1-s)(1-t^2) \quad b_{10}^v(r, s, t) = \frac{1}{2}(1-r^2)(1+s)(1-t^2)$$

The coordinates of all 10 nodes and the values of basis functions at these locations are:

| node # | r | s | t | N_1 | N_2 | N_3 | N_4 | N_5 | N_6 | N_7 | N_8 | b_9^v | b_{10}^v |
|--------|-----|-----|-----|-------|-------|-------|-------|-------|-------|-------|-------|---------|------------|
| 1 | -1 | -1 | -1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | +1 | -1 | -1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | +1 | +1 | -1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | -1 | +1 | -1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | -1 | -1 | +1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 6 | +1 | -1 | +1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 7 | +1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 8 | -1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 9 | 0 | -1 | 0 | 1/4 | 1/4 | 0 | 0 | 1/4 | 1/4 | 0 | 0 | 1 | 0 |
| 10 | 0 | +1 | 0 | 0 | 0 | 1/4 | 1/4 | 0 | 0 | 1/4 | 1/4 | 0 | 1 |

Then

$$\begin{aligned}
N_1^v(r, s, t) &= N_1(r, s, t) - \frac{1}{4}b_9^v(r, s, t) \\
N_2^v(r, s, t) &= N_2(r, s, t) - \frac{1}{4}b_9^v(r, s, t) \\
N_3^v(r, s, t) &= N_3(r, s, t) - \frac{1}{4}b_{10}^v(r, s, t) \\
N_4^v(r, s, t) &= N_4(r, s, t) - \frac{1}{4}b_{10}^v(r, s, t) \\
N_5^v(r, s, t) &= N_5(r, s, t) - \frac{1}{4}b_9^v(r, s, t) \\
N_6^v(r, s, t) &= N_6(r, s, t) - \frac{1}{4}b_9^v(r, s, t) \\
N_7^v(r, s, t) &= N_7(r, s, t) - \frac{1}{4}b_{10}^v(r, s, t) \\
N_8^v(r, s, t) &= N_8(r, s, t) - \frac{1}{4}b_{10}^v(r, s, t) \\
N_9^v(r, s, t) &= b_9^v(r, s, t) \\
N_{10}^v(r, s, t) &= b_{10}^v(r, s, t)
\end{aligned} \tag{5.248}$$

z-component of velocity The bubbles are given by

$$b_9^w(r, s, t) = \frac{1}{2}(1 - r^2)(1 - s^2)(1 - t) \quad b_{10}^w(r, s, t) = \frac{1}{2}(1 - r^2)(1 - s^2)(1 + t)$$

The coordinates of all 10 nodes and the values of basis functions at these locations are:

| node # | r | s | t | N_1 | N_2 | N_3 | N_4 | N_5 | N_6 | N_7 | N_8 | b_9^w | b_{10}^w |
|--------|-----|-----|-----|-------|-------|-------|-------|-------|-------|-------|-------|---------|------------|
| 1 | -1 | -1 | -1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | +1 | -1 | -1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | +1 | +1 | -1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | -1 | +1 | -1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | -1 | -1 | +1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 6 | +1 | -1 | +1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 7 | +1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 8 | -1 | +1 | +1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 9 | 0 | 0 | -1 | 1/4 | 1/4 | 1/4 | 1/4 | 0 | 0 | 0 | 0 | 1 | 0 |
| 10 | 0 | 0 | +1 | 0 | 0 | 0 | 0 | 1/4 | 1/4 | 1/4 | 1/4 | 0 | 1 |

$$\begin{aligned}
N_1^w(r, s, t) &= N_1(r, s, t) - \frac{1}{4}b_9^w(r, s, t) \\
N_2^w(r, s, t) &= N_2(r, s, t) - \frac{1}{4}b_9^w(r, s, t) \\
N_3^w(r, s, t) &= N_3(r, s, t) - \frac{1}{4}b_9^w(r, s, t) \\
N_4^w(r, s, t) &= N_4(r, s, t) - \frac{1}{4}b_9^w(r, s, t) \\
N_5^w(r, s, t) &= N_5(r, s, t) - \frac{1}{4}b_{10}^w(r, s, t) \\
N_6^w(r, s, t) &= N_6(r, s, t) - \frac{1}{4}b_{10}^w(r, s, t) \\
N_7^w(r, s, t) &= N_7(r, s, t) - \frac{1}{4}b_{10}^w(r, s, t) \\
N_8^w(r, s, t) &= N_8(r, s, t) - \frac{1}{4}b_{10}^w(r, s, t) \\
N_9^w(r, s, t) &= b_9^w(r, s, t) \\
N_{10}^w(r, s, t) &= b_{10}^w(r, s, t)
\end{aligned}$$

A word about the B matrix We have

$$u^h(r, s, t) = \sum_{i=1}^{10} N_i^u(r, s, t)u_i \quad (5.249)$$

$$v^h(r, s, t) = \sum_{i=1}^{10} N_i^v(r, s, t)v_i \quad (5.250)$$

$$w^h(r, s, t) = \sum_{i=1}^{10} N_i^w(r, s, t)w_i \quad (5.251)$$

Normally we do not make a distinction between the basis functions associated to u, v, w but because of the bubbles on the faces we now have to.

We have previously established that the strain rate vector $\vec{\varepsilon}$ is:

$$\vec{\varepsilon} = \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix} = \begin{pmatrix} \sum_i \frac{\partial N_i^u}{\partial x} u_i \\ \sum_i \frac{\partial N_i^v}{\partial y} v_i \\ \sum_i \frac{\partial N_i^w}{\partial z} w_i \\ \sum_i \left(\frac{\partial N_i^u}{\partial y} u_i + \frac{\partial N_i^v}{\partial x} v_i \right) \\ \sum_i \left(\frac{\partial N_i^u}{\partial z} u_i + \frac{\partial N_i^w}{\partial x} w_i \right) \\ \sum_i \left(\frac{\partial N_i^v}{\partial z} v_i + \frac{\partial N_i^w}{\partial y} w_i \right) \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial N_1^u}{\partial x} & 0 & 0 & \dots & \frac{\partial N_{10}^u}{\partial x} & 0 & 0 \\ 0 & \frac{\partial N_1^v}{\partial y} & 0 & \dots & 0 & \frac{\partial N_{10}^v}{\partial y} & 0 \\ 0 & 0 & \frac{\partial N_1^w}{\partial z} & \dots & 0 & 0 & \frac{\partial N_{10}^w}{\partial z} \\ \frac{\partial N_1^u}{\partial y} & \frac{\partial N_1^v}{\partial x} & 0 & \dots & \frac{\partial N_{10}^u}{\partial x} & \frac{\partial N_{10}^v}{\partial x} & 0 \\ \frac{\partial N_1^u}{\partial z} & 0 & \frac{\partial N_1^w}{\partial x} & \dots & \frac{\partial N_{10}^u}{\partial z} & 0 & \frac{\partial N_{10}^w}{\partial x} \\ 0 & \frac{\partial N_1^v}{\partial z} & \frac{\partial N_1^w}{\partial y} & \dots & 0 & \frac{\partial N_{10}^v}{\partial z} & \frac{\partial N_{10}^w}{\partial y} \end{pmatrix}}_{\mathbf{B}} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ w_1 \\ u_2 \\ v_2 \\ w_2 \\ u_3 \\ v_3 \\ \dots \\ u_{10} \\ v_{10} \\ w_{10} \end{pmatrix}}_{\vec{V}}$$

5.4.9 The $Q_1^{++} \times Q_1$ of Karabelas et al (2020)

q1q13D_2bubbles.tex

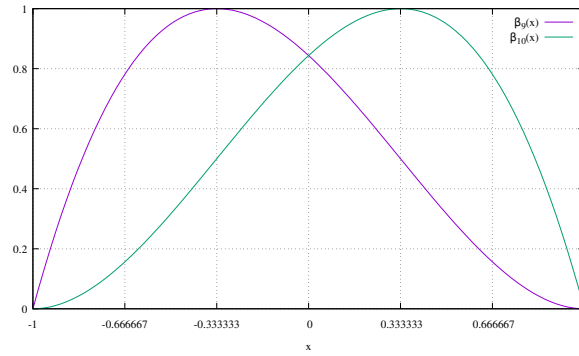
This element is implemented in [STONE](#) 82. The two bubble functions are given in Karabelas *et al.* (2020) [670]:

$$\begin{aligned} b_9(r, s, t) &= \left(\frac{27}{32}\right)^3 (1 - r^2)(1 - s^2)(1 - t^2) \cdot (1 - r)(1 - s)(1 - t) = \beta_9(r) \cdot \beta_9(s) \cdot \beta_9(t) \\ b_{10}(r, s, t) &= \left(\frac{27}{32}\right)^3 (1 - r^2)(1 - s^2)(1 - t^2) \cdot (1 + r)(1 + s)(1 + t) = \beta_{10}(r) \cdot \beta_{10}(s) \cdot \beta_{10}(t) \end{aligned}$$

where I have chosen nodes 1 ($\vec{r}_1 = (-1, -1, -1)$) and 7 ($\vec{r}_7 = (+1, +1, +1)$) as diagonally opposed nodes (a requirement from the paper), and with

$$\beta_9(x) = \frac{27}{32}(1 - x^2)(1 - x) \quad \beta_{10}(x) = \frac{27}{32}(1 - x^2)(1 + x)$$

I have added the $(27/32)^3$ coefficients so that these functions are exactly 1 at their corresponding nodes. The term $(1 - r^2)(1 - s^2)(1 - t^2)$ makes sure that the two bubbles are conforming and exactly zero on the 6 faces of the element. In what follows $\tilde{\mathcal{N}}_{1,8}$ are the standard Q_1 basis functions.



Representation of bubbles $\beta_9(x)$ and $\beta_{10}(x)$

Remark. *Bubble function 9 is not zero at node 10 and vice versa!*

The authors state: "This also allows for a straightforward inclusion in combination with existing finite element codes since all required implementations are purely on the element level". This is especially true if static condensation is used (the authors explain static condensation for the bubbles in the appendix of the paper).

The ten nodes are the standard 8 corners of the Q_1 element as well as $\vec{r}_9 = (-1/3, -1/3, -1/3)$ for b_9 and $\vec{r}_{10} = (1/3, 1/3, 1/3)$ for b_{10} . We have the following approximation of function f inside the element:

$$f^h(r, s, t) = \sum_{i=1}^8 a_i \tilde{\mathcal{N}}_i(r, s, t) + a_9 b_9(r, s, t) + a_{10} b_{10}(r, s, t)$$

We notice that bubble functions are exactly zero at the corners of the reference element and we can

compute the values of the ten polynomials $(\tilde{N}_{1-8}(r, s, t), b_9(r, s, t), b_{10}(r, s, t))$ at the ten nodes:

| | \tilde{N}_1 | \tilde{N}_2 | \tilde{N}_3 | \tilde{N}_4 | \tilde{N}_5 | \tilde{N}_6 | \tilde{N}_7 | \tilde{N}_8 | b_9 | b_{10} |
|-------------------------------------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|-------|----------|
| $\vec{r}_1 = (-1, -1, -1)$ | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $\vec{r}_2 = (+1, -1, -1)$ | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $\vec{r}_3 = (+1, +1, -1)$ | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $\vec{r}_4 = (-1, +1, -1)$ | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| $\vec{r}_5 = (-1, -1, +1)$ | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| $\vec{r}_6 = (+1, -1, +1)$ | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| $\vec{r}_7 = (+1, +1, +1)$ | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| $\vec{r}_8 = (-1, +1, +1)$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| $\vec{r}_9 = (-\frac{1}{3}, -\frac{1}{3}, -\frac{1}{3})$ | 8/27 | 4/27 | 2/27 | 4/27 | 4/27 | 2/27 | 1/27 | 2/27 | 1 | 1/8 |
| $\vec{r}_{10} = (+\frac{1}{3}, +\frac{1}{3}, +\frac{1}{3})$ | 1/27 | 2/17 | 4/27 | 2/27 | 2/27 | 4/27 | 8/27 | 4/27 | 1/8 | 1 |

We then require that the polynomial representation of f^h of f inside the element is such that $f^h(\vec{r}_i) = f_i$, i.e.:

$$\begin{aligned}
f_1 = f^h(r_1, s_1, t_1) &= a_1 \\
f_2 = f^h(r_2, s_1, t_1) &= a_2 \\
f_3 = f^h(r_3, s_1, t_1) &= a_3 \\
f_4 = f^h(r_4, s_1, t_1) &= a_4 \\
f_5 = f^h(r_5, s_1, t_1) &= a_5 \\
f_6 = f^h(r_6, s_1, t_1) &= a_6 \\
f_7 = f^h(r_7, s_1, t_1) &= a_7 \\
f_8 = f^h(r_8, s_1, t_1) &= a_8 \\
f_9 = f^h(r_9, s_9, t_9) &= \frac{1}{27}(8a_1 + 4a_2 + 2a_3 + 4a_4 + 4a_5 + 2a_6 + a_7 + 2a_8) + a_9 + a_{10}/8 \\
f_{10} = f^h(r_{10}, s_{10}, t_{10}) &= \frac{1}{27}(a_1 + 2a_2 + 4a_3 + 2a_4 + 2a_5 + 4a_6 + 8a_7 + 4a_8) + a_9/8 + a_{10}
\end{aligned}$$

or,

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 8/27 & 4/27 & 2/27 & 4/27 & 4/27 & 2/27 & 1/27 & 2/27 & 1 & 1/8 \\ 1/27 & 2/17 & 4/27 & 2/27 & 2/27 & 4/27 & 8/27 & 4/27 & 1/8 & 1 \end{pmatrix} \cdot \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ a_{10} \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \\ f_{10} \end{pmatrix}$$

which yields $a_i = f_i$ for $i = 1, \dots, 8$ and

$$\begin{aligned}
a_9 + a_{10}/8 &= \underbrace{f_9 - \frac{1}{27}(8f_1 + 4f_2 + 2f_3 + 4f_4 + 4f_5 + 2f_6 + f_7 + 2f_8)}_{\tilde{f}_9} \\
a_9/8 + a_{10} &= \underbrace{f_{10} - \frac{1}{27}(f_1 + 2f_2 + 4f_3 + 2f_4 + 2f_5 + 4f_6 + 8f_7 + 4f_8)}_{\tilde{f}_{10}}
\end{aligned}$$

$$\begin{aligned}
8a_9 + a_{10} &= 8f_9 + 8\tilde{f}_9 \\
a_9 + 8a_{10} &= 8f_{10} + 8\tilde{f}_{10}
\end{aligned}$$

and then

$$\begin{aligned}
a_9 &= \frac{1}{63}(64f_9 - 8f_{10}) + \frac{1}{63}(64\tilde{f}_9 - 8\tilde{f}_{10}) \\
&= \frac{8}{63}(8f_9 - f_{10}) + \frac{8}{63}(8\tilde{f}_9 - \tilde{f}_{10}) \\
&= \frac{8}{63}(8f_9 - f_{10}) - \frac{1}{27} \frac{8}{63} [8(8f_1 + 4f_2 + 2f_3 + 4f_4 + 4f_5 + 2f_6 + f_7 + 2f_8) \\
&\quad -(f_1 + 2f_2 + 4f_3 + 2f_4 + 2f_5 + 4f_6 + 8f_7 + 4f_8)] \\
&= \frac{8}{63}(8f_9 - f_{10}) - \frac{1}{27} \frac{8}{63} (63f_1 + 30f_2 + 12f_3 + 30f_4 + 30f_5 + 12f_6 + 12f_8) \\
a_{10} &= \frac{1}{63}(64f_{10} - 8f_9) + \frac{1}{63}(64\tilde{f}_{10} - 8\tilde{f}_9) \\
&= \frac{8}{63}(8f_{10} - f_9) + \frac{8}{63}(8\tilde{f}_{10} - \tilde{f}_9) \\
&= \frac{8}{63}(8f_{10} - f_9) - \frac{1}{27} \frac{8}{63} [8(f_1 + 2f_2 + 4f_3 + 2f_4 + 2f_5 + 4f_6 + 8f_7 + 4f_8) \\
&\quad -(8f_1 + 4f_2 + 2f_3 + 4f_4 + 4f_5 + 2f_6 + f_7 + 2f_8)] \\
&= \frac{8}{63}(8f_{10} - f_9) - \frac{1}{27} \frac{8}{63} (12f_2 + 30f_3 + 12f_4 + 12f_5 + 30f_6 + 63f_7 + 30f_8)
\end{aligned}$$

We can then write

$$\begin{aligned}
f^h(r, s, t) &= f_1 \mathcal{N}_1(r, s, t) + f_2 \mathcal{N}_2(r, s, t) + f_3 \mathcal{N}_3(r, s, t) + f_4 \mathcal{N}_4(r, s, t) \\
&+ f_5 \mathcal{N}_5(r, s, t) + f_6 \mathcal{N}_6(r, s, t) + f_7 \mathcal{N}_7(r, s, t) + f_8 \mathcal{N}_8(r, s, t) \\
&+ \left[\frac{8}{63}(8f_9 - f_{10}) - \frac{1}{27} \frac{8}{63}(63f_1 + 30f_2 + 12f_3 + 30f_4 + 30f_5 + 12f_6 + 12f_8) \right] b_9(r, s, t) \\
&\quad \left[\frac{8}{63}(8f_{10} - f_9) - \frac{1}{27} \frac{8}{63}(12f_2 + 30f_3 + 12f_4 + 12f_5 + 30f_6 + 63f_7 + 30f_8) \right] b_{10}(r, s, t) \\
&= \left(\mathcal{N}_1(r, s, t) - \frac{2^3}{3^3} b_9(r, s, t) \right) f_1 \\
&+ \left(\mathcal{N}_2(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \right) f_2 \\
&+ \left(\mathcal{N}_3(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \right) f_3 \\
&+ \left(\mathcal{N}_4(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \right) f_4 \\
&+ \left(\mathcal{N}_5(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \right) f_5 \\
&+ \left(\mathcal{N}_6(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \right) f_6 \\
&+ \left(\mathcal{N}_7(r, s, t) - \frac{2^3}{3^3} b_{10}(r, s, t) \right) f_7 \\
&+ \left(\mathcal{N}_8(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \right) f_8 \\
&+ \left(\frac{64}{63} b_9(r, s, t) - \frac{8}{63} b_{10}(r, s, t) \right) f_9 + \left(-\frac{8}{63} b_9(r, s, t) + \frac{64}{63} b_{10}(r, s, t) \right) f_{10} \quad (5.252)
\end{aligned}$$

and finally arrive at the basis functions:

$$\begin{aligned}
\mathcal{N}_1(r, s, t) &= \tilde{\mathcal{N}}_1(r, s, t) - \frac{2^3}{3^3} b_9(r, s, t) \\
\mathcal{N}_2(r, s, t) &= \tilde{\mathcal{N}}_2(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \\
\mathcal{N}_3(r, s, t) &= \tilde{\mathcal{N}}_3(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \\
\mathcal{N}_4(r, s, t) &= \tilde{\mathcal{N}}_4(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \\
\mathcal{N}_5(r, s, t) &= \tilde{\mathcal{N}}_5(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t) \\
\mathcal{N}_6(r, s, t) &= \tilde{\mathcal{N}}_6(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \\
\mathcal{N}_7(r, s, t) &= \tilde{\mathcal{N}}_7(r, s, t) - \frac{2^3}{3^3} b_{10}(r, s, t) \\
\mathcal{N}_8(r, s, t) &= \tilde{\mathcal{N}}_8(r, s, t) - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t) \\
\mathcal{N}_9(r, s, t) &= \frac{64}{63} b_9(r, s, t) - \frac{8}{63} b_{10}(r, s, t) \\
\mathcal{N}_{10}(r, s, t) &= -\frac{8}{63} b_9(r, s, t) + \frac{64}{63} b_{10}(r, s, t)
\end{aligned}$$

These are somewhat complex forms for the basis functions so we wish to verify the simple property $\sum \mathcal{N}_i(r, s, t) = 1$ for all (r, s, t) inside the element:

$$\begin{aligned}
\sum_{i=1}^{10} \mathcal{N}_i(r, s, t) &= \sum_{i=1}^8 \tilde{\mathcal{N}}_i(r, s, t) \\
&+ \left[\frac{2^3}{3^3} \left(-1 - \frac{10}{21} - \frac{4}{21} - \frac{10}{21} - \frac{10}{21} - \frac{4}{21} - \frac{4}{21} \right) + \frac{64}{63} - \frac{8}{63} \right] b_9(r, s, t) \\
&+ \left[\frac{2^3}{3^3} \left(-\frac{4}{21} - \frac{10}{21} - \frac{4}{21} - \frac{4}{21} - \frac{10}{21} - 1 - \frac{10}{21} \right) - \frac{8}{63} + \frac{64}{63} \right] b_{10}(r, s, t) \\
&= 1 + \left[\frac{2^3}{3^3} (-1 - 42/21) + \frac{56}{63} \right] b_9(r, s, t) + \left[\frac{2^3}{3^3} (-42/21 - 1) + \frac{56}{63} \right] b_{10}(r, s, t) \\
&= 1 + \left[\frac{2^3}{3^3} (-3) + \frac{8}{9} \right] b_9(r, s, t) + \left[\frac{2^3}{3^3} (-3) + \frac{8}{9} \right] b_{10}(r, s, t) \\
&= 1
\end{aligned} \tag{5.253}$$

Let us move to first order consistency with $f(r) = r$:

$$f^h(r, s, t) = \sum_{i=1}^{10} \mathcal{N}_i(r, s, t) f_i = \sum_{i=1}^{10} \mathcal{N}_i(r, s, t) r_i \tag{5.254}$$

It has been established for the $\tilde{\mathcal{N}}_i$ functions so we are left with


$$\begin{aligned}
f^h(r, s, t) &= \underbrace{\sum_{i=1}^8 \tilde{\mathcal{N}}_i r_i}_{=r} \\
&\quad - \frac{2^3}{3^3} b_9(r, s, t)(-1) \\
&\quad - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t)(+1) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t)(+1) \\
&\quad - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t)(+1) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t)(+1) \\
&\quad - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t)(-1) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t)(-1) \\
&\quad - \frac{2^3}{3^3} \frac{10}{21} b_9(r, s, t)(-1) - \frac{2^3}{3^3} \frac{4}{21} b_{10}(r, s, t)(-1) \\
&\quad - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t)(+1) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t)(+1) \\
&\quad - \frac{2^3}{3^3} b_{10}(r, s, t)(+1) \\
&\quad - \frac{2^3}{3^3} \frac{4}{21} b_9(r, s, t)(-1) - \frac{2^3}{3^3} \frac{10}{21} b_{10}(r, s, t)(-1) \\
&\quad + \frac{64}{63} b_9(r, s, t)(-1/3) - \frac{8}{63} b_{10}(r, s, t)(-1/3) \\
&\quad - \frac{8}{63} b_9(r, s, t)(+1/3) + \frac{64}{63} b_{10}(r, s, t)(+1/3) \\
&= r + b_9(r, s, t) \left(\frac{8}{27} - \frac{8}{27} \frac{10}{21} - \frac{8}{27} \frac{4}{21} + \frac{8}{27} \frac{10}{21} + \frac{8}{27} \frac{10}{21} - \frac{8}{27} \frac{4}{21} + \frac{8}{27} \frac{4}{21} - \frac{64}{189} - \frac{8}{189} \right) \\
&\quad + b_{10}(r, s, t) \left(-\frac{8}{27} \frac{4}{21} - \frac{8}{27} \frac{10}{21} + \frac{8}{27} \frac{4}{21} + \frac{8}{27} \frac{4}{21} - \frac{8}{27} \frac{10}{21} - \frac{8}{27} + \frac{8}{27} \frac{10}{21} + \frac{8}{189} + \frac{64}{189} \right) \\
&= r + b_9(r, s, t) \frac{8}{27} \underbrace{\left(1 - \frac{10}{21} - \frac{4}{21} + \frac{10}{21} + \frac{10}{21} - \frac{9}{7} \right)}_{=0} + b_{10}(r, s, t) \frac{8}{27} \underbrace{\left(-\frac{4}{21} - \frac{10}{21} + \frac{4}{21} + \frac{4}{21} - \frac{10}{21} - \right)}_{=0} \\
&= r
\end{aligned}$$

which proves first-order consistency.

The derivatives of the $\tilde{\mathcal{N}}_i$ basis functions are already established so we only focus on the spatial

derivatives of the bubble functions:

$$\begin{aligned}
\frac{\partial b_9}{\partial r} &= \left(\frac{27}{32}\right)^3 (1-s^2)(1-t^2)(1-s)(1-t)(-1-2r+3r^2) \\
\frac{\partial b_9}{\partial s} &= \left(\frac{27}{32}\right)^3 (1-r^2)(1-t^2)(1-r)(1-t)(-1-2s+3s^2) \\
\frac{\partial b_9}{\partial t} &= \left(\frac{27}{32}\right)^3 (1-r^2)(1-s^2)(1-r)(1-s)(-1-2t+3t^2) \\
\frac{\partial b_{10}}{\partial r} &= \left(\frac{27}{32}\right)^3 (1-s^2)(1-t^2)(1+s)(1+t)(1-2r-3r^2) \\
\frac{\partial b_{10}}{\partial s} &= \left(\frac{27}{32}\right)^3 (1-r^2)(1-t^2)(1+r)(1+t)(1-2s-3s^2) \\
\frac{\partial b_{10}}{\partial t} &= \left(\frac{27}{32}\right)^3 (1-r^2)(1-s^2)(1+r)(1+s)(1-2t-3t^2)
\end{aligned}$$

 **Relevant Literature:** Fortin and Fortin [403] (1985), Soulaïmani, Fortin, Ouellet, Dhatt, and Bertrand [1181] (1987)

Bishnu talksabout Nitsche bc ? press error near boundary in fof085

5.4.10 The DSSY element

What follows is mostly from From Jang *et al.* (2005) [633].

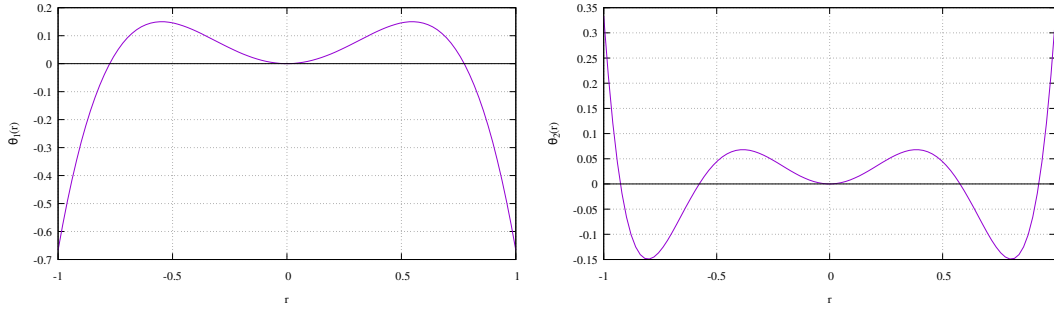
The non-conforming finite element space is defined based on the reference cubic element Q_l on $[-1, 1]^3$:

$$Q_l = \text{Span} \{1, r, s, t, \theta_l(r) - \theta_l(s), \theta_l(r) - \theta_l(t)\} \quad l = 1, 2$$

with¹¹

$$\begin{aligned} \theta_1(r) &= r^2 - \frac{5}{3}r^4 \\ \theta_2(r) &= r^2 - \frac{25}{6}r^4 + \frac{7}{2}r^6 \end{aligned} \quad (5.256)$$

The dimension of Q_l is six and the θ_l functions are as follows:



Representation of functions θ_1 (left) and θ_2 (right).

We have:

- $\theta_1(r = -1) = \theta_1(r = +1) = -\frac{2}{3}$, $\theta_1(r = 0) = 0$
- $\theta_2(r = -1) = \theta_2(r = +1) = \frac{1}{3}$, $\theta_2(r = 0) = 0$

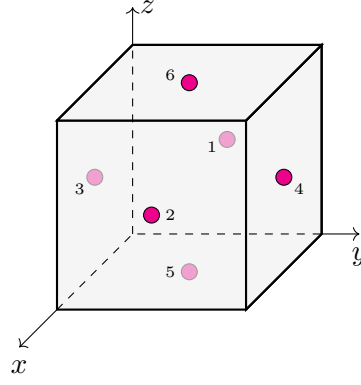
These functions have the property that their average on an edge is zero:

$$\begin{aligned} \frac{1}{1 - (-1)} \int_{-1}^{+1} \theta_1(r) dr &= \frac{1}{2} \int_{-1}^{+1} \left(r^2 - \frac{5}{3}r^4 \right) dr \\ &= \frac{1}{2} \left(\frac{2}{3} - \frac{5}{3} \frac{2}{5} \right) \\ &= 0 \\ \frac{1}{1 - (-1)} \int_{-1}^{+1} \theta_2(r) dr &= \frac{1}{2} \int_{-1}^{+1} \left(r^2 - \frac{25}{6}r^4 + \frac{7}{2}r^6 \right) dr \\ &= \frac{1}{2} \left(\frac{2}{3} - \frac{25}{6} \frac{2}{5} + \frac{7}{2} \frac{2}{7} \right) dr \\ &= 0 \end{aligned}$$

The element has 6 nodes that are located at the face centers of a cube or a brick. The pointwise continuity between interfacing elements is guaranteed only at the face centers, so the field quantities are not conforming along the interface [633].

(tikz_dssy3D.tex)

¹¹Douglas *et al.* [345], Eq. 2.20



(axis should actually be in the middle of the cube!)

The six nodes of the reference element are

| | | | |
|-------------|----|----|----|
| \vec{r}_1 | -1 | 0 | 0 |
| \vec{r}_2 | +1 | 0 | 0 |
| \vec{r}_3 | 0 | -1 | 0 |
| \vec{r}_4 | 0 | + | 0 |
| \vec{r}_5 | 0 | 0 | -1 |
| \vec{r}_6 | 0 | 0 | +1 |

and their corresponding basis functions:

$$\begin{aligned}
N_1^{(l)}(r, s, t) &= \frac{1}{6} - \frac{1}{2}r + \frac{1}{6\theta_l(1)}(2\theta_l(r) - \theta_l(s) - \theta_l(t)) \\
N_2^{(l)}(r, s, t) &= \frac{1}{6} + \frac{1}{2}r + \frac{1}{6\theta_l(1)}(2\theta_l(r) - \theta_l(s) - \theta_l(t)) \\
N_3^{(l)}(r, s, t) &= \frac{1}{6} - \frac{1}{2}s + \frac{1}{6\theta_l(1)}(2\theta_l(s) - \theta_l(t) - \theta_l(r)) \\
N_4^{(l)}(r, s, t) &= \frac{1}{6} + \frac{1}{2}s + \frac{1}{6\theta_l(1)}(2\theta_l(s) - \theta_l(t) - \theta_l(r)) \\
N_5^{(l)}(r, s, t) &= \frac{1}{6} - \frac{1}{2}t + \frac{1}{6\theta_l(1)}(2\theta_l(t) - \theta_l(r) - \theta_l(s)) \\
N_6^{(l)}(r, s, t) &= \frac{1}{6} + \frac{1}{2}t + \frac{1}{6\theta_l(1)}(2\theta_l(t) - \theta_l(r) - \theta_l(s))
\end{aligned} \tag{5.257}$$

These basis functions are also in [345]. We can easily verify that $\sum_i N_i(r, s, t) = 1$ and that $N_i(\vec{r}_j) = \delta_{ij}$:

$$\begin{aligned}
N_1^{(l)}(r_1, s_1, t_1) &= \frac{1}{6} - \frac{1}{2}(-1) + \frac{1}{6\theta_l(1)}(2\theta_l(-1) - \theta_l(0) - \theta_l(0)) = 1 \\
N_1^{(l)}(r_2, s_2, t_2) &= \frac{1}{6} - \frac{1}{2}(+1) + \frac{1}{6\theta_l(1)}(2\theta_l(+1) - \theta_l(0) - \theta_l(0)) = 0 \\
N_1^{(l)}(r_3, s_3, t_3) &= \frac{1}{6} - \frac{1}{2}(0) + \frac{1}{6\theta_l(1)}(2\theta_l(0) - \theta_l(-1) - \theta_l(0)) = 0 \\
N_1^{(l)}(r_4, s_4, t_4) &= \frac{1}{6} - \frac{1}{2}(0) + \frac{1}{6\theta_l(1)}(2\theta_l(0) - \theta_l(+1) - \theta_l(0)) = 0 \\
N_1^{(l)}(r_5, s_5, t_5) &= \frac{1}{6} - \frac{1}{2}(0) + \frac{1}{6\theta_l(1)}(2\theta_l(0) - \theta_l(0) - \theta_l(-1)) = 0 \\
N_1^{(l)}(r_6, s_6, t_6) &= \frac{1}{6} - \frac{1}{2}(0) + \frac{1}{6\theta_l(1)}(2\theta_l(0) - \theta_l(0) - \theta_l(+1)) = 0
\end{aligned}$$

etc ...

$$\begin{aligned}
\partial_r N_1^{(l)}(r, s, t) &= -\frac{1}{2} + \frac{1}{3\theta_l(1)}\theta'_l(r) \\
\partial_r N_2^{(l)}(r, s, t) &= +\frac{1}{2} + \frac{1}{3\theta_l(1)}\theta'_l(r) \\
\partial_r N_3^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(r) \\
\partial_r N_4^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(r) \\
\partial_r N_5^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(r) \\
\partial_r N_6^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(r)
\end{aligned} \tag{5.258}$$

$$\begin{aligned}
\partial_s N_1^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(s) \\
\partial_s N_2^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta'_l(s) \\
\partial_s N_3^{(l)}(r, s, t) &= -\frac{1}{2} + \frac{1}{3\theta_l(1)}\theta'_l(s) \\
\partial_s N_4^{(l)}(r, s, t) &= +\frac{1}{2} + \frac{1}{3\theta_l(1)}\theta'_l(s) \\
\partial_s N_5^{(l)}(r, s, t) &= \frac{1}{6\theta_l(1)}\theta'_l(s) \\
\partial_s N_6^{(l)}(r, s, t) &= \frac{1}{6\theta_l(1)}\theta'_l(s)
\end{aligned} \tag{5.259}$$

$$\begin{aligned}
\partial_t N_1^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta_l(t) \\
\partial_t N_2^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta_l(t) \\
\partial_t N_3^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta_l(t) \\
\partial_t N_4^{(l)}(r, s, t) &= -\frac{1}{6\theta_l(1)}\theta_l(t) \\
\partial_t N_5^{(l)}(r, s, t) &= -\frac{1}{2}t + \frac{1}{3\theta_l(1)}\theta_l(t) \\
\partial_t N_6^{(l)}(r, s, t) &= +\frac{1}{2}t + \frac{1}{3\theta_l(1)}\theta_l(t)
\end{aligned} \tag{5.260}$$

 Relevant Literature: Douglas *et al.* (1999) [345],

5.5 Low order elements recap

Let us assume a Cartesian domain discretised in $nel_x \times nel_y$ elements in 2D and $nel_x \times nel_y \times nel_z$ elements in 3D. Focusing only on the total number of velocity dofs (the values indicated after the arrows are the limits when $nel_x = nel_y = nel_z \gg 1$):

- $Q_1 \times P_0, Q_1 \times Q_1$

$$ndof_{2D} = 2(nel_x + 1) \cdot (nel_y + 1) \rightarrow 2 \cdot nel_x^2 \quad (5.261)$$

$$ndof_{3D} = 3(nel_x + 1) \cdot (nel_y + 1) \cdot (nel_z + 1) \rightarrow 3 \cdot nel_x^3 \quad (5.262)$$

- $Q_2 \times P_{-1}, Q_2 \times Q_1$

$$ndof_{2D} = 2(2nel_x + 1) \cdot (2nel_y + 1) \rightarrow 8 \cdot nel_x^2$$

$$ndof_{3D} = 3(2nel_x + 1) \cdot (2nel_y + 1) \cdot (2nel_z + 1) \rightarrow 24 \cdot nel_x^3$$

- $Q_1^+ \times P_0$

$$ndof_{2D} = 2(nel_x + 1) \cdot (nel_y + 1) + (nel_x + 1) \cdot nel_y + nel_x \cdot (nel_y + 1) \rightarrow 4 \cdot nel_x^2$$

$$ndof_{3D} = 3(nel_x + 1) \cdot (nel_y + 1) \cdot (nel_z + 1) + (nel_x + 1) \cdot nel_y \cdot nel_z + nel_x \cdot (nel_y + 1) \cdot nel_z + nel_x \cdot nel_y \cdot (nel_z + 1) \rightarrow 6 \cdot nel_x^3$$

- $Q_1 \times Q_1 + 1$ bubble

$$ndof_{2D} = 2[(nel_x + 1) \cdot (nel_y + 1) + nel_x \cdot nel_y] \rightarrow 4 \cdot nel_x^2$$

- $Q_1 \times Q_1 + 2$ bubbles

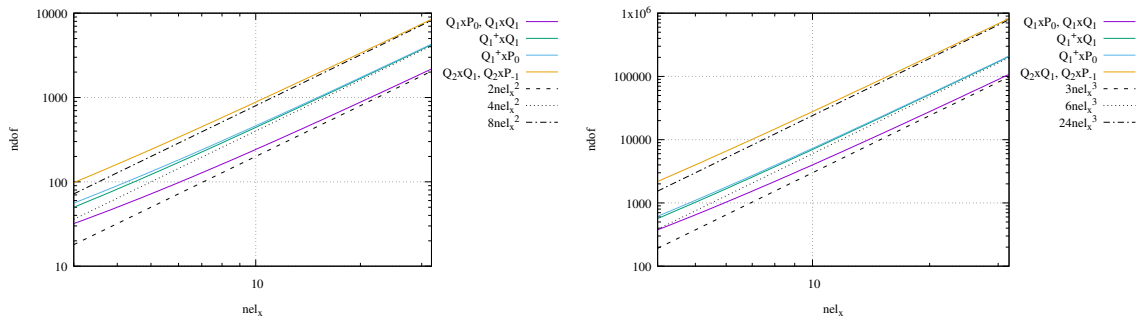
$$ndof_{3D} = 3[(nel_x + 1) \cdot (nel_y + 1) \cdot (nel_z + 1) + 2 \cdot nel_x \cdot nel_y \cdot nel_z] \rightarrow 9 \cdot nel_x^3$$

- Rannacher-turek or DSSY:

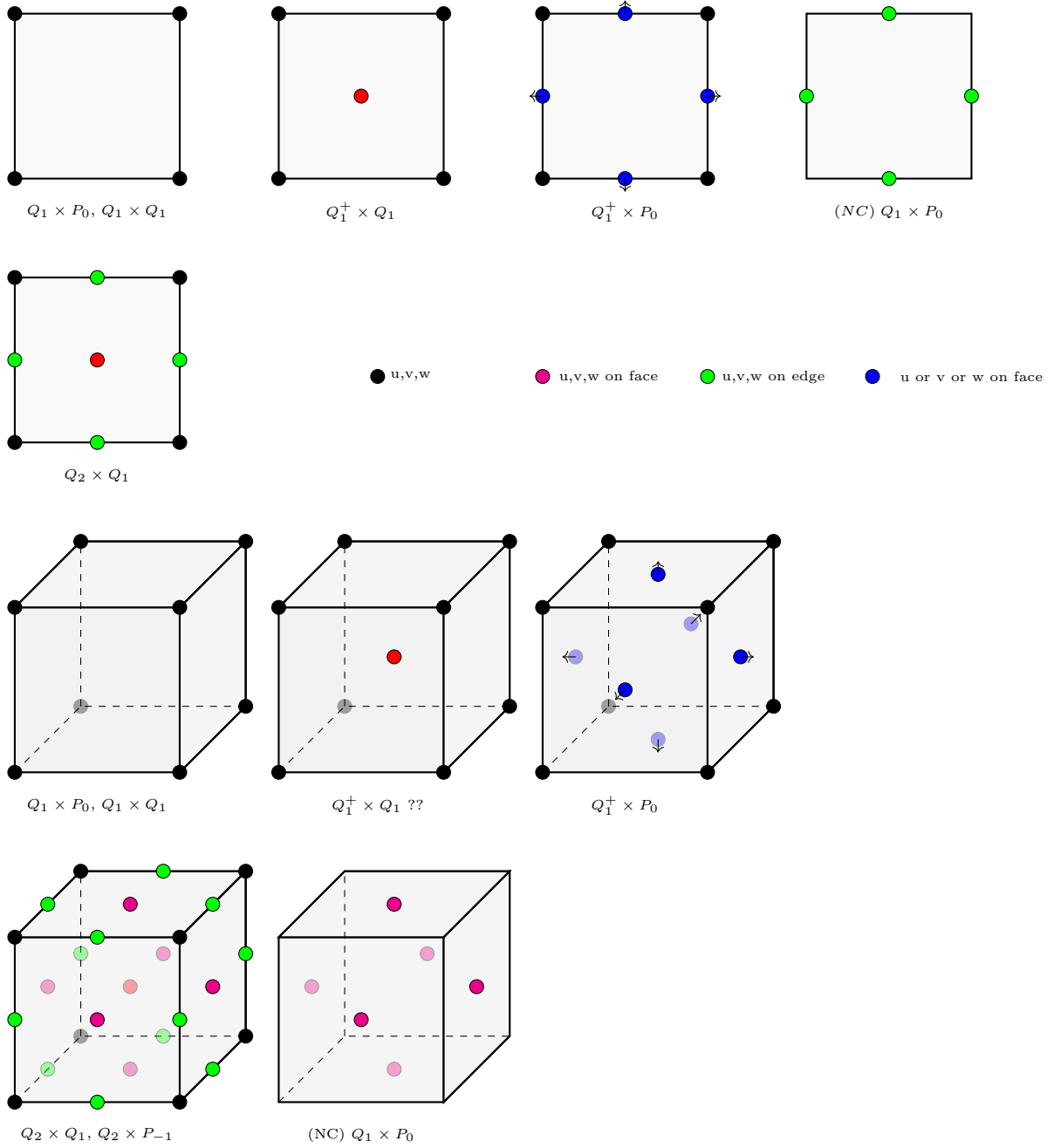
$$ndof_{2D} = 2[(nel_y + 1) \cdot nel_x + nel_y \cdot (nel_x + 1)] \rightarrow 4 \cdot nel_x^2$$

$$ndof_{3D} = 3[nel_x \cdot (nel_y + 1) \cdot nel_z + (nel_x + 1) \cdot nel_y \cdot nel_z + nel_x \cdot nel_y \cdot (nel_z + 1)] \rightarrow 9 \cdot nel_x^3$$

If we now assume $nel_x = nel_y = nel_z$, we can then plot the values above as a function of nel_x :



We see that the $Q_1^+ \times Q_1$ and $Q_1^+ \times P_0$ actually yield the same number of velocity dofs.
Simply based on the dof count and wishing for a (bi/tri)linear approximation for pressure, we must conclude that the $Q_1 \times Q_{1+2}$ bubbles is the most desirable since it is also LBB stable.



Add DSSY, RT, Q1Q1+2 bubbles

2D

| | | | |
|-----------------|--------|--------------------|------------------------------|
| Rannacher-Turek | NCQ1P0 | STONE 77 | pb with buoyancy-driven flow |
| Lamichhane | Q1+Q1 | STONE 72, STONE 74 | |
| DSSY | | STONE 77 | pb with buoyancy-driven flow |
| Fortin | | STONE 80 | |

3D

| | | | |
|-----------------|--------|----------|---------------|
| Rannacher-Turek | NCQ1P0 | ELEFANT | Does not work |
| Lamichhane | Q1+Q1 | STONE | |
| DSSY | | | |
| Fortin | | STONE 81 | |

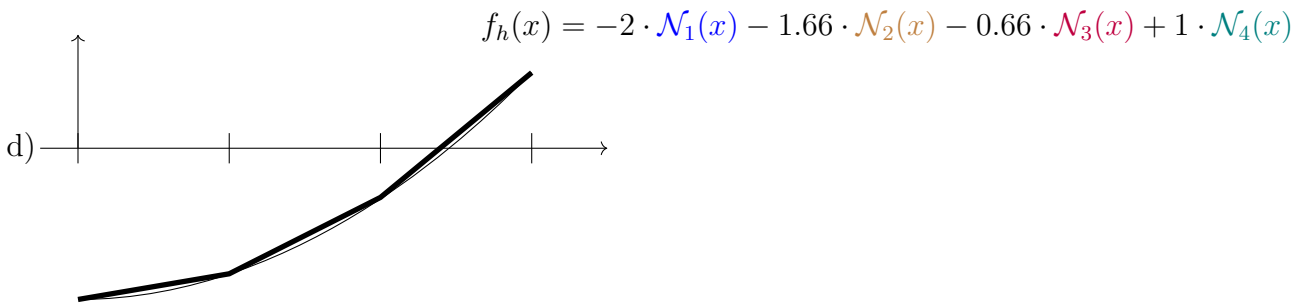
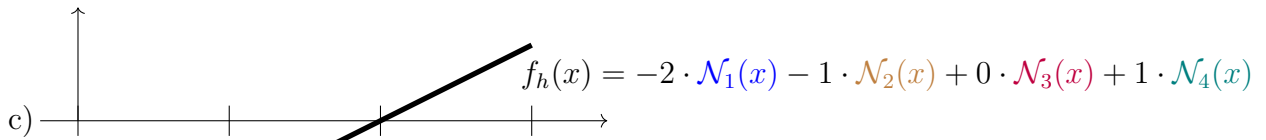
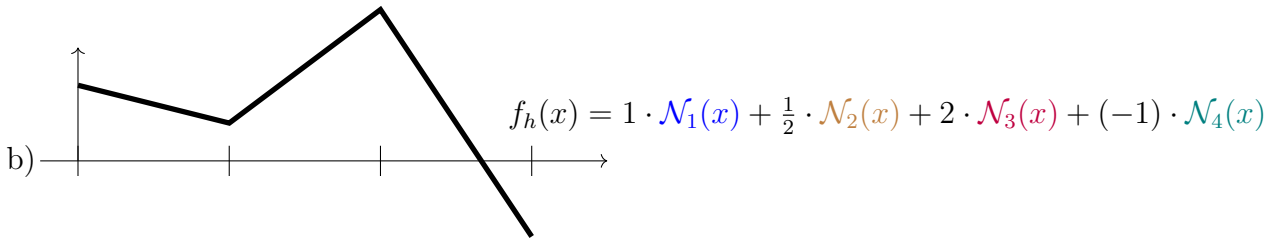
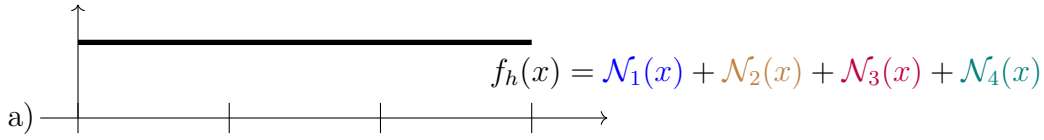
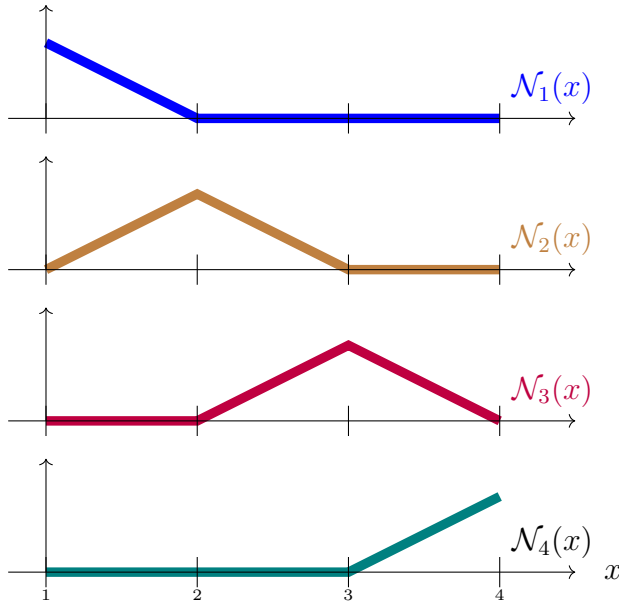
5.6 On the meaning of basis functions

basis_functions_meaning.tex.tex

In one dimension

Let us consider a 1D domain subdivided in 3 elements. We then consider the four linear basis functions attached to each node. On the following sketch these are not depicted on a single element with reduced coordinates but instead for the whole domain in the natural coordinate x :

(tikz_basisfunctions.tex)



The four cases a,b,c,d are examples of combinations of these basis functions:

$$f^h(x) = \sum_{i=1}^4 \mathcal{N}_i(x) f_i$$

Where f_i are the values associated to the four nodes. We assume that the distance h between nodes

is 1.

Example a) illustrates the fact that the sum of all basis functions must be strictly equal to one everywhere in the domain. Failing to do so would mean that the basis functions cannot represent a constant field (see Section 5.3.1).

Example b) illustrates a somewhat random combination of the basis functions, yielding a broken line.

Example c) illustrates the fact that these linear basis functions can exactly represent a linear function. When $f(x) = x - 2$, then $f_1 = f(0) = -2$, $f_2 = f(1) = -1$, $f_3 = f(2) = 0$ and $f_4 = f(3) = +1$, then $f^h(x)$ is exactly $f(x)$ on the domain.

Example d) illustrates the fact that linear basis functions cannot represent a parabola. Smaller and smaller elements will do an increasingly better job and will get closer to the curve but a systematic error will subsist.

Note that these drawings are trivial to produce since $\mathcal{N}_i(x_j) = \delta_{ij}$ by definition, so that $f^h(x_j) = f_j$.

In two dimensions

Chapter 6

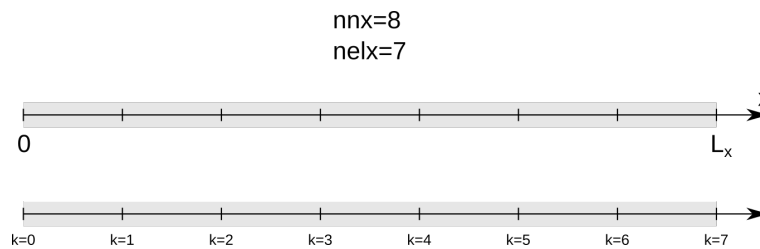
Solving the heat transport equation with linear Finite Elements

chapter5.tex

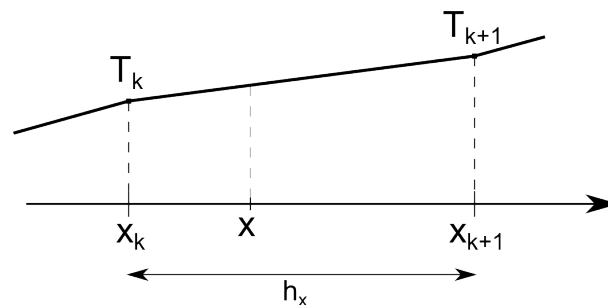
6.1 The diffusion equation in 1D

diff1D.tex

Let us consider the following one-dimensional grid:



It spans the domain Ω of length L_x . It is discretised by means of nnx nodes and $nelx = nnx - 1$ elements. Zooming in on the element e which is bounded by two nodes k and $k + 1$, its size (also sometimes called diameter) is $h_x = x_{k+1} - x_k$, and the temperature field we wish to compute is located on those nodes so that they are logically called T_k and T_{k+1} :



Remark. In what follows I will indicate mathematical functions¹ by this color.

Remark. In what follows \int_{Ω_e} should be understood as $\int_{x_k}^{x_{k+1}}$.

¹[https://en.wikipedia.org/wiki/Function_\(mathematics\)](https://en.wikipedia.org/wiki/Function_(mathematics))

From the strong form to the weak form

We focus here on the 1D diffusion equation (no advection, no heat sources, see Section 2.6):

$$\rho C_p \frac{\partial T}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) \quad (6.1)$$

where ρ is the density, C_p the heat capacity and k the heat conductivity. All three coefficients are assumed to be constant in space and time.

This is the **strong form** of the PDE to solve. I can multiply this equation by a function² $f(x)$ and integrate it over Ω :

$$\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx = \int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx \quad (6.2)$$

Looking at the right hand side, it is of the form $\int uv'$ so that I integrate it by parts³:

$$\int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx = \left[f(x) k \frac{\partial T}{\partial x} \right]_{\partial\Omega} - \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx \quad (6.3)$$

Assuming there is no heat flux on the boundary⁴ (i.e. $q_x = -k\partial T/\partial x = 0$), then

$$\int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx = - \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx \quad (6.4)$$

We then obtain the **weak form** of the diffusion equation in 1D:

$$\boxed{\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx + \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx = 0} \quad (6.5)$$

We then use the additive property of the integral $\int_{\Omega} \dots = \sum_{elts} \int_{\Omega_e} \dots$ so that the equation above becomes

$$\sum_{elts} \left(\underbrace{\int_{\Omega_e} f(x) \rho C_p \frac{\partial T}{\partial t} dx}_{\Lambda_f^e} + \underbrace{\int_{\Omega_e} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx}_{\Upsilon_f^e} \right) = 0 \quad (6.6)$$

From the weak form to a linear system

In order to compute these integrals (analytically or by means of a numerical quadrature), we will need to evaluate T inside the element. However, inside the element, the temperature is not known: all we hope to have at some point is the temperature at the nodes (this is what the code will compute).

For $x \in [x_k, x_{k+1}]$ we need to come up with a way to formulate the temperature in this element and we coin this T^h . It makes sense to think that $T^h(x)$ will then be a function of the temperature at the nodes, i.e. $T^h(x) = \alpha T_k + \beta T_{k+1}$ where α and β are coefficients. One over-simplified approach would be to assign $T^h(x) = (T_k + T_{k+1})/2$ (i.e. T^h is a zero-th order polynomial) but this would make the temperature discontinuous from element to element so we discard this option. A rather logical and simple solution to this problem is a linear temperature field between T_k and T_{k+1} :

²This function should be well-behaved with special properties, but we here assume it is a polynomial function.

³ $\int_{\Omega} uv' = \int_{\partial\Omega} uv - \int_{\Omega} u'v$

⁴This is of course not always the case and we will revisit this at a later stage in the course.

$$T^h(x) = \underbrace{\frac{x_{k+1} - x}{h_x}}_{\mathcal{N}_k^\theta(x)} T_k + \underbrace{\frac{x - x_k}{h_x}}_{\mathcal{N}_{k+1}^\theta(x)} T_{k+1}$$

where $\mathcal{N}_k^\theta(x)$ is the (temperature) basis function associated to node k and $\mathcal{N}_{k+1}^\theta(x)$ is the basis function associated to node $k+1$.

Rather reassuringly, we have:

- $x = x_k$ yields $T^h(x_k) = T_k$
- $x = x_{k+1}$ yields $T^h(x_{k+1}) = T_{k+1}$
- $x = x_{1/2} = (x_k + x_{k+1})/2$ yields $T^h(x_{1/2}) = (T_k + T_{k+1})/2$

In what follows we abbreviate $\partial T^h / \partial t$ by $\dot{T}^h(x)$. Let us compute Λ_f^e and Υ_f^e separately.

$$\begin{aligned} \Lambda_f^e &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p \dot{T}^h(x) dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p [\mathcal{N}_k^\theta(x) \dot{T}_k + \mathcal{N}_{k+1}^\theta(x) \dot{T}_{k+1}] dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p \mathcal{N}_k^\theta(x) \dot{T}_k dx + \int_{x_k}^{x_{k+1}} f(x) \rho C_p \mathcal{N}_{k+1}^\theta(x) \dot{T}_{k+1} dx \\ &= \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p \mathcal{N}_k^\theta(x) dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p \mathcal{N}_{k+1}^\theta(x) dx \right) \dot{T}_{k+1} \end{aligned}$$

Taking $f(x) = \mathcal{N}_k^\theta(x)$ and omitting ' (x) ' in the rhs:

$$\Lambda_{\mathcal{N}_k^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_k^\theta \mathcal{N}_k^\theta dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_k^\theta \mathcal{N}_{k+1}^\theta dx \right) \dot{T}_{k+1}$$

Taking $f(x) = \mathcal{N}_{k+1}^\theta(x)$ and omitting ' (x) ' in the rhs:

$$\Lambda_{\mathcal{N}_{k+1}^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_{k+1}^\theta \mathcal{N}_k^\theta dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_{k+1}^\theta \mathcal{N}_{k+1}^\theta dx \right) \dot{T}_{k+1}$$

We can rearrange these last two equations as follows:

$$\begin{pmatrix} \Lambda_{\mathcal{N}_k^\theta}^e \\ \Lambda_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \mathcal{N}_k^\theta \rho C_p \mathcal{N}_k^\theta dx & \int_{x_k}^{x_{k+1}} \mathcal{N}_k^\theta \rho C_p \mathcal{N}_{k+1}^\theta dx \\ \int_{x_k}^{x_{k+1}} \mathcal{N}_{k+1}^\theta \rho C_p \mathcal{N}_k^\theta dx & \int_{x_k}^{x_{k+1}} \mathcal{N}_{k+1}^\theta \rho C_p \mathcal{N}_{k+1}^\theta dx \end{pmatrix} \cdot \begin{pmatrix} \dot{T}_k \\ \dot{T}_{k+1} \end{pmatrix}$$

and we can take the integrals outside of the matrix:

$$\begin{pmatrix} \Lambda_{\mathcal{N}_k^\theta}^e \\ \Lambda_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} \rho C_p \begin{pmatrix} \mathcal{N}_k^\theta \mathcal{N}_k^\theta & \mathcal{N}_k^\theta \mathcal{N}_{k+1}^\theta \\ \mathcal{N}_{k+1}^\theta \mathcal{N}_k^\theta & \mathcal{N}_{k+1}^\theta \mathcal{N}_{k+1}^\theta \end{pmatrix} dx \right] \cdot \begin{pmatrix} \dot{T}_k \\ \dot{T}_{k+1} \end{pmatrix}$$

Finally, we can define the vectors

$$\vec{\mathcal{N}}^T = \begin{pmatrix} \mathcal{N}_k^\theta(x) \\ \mathcal{N}_{k+1}^\theta(x) \end{pmatrix}$$

and

$$\vec{T}^e = \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix} \quad \dot{\vec{T}}^e = \begin{pmatrix} \dot{T}_k \\ \dot{T}_{k+1} \end{pmatrix}$$

so that

$$\begin{pmatrix} \Lambda_{\mathcal{N}_k^\theta}^e \\ \Lambda_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{\mathcal{N}}^T \rho C_p \vec{\mathcal{N}} dx \right) \cdot \dot{\vec{T}}^e$$

Let us now go back to the diffusion term:

$$\begin{aligned} \Upsilon_f^e &= \int_{x_k}^{x_{k+1}} \frac{\partial f}{\partial x} k \frac{\partial T^h}{\partial x} dx \\ &= \int_{x_k}^{x_{k+1}} \frac{\partial f}{\partial x} k \frac{\partial (\mathcal{N}_k^\theta(x) T_k + \mathcal{N}_{k+1}^\theta(x) T_{k+1})}{\partial x} dx \\ &= \left(\int_{x_k}^{x_{k+1}} \frac{\partial f}{\partial x} k \frac{\partial \mathcal{N}_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} \frac{\partial f}{\partial x} k \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} dx \right) T_{k+1} \end{aligned}$$

Taking $f(x) = \mathcal{N}_k^\theta(x)$

$$\Upsilon_{\mathcal{N}_k^\theta}^e = \left(\int_{x_k}^{x_{k+1}} k \frac{\partial \mathcal{N}_k^\theta}{\partial x} \frac{\partial \mathcal{N}_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} k \frac{\partial \mathcal{N}_k^\theta}{\partial x} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} dx \right) T_{k+1}$$

Taking $f(x) = \mathcal{N}_{k+1}^\theta(x)$

$$\Upsilon_{\mathcal{N}_{k+1}^\theta}^e = \left(\int_{x_k}^{x_{k+1}} k \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \frac{\partial \mathcal{N}_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} k \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} dx \right) T_{k+1}$$

$$\begin{pmatrix} \Upsilon_{\mathcal{N}_k^\theta}^e \\ \Upsilon_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \frac{\partial \mathcal{N}_k^\theta}{\partial x} k \frac{\partial \mathcal{N}_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \frac{\partial \mathcal{N}_k^\theta}{\partial x} k \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} dx \\ \int_{x_k}^{x_{k+1}} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} k \frac{\partial \mathcal{N}_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} k \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} dx \end{pmatrix} \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

or,

$$\begin{pmatrix} \Upsilon_{\mathcal{N}_k^\theta}^e \\ \Upsilon_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} k \begin{pmatrix} \frac{\partial \mathcal{N}_k^\theta}{\partial x} \frac{\partial \mathcal{N}_k^\theta}{\partial x} & \frac{\partial \mathcal{N}_k^\theta}{\partial x} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \\ \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \frac{\partial \mathcal{N}_k^\theta}{\partial x} & \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \end{pmatrix} dx \right] \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

Finally, we can define the vector

$$\vec{B}^T = \begin{pmatrix} \frac{\partial \mathcal{N}_k^\theta}{\partial x} \\ \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \end{pmatrix}$$

so that

$$\begin{pmatrix} \Upsilon_{\mathcal{N}_k^\theta}^e \\ \Upsilon_{\mathcal{N}_{k+1}^\theta}^e \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx \right) \cdot \vec{T}^e$$

The weak form discretised over 1 element becomes

$$\underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{\mathcal{N}}^T \rho C_p \vec{\mathcal{N}} dx \right)}_{\mathbf{M}^e} \cdot \dot{\vec{T}}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx \right)}_{\mathbf{K}_d^e} \cdot \vec{T}^e = 0$$

or,

$$\boxed{\mathbf{M}^e \cdot \dot{\vec{T}}^e + \mathbf{K}_d^e \cdot \vec{T}^e = 0}$$

or,

$$\boxed{\mathbf{M}^e \cdot \frac{\partial \vec{T}^e}{\partial t} + \mathbf{K}_d^e \cdot \vec{T}^e = 0}$$

\mathbf{M}^e is commonly called the **mass matrix**, or capacitance matrix [1051, p103]. Note that the matrices are not coloured: the x dependence has disappeared when the integration was carried out.

In what follows I will omit the e superscript on the \vec{T} term to simplify notations.

We use a first order in time discretisation for the time derivative:

$$\dot{\vec{T}} = \frac{\partial \vec{T}}{\partial t} = \frac{\vec{T}^{new} - \vec{T}^{old}}{\delta t}$$

and in the context of an implicit scheme we get

$$\mathbf{M}^e \cdot \frac{\vec{T}^{new} - \vec{T}^{old}}{\delta t} + \mathbf{K}_d^e \cdot \vec{T}^{new} = 0$$

or,

$$\boxed{(\mathbf{M}^e + \mathbf{K}_d^e \delta t) \cdot \vec{T}^{new} = \mathbf{M}^e \cdot \vec{T}^{old}}$$

with

$$\mathbf{M}^e = \int_{x_k}^{x_{k+1}} \vec{\mathcal{N}}^T \rho C_p \vec{\mathcal{N}} dx \quad \mathbf{K}_d^e = \int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx$$

Computing the elemental matrices

Let us compute \mathbf{M}^e for an element:

$$\mathbf{M}^e = \int_{x_k}^{x_{k+1}} \vec{\mathcal{N}}^T \rho C_p \vec{\mathcal{N}} dx = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_k^\theta \mathcal{N}_k^\theta dx & \int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_k^\theta \mathcal{N}_{k+1}^\theta dx \\ \int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_{k+1}^\theta \mathcal{N}_k^\theta dx & \int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_{k+1}^\theta \mathcal{N}_{k+1}^\theta dx \end{pmatrix} = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix}$$

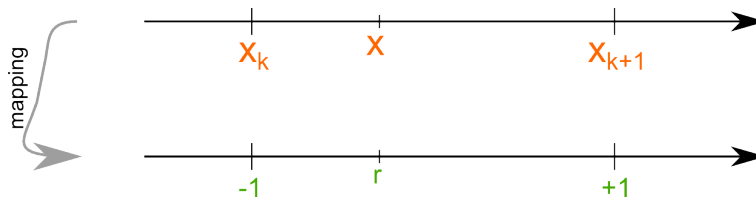
with

$$\vec{\mathcal{N}}^T = \begin{pmatrix} \mathcal{N}_k^\theta(x) \\ \mathcal{N}_{k+1}^\theta(x) \end{pmatrix} = \begin{pmatrix} \frac{x_{k+1}-x}{h_x} \\ \frac{x-x_k}{h_x} \end{pmatrix}$$

We only need to compute 3 integrals since $M_{12} = M_{21}$. Let us start with M_{11} :

$$M_{11} = \int_{x_k}^{x_{k+1}} \rho C_p \mathcal{N}_k^\theta(x) \mathcal{N}_k^\theta(x) dx = \int_{x_k}^{x_{k+1}} \rho C_p \frac{x_{k+1}-x}{h_x} \frac{x_{k+1}-x}{h_x} dx$$

It is then customary to carry out the change of variable $x \rightarrow r$ where $r \in [-1 : 1]$ as shown hereunder:



The relationships between x and r are:

$$r = \frac{2}{h_x}(x - x_k) - 1 \quad x = \frac{h_x}{2}(1 + r) + x_k$$

In what follows we assume for simplicity that ρ and C_p are constant within each element so that:

$$M_{11} = \rho C_p \int_{x_k}^{x_{k+1}} \frac{x_{k+1} - x}{h_x} \frac{x_{k+1} - x}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 - r)(1 - r) dr = \rho C_p \frac{h_x}{3}$$

Similarly we arrive at

$$\begin{aligned} M_{12} &= \rho C_p \int_{x_k}^{x_{k+1}} \frac{x_{k+1} - x}{h_x} \frac{x - x_k}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 - r)(1 + r) dr = \rho C_p \frac{h_x}{6} \\ M_{22} &= \rho C_p \int_{x_k}^{x_{k+1}} \frac{x - x_k}{h_x} \frac{x - x_k}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 + r)(1 + r) dr = \rho C_p \frac{h_x}{3} \end{aligned}$$

Finally

$$\boxed{\mathbf{M}^e = \rho C_p \frac{h_x}{3} \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}}$$

In the new coordinate system, the **basis functions**

$$\mathcal{N}_k^\theta(x) = \frac{x_{k+1} - x}{h_x} \quad \mathcal{N}_{k+1}^\theta(x) = \frac{x - x_k}{h_x}$$

become

$$\mathcal{N}_k^\theta(r) = \frac{1}{2}(1 - r) \quad \mathcal{N}_{k+1}^\theta(r) = \frac{1}{2}(1 + r)$$

Also,

$$\frac{\partial \mathcal{N}_k^\theta}{\partial x} = -\frac{1}{h_x} \quad \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} = \frac{1}{h_x}$$

so that

$$\vec{B}^T = \begin{pmatrix} \frac{\partial \mathcal{N}_k^\theta}{\partial x} \\ \frac{\partial \mathcal{N}_{k+1}^\theta}{\partial x} \end{pmatrix} = \begin{pmatrix} -\frac{1}{h_x} \\ \frac{1}{h_x} \end{pmatrix}$$

We here also assume that the heat conductivity k is constant within the element:

$$\mathbf{K}_d^e = \int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx = k \int_{x_k}^{x_{k+1}} \vec{B}^T \vec{B} dx$$

simply becomes

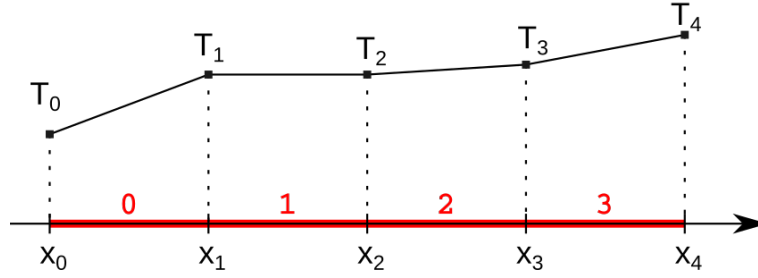
$$\mathbf{K}_d^e = k \int_{x_k}^{x_{k+1}} \frac{1}{h_x^2} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} dx$$

and then

$$\boxed{\mathbf{K}_d^e = \frac{k}{h_x} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}}$$

In practice

Let us consider this very simple grid consisting of 4 elements/5 nodes⁵:



For each element we have

$$\underbrace{(M^e + K_d^e \delta t)}_{A^e} \cdot \vec{T}^{new} = \underbrace{M^e \cdot \vec{T}^{old}}_{\vec{b}^e}$$

with

$$\vec{T}^{new} = \begin{pmatrix} T_k^{new} \\ T_{k+1}^{new} \end{pmatrix}$$

We can write this equation very explicitly for each element:

- element 0

$$A^0 \cdot \begin{pmatrix} T_0^{new} \\ T_1^{new} \end{pmatrix} = \vec{b}^0 \quad \rightarrow \quad \begin{cases} A_{00}^0 T_0^{new} + A_{01}^0 T_1^{new} = b_0^0 \\ A_{10}^0 T_0^{new} + A_{11}^0 T_1^{new} = b_1^0 \end{cases}$$

Element 0 is made of nodes 0 and 1, so it will contribute to lines and columns 0,1 of the FE matrix:

$$\begin{pmatrix} A_{00}^0 & A_{01}^0 & 0 & 0 & 0 \\ A_{10}^0 & A_{11}^0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} T_0^{new} \\ T_1^{new} \\ T_2^{new} \\ T_3^{new} \\ T_4^{new} \end{pmatrix} = \begin{pmatrix} b_0^0 \\ b_1^0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

- element 1

$$A^1 \cdot \begin{pmatrix} T_1^{new} \\ T_2^{new} \end{pmatrix} = \vec{b}^1 \quad \rightarrow \quad \begin{cases} A_{00}^1 T_1^{new} + A_{01}^1 T_2^{new} = b_0^1 \\ A_{10}^1 T_1^{new} + A_{11}^1 T_2^{new} = b_1^1 \end{cases}$$

Element 1 is made of nodes 1 and 2, so it will contribute to lines and columns 1,2 of the FE matrix:

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & A_{00}^1 & A_{01}^1 & 0 & 0 \\ 0 & A_{10}^1 & A_{11}^1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} T_0^{new} \\ T_1^{new} \\ T_2^{new} \\ T_3^{new} \\ T_4^{new} \end{pmatrix} = \begin{pmatrix} 0 \\ b_0^1 \\ b_1^1 \\ 0 \\ 0 \end{pmatrix}$$

- element 2

$$A^2 \cdot \begin{pmatrix} T_2^{new} \\ T_3^{new} \end{pmatrix} = \vec{b}^2 \quad \rightarrow \quad \begin{cases} A_{00}^2 T_2^{new} + A_{01}^2 T_3^{new} = b_0^2 \\ A_{10}^2 T_2^{new} + A_{11}^2 T_3^{new} = b_1^2 \end{cases}$$

Element 2 is made of nodes 2 and 3, so it will contribute to lines and columns 2,3 of the FE matrix:

⁵I have here adopted the illogical numbering of python – the language used by my students– so that python programmers can benchmark their code against this simple example

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & A_{00}^2 & A_{01}^2 & 0 \\ 0 & 0 & A_{10}^2 & A_{11}^2 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} T_0^{new} \\ T_1^{new} \\ T_2^{new} \\ T_3^{new} \\ T_4^{new} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ b_0^2 \\ b_1^2 \\ 0 \end{pmatrix}$$

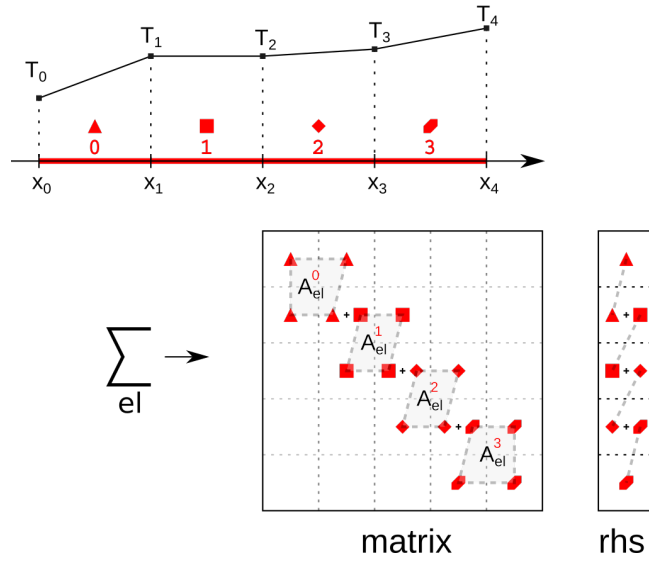
- element 3

$$\mathbf{A}^3 \cdot \begin{pmatrix} T_3^{new} \\ T_4^{new} \end{pmatrix} = \tilde{b}^3 \rightarrow \begin{cases} A_{00}^3 T_3^{new} + A_{01}^3 T_4^{new} = b_0^3 \\ A_{10}^3 T_3^{new} + A_{11}^3 T_4^{new} = b_1^3 \end{cases}$$

Element 3 is made of nodes 3 and 4, so it will contribute to lines and columns 3,4 of the FE matrix:

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & A_{00}^3 & A_{01}^3 \\ 0 & 0 & 0 & A_{10}^3 & A_{11}^3 \end{pmatrix} \cdot \begin{pmatrix} T_0^{new} \\ T_1^{new} \\ T_2^{new} \\ T_3^{new} \\ T_4^{new} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ b_0^3 \\ b_1^3 \end{pmatrix}$$

All equations can be cast into a single linear system: this is the **assembly** phase. The process can also be visualised as shown hereunder. Because nodes 2,3,4 belong to two elements elemental contributions will be summed in the matrix and the rhs:



The assembled matrix and rhs are then (I have dropped the *new* superscripts):

$$\begin{pmatrix} A_{00}^0 & A_{01}^0 & 0 & 0 & 0 \\ A_{10}^0 & A_{11}^0 + A_{00}^1 & A_{01}^1 & 0 & 0 \\ 0 & A_{10}^1 & A_{11}^1 + A_{00}^2 & A_{01}^2 & 0 \\ 0 & 0 & A_{10}^2 & A_{11}^2 + A_{00}^3 & A_{01}^3 \\ 0 & 0 & 0 & A_{10}^3 & A_{11}^3 \end{pmatrix} \begin{pmatrix} T_0 \\ T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix} = \begin{pmatrix} b_0^0 \\ b_1^0 + b_0^1 \\ b_1^1 + b_0^2 \\ b_1^2 + b_0^3 \\ b_1^3 \end{pmatrix}$$

Ultimately the assembled matrix system also takes the form

$$\begin{pmatrix} \mathcal{A}_{00} & \mathcal{A}_{01} & 0 & 0 & 0 \\ \mathcal{A}_{10} & \mathcal{A}_{11} & \mathcal{A}_{12} & 0 & 0 \\ 0 & \mathcal{A}_{21} & \mathcal{A}_{22} & \mathcal{A}_{23} & 0 \\ 0 & 0 & \mathcal{A}_{32} & \mathcal{A}_{33} & \mathcal{A}_{34} \\ 0 & 0 & 0 & \mathcal{A}_{43} & \mathcal{A}_{44} \end{pmatrix} \begin{pmatrix} T_0 \\ T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ b_4 \end{pmatrix}$$

and we see that it is sparse. Its sparsity structure is easy to derive: each row corresponds to a dof, and since nodes 1 and 2 'see' each other (they belong to the same element) there will be non-zero entries in the first and second column. Likewise, node 2 'sees' node 1 (in other words, there is an edge linking nodes 1 and 2), itself, and node 3, so that there are non-zero entries in the second row at columns 1, 2, and 3.

Before we solve the system, we need to take care of boundary conditions. Let us assume that we wish to fix the temperature at node 2, or in other words we wish to set

$$T_2 = T^o$$

This equation can be cast as

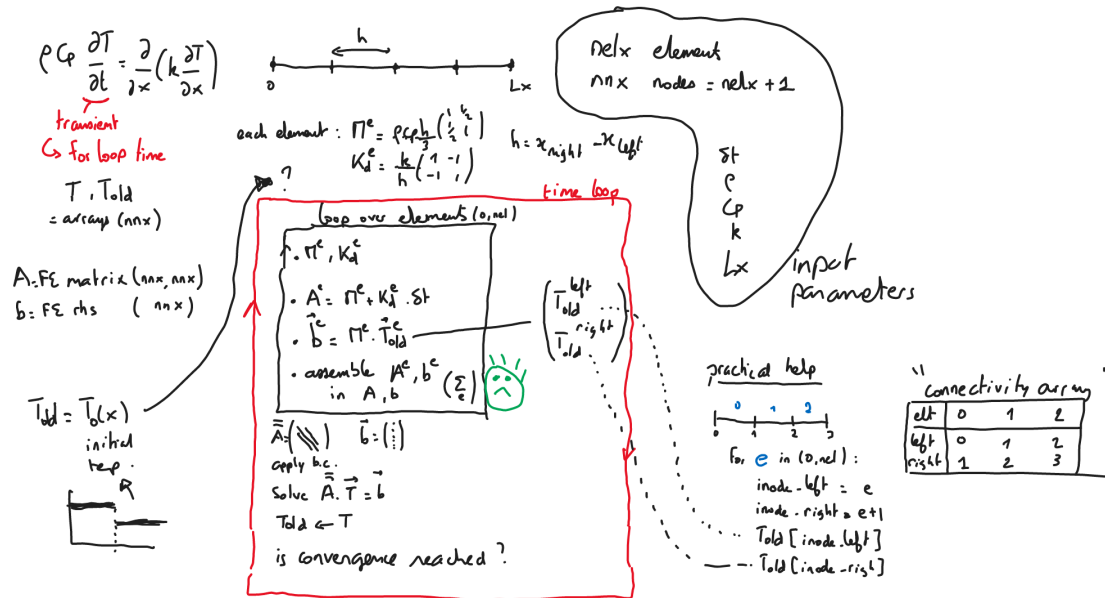
$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} 0 \\ T^o \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

This replaces the second line in the previous matrix equation:

$$\begin{pmatrix} \mathcal{A}_{11} & \mathcal{A}_{12} & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & \mathcal{A}_{32} & \mathcal{A}_{33} & \mathcal{A}_{34} & 0 \\ 0 & 0 & \mathcal{A}_{43} & \mathcal{A}_{44} & \mathcal{A}_{45} \\ 0 & 0 & 0 & \mathcal{A}_{54} & \mathcal{A}_{55} \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} b_1 \\ T^o \\ b_3 \\ b_4 \\ b_5 \end{pmatrix}$$

That's it, we have a linear system of equations which can be solved!

The following figure presents a hand-drawn template of how a typical 1D FE code is structured:



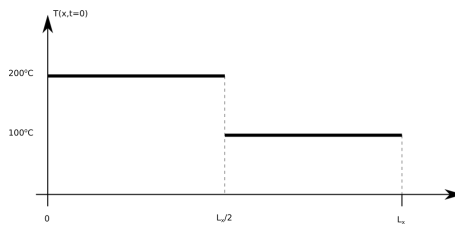
Additional comments: A and b should be zeroed at every time step. Convergence should be tested before Told received T.

redo this for 2D



Exercise FEM-01

Write a code which solves the 1D diffusion equation in time. The initial temperature field is as follows:



$$T(x, t = 0) = 200 \quad x < L_x/2 \quad T(x, t = 0) = 100 \quad x \geq L_x/2$$

The domain is $L_x = 100\text{km}$ and the properties of the material are $\rho = 3000\text{kg/m}^3$, $k = 3\text{W/m/K}$, $C_p = 1000\text{J/K}$. Boundary conditions are:

$$T(t, x = 0) = 200^\circ\text{C} \quad T(t, x = L_x) = 100^\circ\text{C}$$

There are `nelx` elements and `nnx` nodes. All elements are `hx` long. The code will carry out `nstep` timesteps of length `dt` with $\delta t = 0.5 \frac{h_x^2}{\kappa}$.

Bonus: add a small random perturbation up to $\pm 20\%$ of h_x to each node position inside the domain and make sure that the recovered steady state solution is unchanged.

6.2 The advection-diffusion equation in 1D

We start with the 1D advection-diffusion equation

$$\rho C_p \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} \right) = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) + H \quad (6.7)$$

This is the **strong form** of the ODE to solve. As in the previous section, I multiply this equation by a function $f(x)$ and integrate it over the domain Ω :

$$\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx + \int_{\Omega} f(x) \rho C_p u \frac{\partial T}{\partial x} dx = \int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx + \int_{\Omega} f(x) H dx$$

As in the previous section I integrate the r.h.s. by parts:

$$\int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx = \left[f(x) k \frac{\partial T}{\partial x} \right]_{\partial \Omega} - \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx$$

Disregarding the boundary term for now, we then obtain the **weak form** of the diffusion equation in 1D:

$$\boxed{\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx + \int_{\Omega} f(x) \rho C_p u \frac{\partial T}{\partial x} dx + \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx = \int_{\Omega} f(x) H dx}$$

We then use the additive property of the integral $\int_{\Omega} \dots = \sum_{elts} \int_{\Omega_e} \dots$

$$\sum_{elts} \left(\underbrace{\int_{\Omega_e} f(x) \rho C_p \frac{\partial T}{\partial t} dx}_{\Lambda_f^e} + \underbrace{\int_{\Omega_e} f(x) \rho C_p u \frac{\partial T}{\partial x} dx}_{\Sigma_f^e} + \underbrace{\int_{\Omega_e} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx}_{\Upsilon_f^e} - \underbrace{\int_{\Omega_e} f(x) H dx}_{\Omega_f^e} \right) = 0$$

replace omega by other letter for source term

In the element, we have seen that the temperature can be written:

$$T^h(x) = \mathbf{N}_k^{\theta}(x) T_k + \mathbf{N}_{k+1}^{\theta}(x) T_{k+1}$$

In the previous presentation we have computed Λ_f^e and Υ_f^e . Let us now turn to Σ_f^e and Ω_f^e .

$$\begin{aligned} \Sigma_f^e &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial T^h}{\partial x} dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial [\mathbf{N}_k^{\theta}(x) T_k + \mathbf{N}_{k+1}^{\theta}(x) T_{k+1}]}{\partial x} dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial \mathbf{N}_k^{\theta}}{\partial x} T_k dx + \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial \mathbf{N}_{k+1}^{\theta}}{\partial x} T_{k+1} dx \\ &= \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial \mathbf{N}_k^{\theta}}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial \mathbf{N}_{k+1}^{\theta}}{\partial x} dx \right) T_{k+1} \end{aligned}$$

Taking $f(x) = \mathbf{N}_k^{\theta}(x)$ and omitting ' (x) ' in the rhs:

$$\Sigma_{\mathbf{N}_k^{\theta}}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_k^{\theta} \frac{\partial \mathbf{N}_k^{\theta}}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_k^{\theta} \frac{\partial \mathbf{N}_{k+1}^{\theta}}{\partial x} dx \right) T_{k+1}$$

Taking $f(x) = \mathbf{N}_{k+1}^{\theta}(x)$ and omitting ' (x) ' in the rhs:

$$\Sigma_{\mathbf{N}_{k+1}^{\theta}}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_{k+1}^{\theta} \frac{\partial \mathbf{N}_k^{\theta}}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_{k+1}^{\theta} \frac{\partial \mathbf{N}_{k+1}^{\theta}}{\partial x} dx \right) T_{k+1}$$

$$\begin{pmatrix} \Sigma_{N_k^\theta} \\ \Sigma_{N_{k+1}^\theta} \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_k^\theta \frac{\partial \mathbf{N}_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_k^\theta \frac{\partial \mathbf{N}_{k+1}^\theta}{\partial x} dx \\ \int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_{k+1}^\theta \frac{\partial \mathbf{N}_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \rho C_p u \mathbf{N}_{k+1}^\theta \frac{\partial \mathbf{N}_{k+1}^\theta}{\partial x} dx \end{pmatrix} \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

or,

$$\begin{pmatrix} \Sigma_{N_k^\theta} \\ \Sigma_{N_{k+1}^\theta} \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} \rho C_p u \begin{pmatrix} \mathbf{N}_k^\theta \frac{\partial \mathbf{N}_k^\theta}{\partial x} & \mathbf{N}_k^\theta \frac{\partial \mathbf{N}_{k+1}^\theta}{\partial x} \\ \mathbf{N}_{k+1}^\theta \frac{\partial \mathbf{N}_k^\theta}{\partial x} & \mathbf{N}_{k+1}^\theta \frac{\partial \mathbf{N}_{k+1}^\theta}{\partial x} \end{pmatrix} dx \right] \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

Finally, we have already defined the vectors

$$\vec{N}^T = \begin{pmatrix} \mathbf{N}_k^\theta(x) \\ \mathbf{N}_{k+1}^\theta(x) \end{pmatrix} \quad \vec{B}^T = \begin{pmatrix} \frac{\partial \mathbf{N}_k^\theta}{\partial x} \\ \frac{\partial \mathbf{N}_{k+1}^\theta}{\partial x} \end{pmatrix} \quad \vec{T}^e = \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

so that

$$\begin{pmatrix} \Sigma_{N_k^\theta} \\ \Sigma_{N_{k+1}^\theta} \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p u \vec{B} dx \right) \cdot \vec{T}^e = \mathbf{K}_a \cdot \vec{T}^e$$

One can easily show that

$$\mathbf{K}_a^e = \rho C_p u \begin{pmatrix} -1/2 & 1/2 \\ -1/2 & 1/2 \end{pmatrix}$$

Note that the matrix \mathbf{K}_a^e is *not* symmetric.

Let us now look at the source term:

$$\Omega_f^e = \int_{x_k}^{x_{k+1}} f(x) H(x) dx$$

Taking $f(x) = \mathbf{N}_k^\theta(x)$:

$$\Omega_{N_k^\theta} = \int_{x_k}^{x_{k+1}} \mathbf{N}_k^\theta(x) H(x) dx$$

Taking $f(x) = \mathbf{N}_{k+1}^\theta(x)$:

$$\Omega_{N_{k+1}^\theta} = \int_{x_k}^{x_{k+1}} \mathbf{N}_{k+1}^\theta(x) H(x) dx$$

We can rearrange both equations as follows:

$$\begin{pmatrix} \Omega_{N_k^\theta} \\ \Omega_{N_{k+1}^\theta} \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \mathbf{N}_k^\theta(x) H(x) dx \\ \int_{x_k}^{x_{k+1}} \mathbf{N}_{k+1}^\theta(x) H(x) dx \end{pmatrix}$$

or,

$$\begin{pmatrix} \Omega_{N_k^\theta} \\ \Omega_{N_{k+1}^\theta} \end{pmatrix} = \int_{x_k}^{x_{k+1}} \begin{pmatrix} \mathbf{N}_k^\theta(x) H(x) \\ \mathbf{N}_{k+1}^\theta(x) H(x) \end{pmatrix} dx = \left(\int_{x_k}^{x_{k+1}} \vec{N}^T H(x) dx \right)$$

The weak form discretised over 1 element becomes

$$\underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \vec{N} dx \right)}_{\mathbf{M}^e} \cdot \dot{\vec{T}}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p u \vec{B} dx \right)}_{\mathbf{K}_a^e} \cdot \vec{T}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{B}^T \rho C_p u \vec{B} dx \right)}_{\mathbf{K}_d^e} \cdot \vec{T}^e = \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T H(x) dx \right)}_{\vec{F}^e}$$

or,

$$\mathbf{M}^e \cdot \dot{\vec{T}}^e + (\mathbf{K}_d^e + \mathbf{K}_a^e) \cdot \vec{T}^e = \vec{F}^e$$

or,

$$\mathbf{M}^e \cdot \frac{\partial \vec{T}^e}{\partial t} + (\mathbf{K}_a^e + \mathbf{K}_d^e) \cdot \vec{T}^e = \vec{F}^e$$

As in the diffusion case of the previous section these matrices and vectors will need to be assembled into \mathbf{M} , \mathbf{K}_a , \mathbf{K}_d , \vec{T} and \vec{F} :

$$\mathbf{M} \cdot \frac{\partial \vec{T}}{\partial t} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = \vec{F}$$

We can revisit the time discretisation again, assuming for simplicity that the coefficients of the PDE are not time-dependent. Choosing a fully explicit approach would have us write

$$\mathbf{M} \cdot \frac{\vec{T}^{n+1} - \vec{T}^n}{\delta t} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T}^n = \vec{F} \quad \Rightarrow \quad \mathbf{M} \cdot \vec{T}^{n+1} = [\mathbf{M} - (\mathbf{K}_a + \mathbf{K}_d)\delta t] \cdot \vec{T}^n + \vec{F} \delta t \quad (6.8)$$

while choosing a fully implicit approach would have us write

$$\mathbf{M} \cdot \frac{\vec{T}^{n+1} - \vec{T}^n}{\delta t} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T}^{n+1} = \vec{F} \quad \Rightarrow \quad [\mathbf{M} + (\mathbf{K}_a + \mathbf{K}_d)\delta t] \cdot \vec{T}^{n+1} = \mathbf{M} \cdot \vec{T}^n + \vec{F} \delta t \quad (6.9)$$

We can also consider a more generic approach and write:

$$\mathbf{M} \cdot \frac{\vec{T}^{n+1} - \vec{T}^n}{\delta t} + (\mathbf{K}_a + \mathbf{K}_d) \cdot (\alpha \vec{T}^{n+1} + (1 - \alpha) \vec{T}^n) = \vec{F} \quad (6.10)$$

$$[\mathbf{M} + \alpha(\mathbf{K}_a + \mathbf{K}_d)\delta t] \cdot \vec{T}^{n+1} = [\mathbf{M} - (1 - \alpha)(\mathbf{K}_a + \mathbf{K}_d)\delta t] \cdot \vec{T}^n + \vec{F} \delta t \quad (6.11)$$

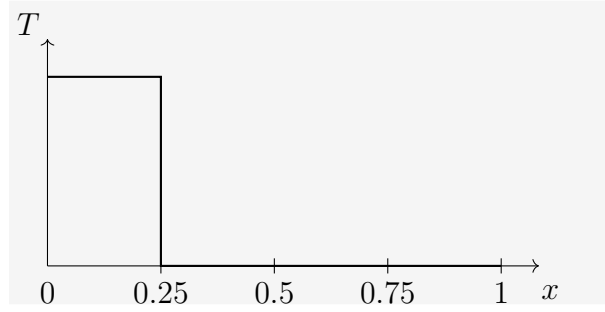
When $\alpha = 0$ we recover the explicit scheme, when $\alpha = 1$ we recover the implicit one, and when $\alpha = 1/2$ we get a so-called mid-point algorithm (Crank-Nicolson).

Write about SUPG

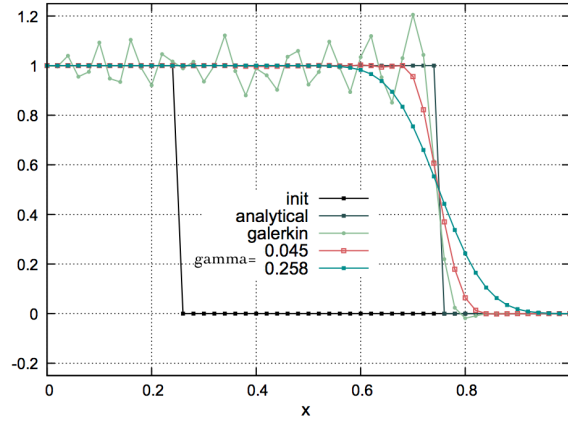


Exercise FEM-2

Let us consider the domain $[0, 1]$. The temperature field at $t = 0$ is given by $T = 1$ for $x < 0.25$ and $T = 0$ otherwise. The prescribed velocity is $u = 1$ and we set $nnx = 51$. Boundary conditions are $T = 1$ at $x = 0$ and $T = 0$ at $x = 1$. Only advection is present, no heat source nor diffusion.



Set $\rho = C_p = 1$. Run the model for 250 time steps with $\delta t = 0.002$. Implement a fully implicit, explicit and Crank-Nicolson time discretisation. When using Crank-Nicolson, you should then be able to recover the green line of the following figure:



Taken from Thieulot (2011) [1258]. Note that $\tau = \gamma h / u$.

Finally, implement the SUPG method and recover the red and turquoise lines.

6.3 The advection-diffusion equation in 2D

We start from the 'bare-bones' heat transport equation (source terms are omitted):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot (k \vec{\nabla} T) \quad (6.12)$$

In what follows we assume that the velocity field \vec{v} is known so that temperature is the only unknown. Let \mathcal{N}^θ be the temperature basis functions so that the temperature inside an element is given by⁶:

$$T^h(\vec{r}) = \sum_{i=1}^{m_T} \mathcal{N}_i^\theta(\vec{r}) T_i = \vec{\mathcal{N}}^\theta \cdot \vec{T} \quad (6.13)$$

where \vec{T} is a vector of length m_T . The weak form is then

$$\int_{\Omega} \mathcal{N}_i^\theta \left[\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) \right] d\Omega = \int_{\Omega} \mathcal{N}_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T d\Omega \quad (6.14)$$

$$\underbrace{\int_{\Omega} \mathcal{N}_i^\theta \rho C_p \frac{\partial T}{\partial t} d\Omega}_I + \underbrace{\int_{\Omega} \mathcal{N}_i^\theta \rho C_p \vec{v} \cdot \vec{\nabla} T d\Omega}_{II} = \underbrace{\int_{\Omega} \mathcal{N}_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T d\Omega}_{III} \quad i = 1, m_T \quad (6.15)$$

Looking at the first term:

$$\int_{\Omega} \mathcal{N}_i^\theta \rho C_p \frac{\partial T}{\partial t} d\Omega = \int_{\Omega} \mathcal{N}_i^\theta \rho C_p \vec{\mathcal{N}}^\theta \cdot \dot{\vec{T}} d\Omega \quad (6.16)$$

$$(6.17)$$

so that when we assemble all contributions for $i = 1, m_T$ we get:

$$I = \int_{\Omega} \vec{\mathcal{N}}^\theta \rho C_p \vec{\mathcal{N}}^\theta \cdot \dot{\vec{T}} d\Omega = \left(\int_{\Omega} \rho C_p \vec{\mathcal{N}}^\theta \vec{\mathcal{N}}^\theta d\Omega \right) \cdot \dot{\vec{T}} = \mathbf{M}^T \cdot \dot{\vec{T}} \quad (6.18)$$

where \mathbf{M}^T is the mass matrix of the system of size $(m_T \times m_T)$ with

$$M_{ij}^T = \int_{\Omega} \rho C_p \mathcal{N}_i^\theta \mathcal{N}_j^\theta d\Omega \quad (6.19)$$

Turning now to the second term:

$$\int_{\Omega} \mathcal{N}_i^\theta \rho C_p \vec{v} \cdot \vec{\nabla} T d\Omega = \int_{\Omega} \mathcal{N}_i^\theta \rho C_p \left(u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y} \right) d\Omega \quad (6.20)$$

$$= \int_{\Omega} \mathcal{N}_i^\theta \rho C_p \left(u \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} + v \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} \right) \cdot \vec{T} d\Omega \quad (6.21)$$

$$(6.22)$$

so that when we assemble all contributions for $i = 1, m_T$ we get:

$$II = \left(\int_{\Omega} \rho C_p \vec{\mathcal{N}}^\theta \left(u \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} + v \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} \right) d\Omega \right) \cdot \vec{T} = \mathbf{K}_a \cdot \vec{T}$$

where \mathbf{K}_a is the advection term matrix of size $(m_T \times m_T)$ with

$$(K_a)_{ij} = \int_{\Omega} \rho C_p \mathcal{N}_i^\theta \left(u \frac{\partial \mathcal{N}_j^\theta}{\partial x} + v \frac{\partial \mathcal{N}_j^\theta}{\partial y} \right) d\Omega$$

⁶the θ superscript has been chosen to denote temperature so as to avoid confusion with the transpose operator

Now looking at the third term, we carry out an integration by part and neglect the surface term for now, so that

$$\int_{\Omega} \mathcal{N}_i^{\theta} \vec{\nabla} \cdot k \vec{\nabla} T d\Omega = - \int_{\Omega} k \vec{\nabla} \mathcal{N}_i^{\theta} \cdot \vec{\nabla} T d\Omega \quad (6.23)$$

$$= - \int_{\Omega} k \vec{\nabla} \mathcal{N}_i^{\theta} \cdot \vec{\nabla} (\vec{\mathcal{N}}^{\theta} \cdot \vec{T}) d\Omega \quad (6.24)$$

$$(6.25)$$

with

$$\vec{\nabla} \vec{\mathcal{N}}^{\theta} = \begin{pmatrix} \partial_x \mathcal{N}_1^{\theta} & \partial_x \mathcal{N}_2^{\theta} & \dots & \partial_x \mathcal{N}_{m_T}^{\theta} \\ \partial_y \mathcal{N}_1^{\theta} & \partial_y \mathcal{N}_2^{\theta} & \dots & \partial_y \mathcal{N}_{m_T}^{\theta} \end{pmatrix}$$

so that finally:

$$III = - \left(\int_{\Omega} k (\vec{\nabla} \vec{\mathcal{N}}^{\theta})^T \cdot \vec{\nabla} \vec{\mathcal{N}}^{\theta} d\Omega \right) \cdot \vec{T} = -\mathbf{K}_d \cdot \vec{T}$$

where \mathbf{K}_d is the diffusion term matrix:

$$\mathbf{K}_d = \int_{\Omega} k (\vec{\nabla} \vec{\mathcal{N}}^{\theta})^T \cdot \vec{\nabla} \vec{\mathcal{N}}^{\theta} d\Omega$$

Ultimately terms I, II, III together yield:

$$\boxed{\mathbf{M}^{\theta} \cdot \dot{\vec{T}} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = \vec{0}}$$

add source term!!

add something about computing $\mathbf{K}_d, \mathbf{K}_a$ with quadrature!

On steady states

It is said that a system is in a steady state if the (state) variables which define the behavior of the system are unchanging in time. In continuous time, this means that the partial derivative with respect to time is zero and remains so:

$$\frac{\partial}{\partial t} = 0 \quad \forall t$$

This is irrelevant for the Stokes equations which do not contain an explicit time dependence but the heat transport equation can reach a steady state. Note that if one is only interested in the steady state solution (and not how the system gets there in time) then the heat transport equation should be solved with $\partial T / \partial t$ set to zero.

Dealing with the time discretisation

time_discretisation.tex

Essentially we have to solve a PDE of the type:

$$\frac{\partial T}{\partial t} = \mathcal{F}(\vec{\mathbf{v}}, T, \vec{\nabla} T, \Delta T)$$

with $\mathcal{F} = \frac{1}{\rho C_p} (-\vec{\mathbf{v}} \cdot \vec{\nabla} T + \vec{\nabla} \cdot k \vec{\nabla} T)$. The (explicit) forward Euler method is:

$$\frac{T^{n+1} - T^n}{\delta t} = \mathcal{F}^n(T, \vec{\nabla} T, \Delta T)$$

The (implicit) backward Euler method is:

$$\frac{T^{n+1} - T^n}{\delta t} = \mathcal{F}^{n+1}(T, \vec{\nabla}T, \Delta T)$$

and the (implicit) Crank-Nicolson algorithm is:

$$\frac{T^{n+1} - T^n}{\delta t} = \frac{1}{2} \left[\mathcal{F}^n(T, \vec{\nabla}T, \Delta T) + \mathcal{F}^{n+1}(T, \vec{\nabla}T, \Delta T) \right]$$

where the superscript n indicates the time step. The Crank-Nicolson is obviously based on the trapezoidal rule, with second-order convergence in time.

In what follows, I omit the superscript on the mass matrix to simplify notations: $\mathbf{M}^\theta = \mathbf{M}$. In terms of Finite Elements, these become:

- Explicit Forward euler:

$$\frac{1}{\delta t}(\mathbf{M}^{n+1} \cdot \vec{T}^{n+1} - \mathbf{M}^n \cdot \vec{T}^n) = -(\mathbf{K}_a^n + \mathbf{K}_d^n) \cdot \vec{T}^n$$

or,

$$\boxed{\mathbf{M}^{n+1} \cdot \vec{T}^{n+1} = (\mathbf{M}^n - (\mathbf{K}_a^n + \mathbf{K}_d^n)\delta t) \cdot \vec{T}^n}$$

- Implicit Backward euler:

$$\frac{1}{\delta t}(\mathbf{M}^{n+1} \cdot \vec{T}^{n+1} - \mathbf{M}^n \cdot \vec{T}^n) = -(\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1}) \cdot \vec{T}^{n+1}$$

or,

$$\boxed{(\mathbf{M}^{n+1} + (\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1})\delta t) \cdot \vec{T}^{n+1} = \mathbf{M}^n \cdot \vec{T}^n} \quad (6.26)$$

- Crank-Nicolson

$$\frac{1}{\delta t}(\mathbf{M}^{n+1} \cdot \vec{T}^{n+1} - \mathbf{M}^n \cdot \vec{T}^n) = \frac{1}{2} \left[-(\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1}) \cdot \vec{T}^{n+1} - (\mathbf{K}_a^n + \mathbf{K}_d^n) \cdot \vec{T}^n \right]$$

or,

$$\boxed{\left(\mathbf{M}^{n+1} + (\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1})\frac{\delta t}{2} \right) \cdot \vec{T}^{n+1} = \left(\mathbf{M}^n - (\mathbf{K}_a^n + \mathbf{K}_d^n)\frac{\delta t}{2} \right) \cdot \vec{T}^n}$$

Note that in benchmarks where the domain/grid does not deform, the coefficients do not change in space and the velocity field is constant in time, or in practice out of convenience, the \mathbf{K} and \mathbf{M} matrices do not change and the r.h.s. can be constructed with the same matrices as the FE matrix.

The Backward differentiation formula (see for instance Hairer & Wanner [520] or Wikipedia⁷. See also step-31 of deal.II⁸. The second-order BDF (or BDF-2) as shown in Kronbichler *et al.* (2012) [732] is as follows: it is a finite-difference quadratic interpolation approximation of the $\partial T / \partial t$ term which involves t^n , t^{n-1} and t^{n-2} :

$$\frac{\partial T}{\partial t}(t^n) \simeq \frac{1}{\delta t_n} \left(\frac{2\delta t_n + \delta t_{n-1}}{\delta t_n + \delta t_{n-1}} T^n - \frac{\delta t_n + \delta t_{n-1}}{\delta t_{n-1}} T^{n-1} + \frac{\delta t_n^2}{\delta t_{n-1}(\delta t_n + \delta t_{n-1})} T^{n-2} \right) \quad (6.27)$$

⁷https://en.wikipedia.org/wiki/Backward_differentiation_formula

⁸https://www.dealii.org/current/doxygen/deal.II/step_31.html

where $\delta t_n = t^n - t^{n-1}$. We also then have the approximation

$$T^n \simeq T^{n-1} + \delta t_n \frac{\partial T}{\partial t} \simeq T^{n-1} + \delta t_n \frac{T^{n-1} - T^{n-2}}{\delta t_{n-1}} = \left(1 + \frac{\delta t_n}{\delta t_{n-1}}\right) T^{n-1} + \frac{\delta t_n}{\delta t_{n-1}} T^{n-2}$$

Starting again from $\mathbf{M}^\theta \cdot \vec{T} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = \vec{0}$, we write

$$\mathbf{M}^\theta \cdot \frac{1}{\delta t_n} \left(\frac{2\delta t_n + \delta t_{n-1}}{\delta t_n + \delta t_{n-1}} \vec{T}^n - \frac{\delta t_n + \delta t_{n-1}}{\delta t_{n-1}} \vec{T}^{n-1} + \frac{\delta t_n^2}{\delta t_{n-1}(\delta t_n + \delta t_{n-1})} \vec{T}^{n-2} \right) + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T}^n = \vec{0}$$

and finally:

$$\left[\frac{2\delta t_n + \delta t_{n-1}}{\delta t_n + \delta t_{n-1}} \mathbf{M}^\theta + \delta t_n (\mathbf{K}_a + \mathbf{K}_d) \right] \cdot \vec{T}^n = \frac{\delta t_n + \delta t_{n-1}}{\delta t_{n-1}} \mathbf{M}^\theta \cdot \vec{T}^{n-1} - \frac{\delta t_n^2}{\delta t_{n-1}(\delta t_n + \delta t_{n-1})} \mathbf{M}^\theta \cdot \vec{T}^{n-2} \quad (6.28)$$

For practical reasons one may wish to bring the advection term to the rhs (i.e. fully implicit) so that the matrix is symmetric. In this case the equation becomes

$$\left[\frac{2\delta t_n + \delta t_{n-1}}{\delta t_n + \delta t_{n-1}} \mathbf{M}^\theta + \delta t_n \mathbf{K}_d \right] \cdot \vec{T}^n = \frac{\delta t_n + \delta t_{n-1}}{\delta t_{n-1}} \mathbf{M}^\theta \cdot \vec{T}^{n-1} - \frac{\delta t_n^2}{\delta t_{n-1}(\delta t_n + \delta t_{n-1})} \mathbf{M}^\theta \cdot \vec{T}^{n-2} - \delta t_n \mathbf{K}_a^\star \cdot \vec{T}^{n,\star}$$

with

$$(\cdot)^\star = \left(1 + \frac{\delta t_n}{\delta t_{n-1}}\right) (\cdot)^{n-1} + \frac{\delta t_n}{\delta t_{n-1}} (\cdot)^{n-2}$$

which denotes the extrapolation of a quantity to time n . Be aware that the \mathbf{K}_a^\star matrix contains the velocity \vec{v}^\star .

Note that if all timesteps are equal, i.e. $\delta t_n = \delta t_{n-1} = \delta t$, Eq. (6.28) becomes:

$$\left[\frac{3}{2} \mathbf{M}^\theta + \delta t (\mathbf{K}_a + \mathbf{K}_d) \right] \cdot \vec{T}^n = \mathbf{M}^\theta \cdot \left(2\vec{T}^{n-1} - \frac{1}{2}\vec{T}^{n-2} \right)$$

or,

$$\left[\mathbf{M}^\theta + \frac{2}{3} \delta t (\mathbf{K}_a + \mathbf{K}_d) \right] \cdot \vec{T}^n = \mathbf{M}^\theta \cdot \left(\frac{4}{3}\vec{T}^{n-1} - \frac{1}{3}\vec{T}^{n-2} \right)$$

When the timestep δt is kept constant (which may be a bad idea with regards to the CFL condition), the backward differencing formula family of implicit methods for the integration of ODEs are simplified. The BDF-1 is simply the backward Euler method as seen above:

$$T^{n+1} - T^n = \delta t \mathcal{F}^{n+1}$$

The BDF-2 is given by

$$T^{n+2} - \frac{4}{3}T^{n+1} + \frac{1}{3}T^n = \frac{2}{3}\delta t \mathcal{F}^{n+2}$$

The BDF-3 is given by

$$T^{n+3} - \frac{18}{11}T^{n+2} + \frac{9}{11}T^{n+1} - \frac{2}{11}T^n = \frac{6}{11}\delta t \mathcal{F}^{n+3}$$

The BDF-4 is given by

$$T^{n+4} - \frac{48}{25}T^{n+3} + \frac{36}{25}T^{n+2} - \frac{16}{25}T^{n+1} + \frac{3}{25}T^n = \frac{12}{25}\delta t \mathcal{F}^{n+4}$$

Each BDF- s method achieves order s .

Anisotropic heat conduction

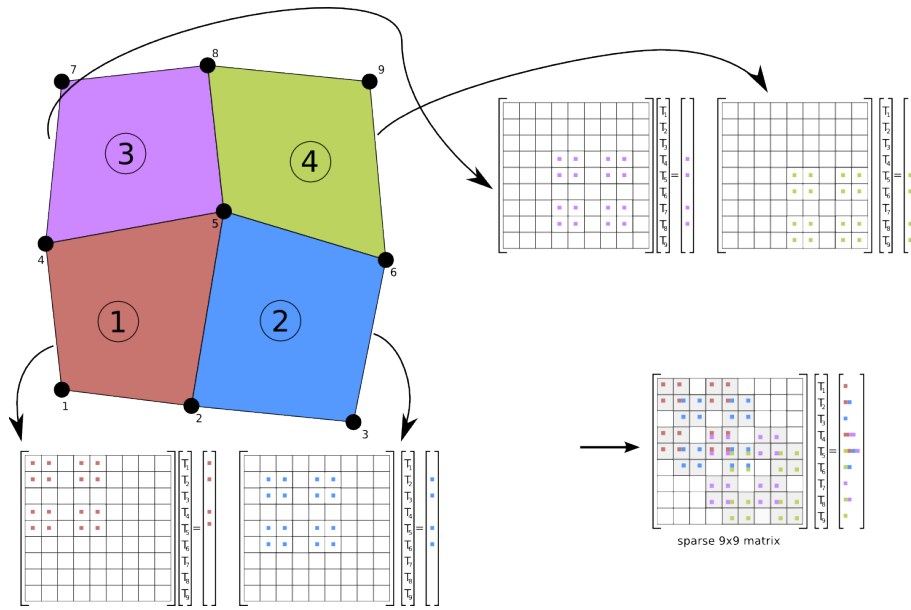
It is most often assumed that the heat conductivity is isotropic so that one speaks of heat conductivity as a scalar k . However many materials are orthotropic and in that case the heat conductivity is a tensor \mathbf{k} which (in 2D) writes (see Reddy [1051, p121]):

$$\mathbf{k} = \begin{pmatrix} k_{xx} & k_{xy} \\ k_{yx} & k_{yy} \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \cdot \begin{pmatrix} k_1 & 0 \\ 0 & k_2 \end{pmatrix} \cdot \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

where k_1 and k_2 are the conductivities in the principal axes system and θ is the local orientation. In that case the diffusion term in the heat transport equation becomes $\vec{\nabla} \cdot (\mathbf{k} \cdot \vec{\nabla} T)$.

About the assembly

Let us consider for simplicity the following grid composed of 9 nodes and 4 Q_1 elements. Each node carries a single degree of freedom.



There are four elements:

- element 1 is composed of nodes $(1, 2, 5, 4) = \vec{T}^{el1}$
- element 2 is composed of nodes $(2, 3, 6, 5) = \vec{T}^{el2}$
- element 3 is composed of nodes $(4, 5, 8, 7) = \vec{T}^{el3}$
- element 4 is composed of nodes $(5, 6, 9, 8) = \vec{T}^{el4}$

For each element one has computed an elemental matrix \mathbf{A}^{el} and a right hand side \vec{b}^{el} .

$$\mathbf{A}^{el1} \cdot \mathbf{T}^{el1} = \mathbf{b}^{el1}$$

$$\mathbf{A}^{el2} \cdot \mathbf{T}^{el2} = \mathbf{b}^{el2}$$

$$\mathbf{A}^{el3} \cdot \mathbf{T}^{el3} = \mathbf{b}^{el3}$$

$$\mathbf{A}^{el4} \cdot \mathbf{T}^{el4} = \mathbf{b}^{el4}$$

As seen in the 1D case, these four linear systems must be assembled in a single large matrix of size 9×9 as shown in the figure above.

6.4 Another approach to solving the advection diffusion

As we have seen above, one usually solves the heat transport equation (i.e. an advection-diffusion equation) in this form (source terms are neglected):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T \quad (6.29)$$

As we have seen in Section 2.6, the diffusion term is actually the divergence of the heat flux $\vec{q} = -k \vec{\nabla} T$. We could then choose to keep the heat flux as an unknown and solve a coupled system of equations instead:

$$\begin{aligned} \rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) &= -\vec{\nabla} \cdot \vec{q} \\ \vec{q} &= -k \vec{\nabla} T \end{aligned}$$

or,

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) + \vec{\nabla} \cdot \vec{q} = 0 \quad (6.30)$$

$$\vec{q} + k \vec{\nabla} T = \vec{0} \quad (6.31)$$

We have seen that the two left hand side terms of the first equation become $\mathbf{M}^\theta \cdot \vec{T}$ and $\mathbf{K}_a \cdot \vec{T}$. Let $\mathcal{N}^\theta(\vec{r})$ be the temperature basis functions so that the temperature inside an element is given by

$$T^h(\vec{r}) = \sum_{i=1}^{m_T} \mathcal{N}_i^\theta(\vec{r}) T_i = \vec{\mathcal{N}}^\theta \cdot \vec{T} \quad (6.32)$$

where \vec{T} is a vector of length m_T . Let $\mathcal{N}^q(\vec{r})$ be the heat flux basis functions, and let us define (in 2D) $\vec{q} = \begin{pmatrix} q_x \\ q_y \end{pmatrix}$ so that

$$q_x^h(\vec{r}) = \sum_{i=1}^{m_q} \mathcal{N}_i^q(\vec{r}) q_{x_i} = \vec{\mathcal{N}}^q \cdot \vec{q}_x \quad (6.33)$$

$$q_y^h(\vec{r}) = \sum_{i=1}^{m_q} \mathcal{N}_i^q(\vec{r}) q_{y_i} = \vec{\mathcal{N}}^q \cdot \vec{q}_y \quad (6.34)$$

where $\vec{\mathcal{N}}^q$, \vec{q}_x and \vec{q}_y are vectors of length m_q . The weak form of the third term of Eq. (6.30) is then

$$\begin{aligned} \int_{\Omega} \mathcal{N}_i^\theta \nabla \cdot \vec{q} d\Omega &= - \int_{\Omega} \mathcal{N}_i^\theta \left(\frac{\partial}{\partial x} q_x + \frac{\partial}{\partial y} q_y \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^\theta \left(\frac{\partial \vec{\mathcal{N}}^q}{\partial x} \cdot \vec{q}_x + \frac{\partial \vec{\mathcal{N}}^q}{\partial y} \cdot \vec{q}_y \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^\theta \frac{\partial \vec{\mathcal{N}}^q}{\partial x} \cdot \vec{q}_x d\Omega + \int_{\Omega} \mathcal{N}_i^\theta \frac{\partial \vec{\mathcal{N}}^q}{\partial y} \cdot \vec{q}_y d\Omega \end{aligned}$$

Writing this last equation for $i = 1, \dots, m_T$ yields

$$\underbrace{\left(\int_{\Omega} \vec{\mathcal{N}}^\theta \frac{\partial \vec{\mathcal{N}}^q}{\partial x} d\Omega \right)}_{\mathbf{H}_x} \cdot \vec{q}_x + \underbrace{\left(\int_{\Omega} \vec{\mathcal{N}}^\theta \frac{\partial \vec{\mathcal{N}}^q}{\partial y} d\Omega \right)}_{\mathbf{H}_y} \cdot \vec{q}_y$$

In the end, we obtain:

$$\mathbf{M}^\theta \cdot \vec{T} + \mathbf{K}_a \cdot \vec{T} + \mathbf{H}_x \cdot \vec{q}_x + \mathbf{H}_y \cdot \vec{q}_y = \vec{0} \quad (6.35)$$

Turning now to Eq. (6.31), its weak form is

$$\int_{\Omega} \mathcal{N}_i^q \left(\vec{q} + k \vec{\nabla} T \right) d\Omega = \vec{0}$$

and we can decompose it in its x and y components:

$$\begin{aligned} 0 &= \int_{\Omega} \mathcal{N}_i^q \left(q_x^h + k \frac{\partial T^h}{\partial x} \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^q \left(\vec{\mathcal{N}}^q \cdot \vec{q}_x + k \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} \cdot \vec{T} \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^q \vec{\mathcal{N}}^q \cdot \vec{q}_x d\Omega + \int_{\Omega} k \mathcal{N}_i^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} \cdot \vec{T} d\Omega \\ 0 &= \int_{\Omega} \mathcal{N}_i^q \left(q_y^h + k \frac{\partial T^h}{\partial y} \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^q \left(\vec{\mathcal{N}}^q \cdot \vec{q}_y + k \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} \cdot \vec{T} \right) d\Omega \\ &= \int_{\Omega} \mathcal{N}_i^q \vec{\mathcal{N}}^q \cdot \vec{q}_y d\Omega + \int_{\Omega} k \mathcal{N}_i^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} \cdot \vec{T} d\Omega \end{aligned} \quad (6.36)$$

Writing these equations for $i = 1, \dots, m_q$ yields:

$$\begin{aligned} 0 &= \int_{\Omega} \vec{\mathcal{N}}^q \vec{\mathcal{N}}^q \cdot \vec{q}_x d\Omega + \int_{\Omega} k \vec{\mathcal{N}}^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} \cdot \vec{T} d\Omega \\ &= \underbrace{\left(\int_{\Omega} \vec{\mathcal{N}}^q \vec{\mathcal{N}}^q d\Omega \right)}_{\mathbf{M}^q} \cdot \vec{q}_x + \underbrace{\left(\int_{\Omega} k \vec{\mathcal{N}}^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial x} d\Omega \right)}_{\mathbf{G}_x} \cdot \vec{T} \end{aligned} \quad (6.37)$$

$$\begin{aligned} 0 &= \int_{\Omega} \vec{\mathcal{N}}^q \vec{\mathcal{N}}^q \cdot \vec{q}_y d\Omega + \int_{\Omega} k \vec{\mathcal{N}}^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} \cdot \vec{T} d\Omega \\ &= \underbrace{\left(\int_{\Omega} \vec{\mathcal{N}}^q \vec{\mathcal{N}}^q d\Omega \right)}_{\mathbf{M}^q} \cdot \vec{q}_y + \underbrace{\left(\int_{\Omega} k \vec{\mathcal{N}}^q \frac{\partial \vec{\mathcal{N}}^\theta}{\partial y} d\Omega \right)}_{\mathbf{G}_y} \cdot \vec{T} \end{aligned} \quad (6.38)$$

Finally Eqs. (6.35, 6.37, 6.38) can be combined and yield the following system (assuming an implicit backward Euler time scheme):

$$\begin{pmatrix} \mathbf{M}^\theta + \mathbf{K}_a \delta t & \mathbf{H}_x \delta t & \mathbf{H}_y \delta t \\ \mathbf{G}_x & \mathbf{M}^q & 0 \\ \mathbf{G}_y & 0 & \mathbf{M}^q \end{pmatrix} \cdot \begin{pmatrix} \vec{T}^{n+1} \\ \vec{q}_x^{n+1} \\ \vec{q}_y^{n+1} \end{pmatrix} = \begin{pmatrix} \mathbf{M}^\theta \cdot \vec{T}^n \\ \vec{0} \\ \vec{0} \end{pmatrix}$$

If we choose $m_q = m_T$ and $\mathcal{N}^q = \mathcal{N}^\theta$ then $\mathbf{M}^\theta = \mathbf{M}^q = \mathbf{M}$ so that

$$\begin{pmatrix} \mathbf{M} + \mathbf{K}_a \delta t & \mathbf{H}_x \delta t & \mathbf{H}_y \delta t \\ \mathbf{G}_x & \mathbf{M} & 0 \\ \mathbf{G}_y & 0 & \mathbf{M} \end{pmatrix} \cdot \begin{pmatrix} \vec{T}^{n+1} \\ \vec{q}_x^{n+1} \\ \vec{q}_y^{n+1} \end{pmatrix} = \begin{pmatrix} \mathbf{M} \cdot \vec{T}^n \\ \vec{0} \\ \vec{0} \end{pmatrix}$$

Also, if k is constant in space then $\mathbf{G}_{x,y} = k\mathbf{H}_{x,y}$. Rather interestingly, one could write Eqs. (6.37,6.38) as

$$\bar{q}_x^{n+1} = -(\mathbf{M}^q)^{-1} \cdot \mathbf{G}_x \cdot \vec{T}^{n+1} \quad (6.39)$$

$$\bar{q}_y^{n+1} = -(\mathbf{M}^q)^{-1} \cdot \mathbf{G}_y \cdot \vec{T}^{n+1} \quad (6.40)$$

and inject it in Eq. (6.35) to yield:

$$\mathbf{M}^\theta \cdot \dot{\vec{T}} + [\mathbf{K}_a - \mathbf{H}_x \cdot (\mathbf{M}^q)^{-1} \cdot \mathbf{G}_x - \mathbf{H}_y \cdot (\mathbf{M}^q)^{-1} \cdot \mathbf{G}_y] \cdot \vec{T}^{n+1} = \vec{0} \quad (6.41)$$

which means that we can directly solve for temperature! Rather interestingly, it is not equivalent to Eq. (6.26). Food for thought ...

We will see that this approach bears a lot of resemblance to the one taken in the context of Discontinuous Galerkin methods.

6.5 The advection-diffusion eq in axisymmetric cylindrical coordinates

hte_axisymm.tex

We start from

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{\mathbf{v}} \cdot \vec{\nabla} T \right) = k \Delta T$$

The temperature gradient in cylindrical coordinates is

$$\vec{\nabla} T = \begin{pmatrix} \partial_r T \\ \frac{1}{r} \partial_\theta T \\ \partial_z T \end{pmatrix}$$

Since $\mathbf{v}_\theta = 0$ and also $\partial_\theta T = 0$ then

$$\vec{\mathbf{v}} \cdot \vec{\nabla} T = \mathbf{v}_r \frac{\partial T}{\partial r} + \mathbf{v}_z \frac{\partial T}{\partial z}$$

and we have the Laplace operator (terms in ∂_θ have been left out):

$$\Delta T = \frac{1}{r} \frac{\partial}{\partial r} \left(r \frac{\partial T}{\partial r} \right) + \frac{\partial^2 T}{\partial z^2}$$

However for the FE formulation we will formulate the equation as

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{\mathbf{v}} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot (k \vec{\nabla} T)$$

After multiplying this equation by a test function and integrating over the domain, the diffusion term is integrated by parts (surface terms are per usual discarded), and we finally obtain

$$\mathbf{M} \cdot \frac{\partial \vec{T}}{\partial t} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = \vec{0}$$

with

$$\mathbf{K}_a = \int \rho C_p \vec{N}^T (\vec{\nu} \cdot \mathbf{B}) dV \quad (6.42)$$

$$\mathbf{K}_d = \int k \mathbf{B}^T \cdot \mathbf{B} dV \quad (6.43)$$

where the matrix \mathbf{B} is identical in this case to the 2D Cartesian one.

It looks like switching from 2D Cartesian to 3D cylindrical axisymmetric does not introduce any change in the formulation. A bit too good to be true ?

 **Relevant Literature:** check section 2.9 of J.N. Reddy and D.K. Gartling. *The Finite Element Method in Heat Transfer and Fluid Dynamics*. CRC Press, 2010. ISBN: 978-1-4200-8598-3.

6.6 The SUPG formulation for the energy equation

supg.tex

As abundantly documented in the literature advection needs to be stabilised as it otherwise showcases non-negligible under- and overshoots. A standard approach is the Streamline Upwind Petrov Galerkin (SUPG) method.

A nice overview of upwind techniques and how SUPG came to be is to be found in Chapter 2 of Donea & Huerta [341] or in Brooks and Hughes (1982) [154]. Hughes *et al.* (1986) [610] present a review of the SUPG method and discuss a additional discontinuity-capturing term. So do Tezduyar and Park [1251] (1986). Note that the idea of upwind goes back to the seventies Heinrich, Huyakorn, Zienkiewicz, and Mitchell [559] (1977)

TODO? It is compared to other methods in the context of mantle convection in Malevsky & Yuen (1991) [826].

Linear elements - artificial diffusion

We have seen in Section 3.9 that the discretised advection-diffusion equation in 1D is given by:

$$\frac{u}{2h} \left[\left(1 - \frac{1}{\text{Pe}}\right) T_{i+1} + \frac{2}{\text{Pe}} T_i - \left(1 + \frac{1}{\text{Pe}}\right) T_{i-1} \right] = f$$

and we show in Stone ?? that the solution is far from accurate for $\text{Pe} > 1$. Let us ask ourselves the following question: could come up with a modified version of the equation above which guarantees an exact solution on a uniform mesh of linear elements?

Essentially, we are looking for A , B and C such that

$$AT_{i+1} + BT_i + CT_{i-1} = f$$

One can show (see Donea & Huerta [341]) that a successful candidate formulation is⁹:

$$\frac{u}{2h} [(1 - \coth(\text{Pe})) T_{i+1} + 2 \coth(\text{Pe}) T_i - (1 + \coth(\text{Pe})) T_{i-1}] = f$$

which can be arranged into the following form:

$$u \frac{T_{i+1} - T_{i-1}}{2h_x} - (\kappa + \tilde{\kappa}) \frac{T_{i-1} - 2T_i + T_{i+1}}{h_x^2} = \frac{f}{\rho_0 C_p} \quad (6.44)$$

⁹ $\coth(x) = (\exp(x) - \exp(-x)) / (\exp(x) + \exp(-x))$

where $\tilde{\kappa}$ is known as artificial (numerical) diffusion (dissipation) given by

$$\tilde{\kappa} = \beta \frac{uh}{2} = \beta \kappa \text{Pe} \quad \text{with} \quad \beta = \coth(\text{Pe}) - \frac{1}{\text{Pe}}$$

where Pe is the (dimensionless) Peclet number defined as:

$$\text{Pe} = \frac{uh}{2\kappa}$$

If $\text{Pe} > 1$ we say the problem is advection-dominated, else if $\text{Pe} < 1$ we say the problem is diffusion-dominated.

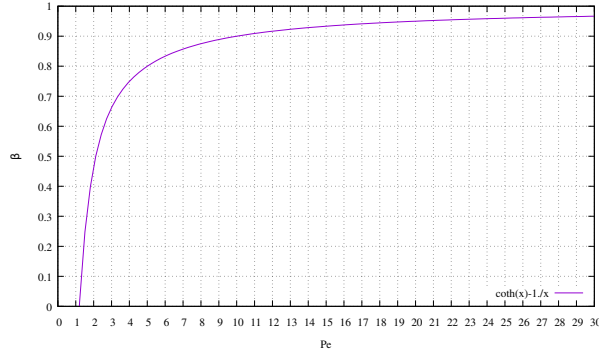
One can also define $\tilde{k} = \rho_0 C_p \tilde{\kappa}$ so that the discretised 1D advection-diffusion becomes:

$$\rho_0 C_p u \frac{T_{i+1} - T_{i-1}}{2h_x} - (k + \tilde{k}) \frac{T_{i-1} - 2T_i + T_{i+1}}{h_x^2} = f \quad (6.45)$$

or, in its continuous formulation:

$$\rho_0 C_p u \frac{dT}{dx} - (k + \tilde{k}) \frac{d^2T}{dx^2} = f$$

The β factor is shown in the following figure as a function of the Peclet number:



Note that the value of β is positive only for $\text{Pe} > 1$.

Linear elements - bubble functions & Petrov-Galerkin formulation

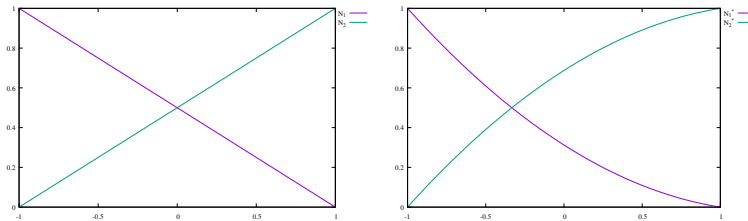
Let us assume that $u > 0$, i.e. advection goes from left to right. Node $i - 1$ is said to be on the upstream side of node i , and node $i + 1$ is on the downstream side of node i . In Petrov FEM, instead of selecting the weight functions to be the same as the standard basis functions, we distort them as shown below.

The distortion is based on so-called bubble functions since they have zero values on the nodes and they are nonzero on elements' interiors. For instance one can take

$$\mathcal{N}_1^*(r) = \frac{1}{2}(1 - r) - \frac{3}{4}\beta(1 - r^2) \quad (6.46)$$

$$\mathcal{N}_2^*(r) = \frac{1}{2}(1 + r) + \frac{3}{4}\beta(1 - r^2) \quad (6.47)$$

as test/weight functions where β is a parameter that controls the amount of upwinding.



Left: standard Q_1 test functions; Right: modified (Q_1 +bubble) functions for $\beta = 0.25$

$$\begin{aligned}
\mathbf{K}_a^e &= \int_{x_k}^{x_{k+1}} (\vec{\mathcal{N}}^\star)^T \rho C_p \vec{\mathbf{v}} \cdot \mathbf{B} \, dx \\
&= \rho C_p u \begin{pmatrix} -1/2 & 1/2 \\ -1/2 & 1/2 \end{pmatrix} + \frac{h}{2} \int_{-1}^{+1} \rho C_p u \begin{pmatrix} -\frac{3}{4}\beta(1-r^2) \\ +\frac{3}{4}\beta(1-r^2) \end{pmatrix} \begin{pmatrix} \frac{d\mathcal{N}_1}{dr} & \frac{d\mathcal{N}_2}{dr} \end{pmatrix} dr \quad (6.48)
\end{aligned}$$

$$= \frac{\rho C_p u}{2} \begin{pmatrix} -1 & 1 \\ -1 & 1 \end{pmatrix} + \frac{u}{2} \begin{pmatrix} \beta & -\beta \\ -\beta & \beta \end{pmatrix} \quad (6.49)$$

VERIFY

We see that this additional matrix is akin to \mathbf{K}_d^e so that

$$\mathbf{K}^e = \mathbf{K}_a^e + \mathbf{K}_d^e = \frac{\rho C_p u}{2} \begin{pmatrix} -1 & 1 \\ -1 & 1 \end{pmatrix} + \left(\frac{k}{h} + \frac{u\beta\rho C_p}{2} \right) \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

VERIFY!!

The SUPG method

We start from the 'bare-bones' heat transport equation (source terms are omitted):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T \quad (6.50)$$

which we once again write

$$\rho C_p \left(\dot{T} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T \quad (6.51)$$

In what follows we assume that the velocity field \vec{v} is known so that temperature is the only unknown. Let N_i^θ be the temperature basis function at node i so that the temperature inside an element is given by¹⁰:

$$T^h(\vec{r}) = \sum_{i=1}^{m_T} \mathcal{N}_i^\theta(\vec{r}) T_i = \vec{\mathcal{N}}^\theta \cdot \vec{T} \quad (6.52)$$

where \vec{T} and $\vec{\mathcal{N}}^\theta$ are vectors of length m_T . Also:

$$\vec{\nabla} \vec{\mathcal{N}}^\theta = \begin{pmatrix} \partial_x \mathcal{N}_1^\theta & \partial_x \mathcal{N}_2^\theta & \dots & \partial_x \mathcal{N}_{m_T}^\theta \\ \partial_y \mathcal{N}_1^\theta & \partial_y \mathcal{N}_2^\theta & \dots & \partial_y \mathcal{N}_{m_T}^\theta \end{pmatrix}$$

The weak form is then

$$\int_{\Omega_e} \mathcal{N}_i^\theta \left[\rho C_p \left(\dot{T} + \vec{v} \cdot \vec{\nabla} T \right) \right] d\Omega = \int_{\Omega_e} \mathcal{N}_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T d\Omega \quad (6.53)$$

or,

$$\underbrace{\int_{\Omega_e} \mathcal{N}_i^\theta \rho C_p \dot{T}^h d\Omega}_I + \underbrace{\int_{\Omega_e} \mathcal{N}_i^\theta \rho C_p \vec{v} \cdot \vec{\nabla} T^h d\Omega}_{II} = \underbrace{\int_{\Omega_e} \mathcal{N}_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T^h d\Omega}_{III} \quad i = 1, m_T$$

The streamline upwind Petrov-Galerkin (SUPG) method adds the following stabilisation term to the left hand side of this equation (see Section 2.4 in Donea & Huerta [341]):

$$\int_{\Omega_e} (\vec{v} \cdot \vec{\nabla} \mathcal{N}_i^\theta) \tau \mathcal{R}(\vec{v}) d\Omega$$

where \mathcal{R} is the residual defined as

$$\mathcal{R}(T^h) = \rho C_p (\dot{T}^h + \vec{v} \cdot \vec{\nabla} T^h) - \vec{\nabla} \cdot k \vec{\nabla} T^h$$

and τ is a stabilisation parameter (see Section 6.6).

We have already worked out in Section 6.3 the final forms of I , II and III , so we here focus on the additional SUPG terms.

We then have three terms to deal with:

$$IV = \int_{\Omega_e} \tau (\vec{v} \cdot \nabla \mathcal{N}_i^\theta) (\rho C_p \dot{T}^h) d\Omega \quad (6.54)$$

$$V = \int_{\Omega_e} \tau (\vec{v} \cdot \nabla \mathcal{N}_i^\theta) (\rho C_p \vec{v} \cdot \vec{\nabla} T^h) d\Omega \quad (6.55)$$

$$VI = \int_{\Omega_e} \tau (\vec{v} \cdot \nabla \mathcal{N}_i^\theta) (-\vec{\nabla} \cdot k \vec{\nabla} T^h) d\Omega \quad (6.56)$$

¹⁰the θ superscript has been chosen to denote temperature so as to avoid confusion with the transpose operator

We can compute the quantity $I + IV$:

$$\begin{aligned} I + IV &= \int_{\Omega_e} \mathcal{N}_i^\theta \rho C_p \dot{T}^h d\Omega + \int_{\Omega_e} (\vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta) \tau (\rho C_p \dot{T}^h) d\Omega \\ &= \int_{\Omega_e} (\mathcal{N}_i^\theta + \tau \vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta) (\rho C_p \dot{T}^h) d\Omega \end{aligned} \quad (6.57)$$

and also $II + V$:

$$II + V = \int_{\Omega_e} \mathcal{N}_i^\theta \rho C_p \vec{\nu} \cdot \vec{\nabla} T^h d\Omega + \int_{\Omega_e} (\vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta) \tau (\rho C_p \vec{\nu} \cdot \vec{\nabla} T^h) d\Omega \quad (6.58)$$

$$= \int_{\Omega_e} (\mathcal{N}_i^\theta + \tau \vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta) (\rho C_p \vec{\nu} \cdot \vec{\nabla} T^h) d\Omega \quad (6.59)$$

Remark. *Because of the integration by parts which will be applied to III, the terms III and VI cannot be summed together.*

Remark. *If the equation is a pure advection equation, then $k = 0$ so $III = VI = 0$.*

We see that both $I + IV$ and $II + V$ contain the term $\mathcal{N}_i^\theta + \tau \vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta$ which we can interpret as a 'modified' basis function:

$$\boxed{\underline{\mathcal{N}}_i^\theta = \mathcal{N}_i^\theta + \tau \vec{\nu} \cdot \vec{\nabla} \mathcal{N}_i^\theta}$$

This yields the following modified elemental matrices

$$\mathbf{M}^\theta \rightarrow \underline{\mathbf{M}}^\theta = \int \rho C_p \underline{\mathcal{N}}^\theta \vec{\mathcal{N}}^\theta d\Omega$$

$$\mathbf{K}_a \rightarrow \underline{\mathbf{K}}_a = \int \rho C_p \underline{\mathcal{N}}^\theta (\vec{\nu} \cdot \vec{\nabla} \vec{\mathcal{N}}^\theta) d\Omega$$

Remark. *The modified mass matrix is not symmetrical anymore.*

Under the assumption that k is constant within the element, we have

$$VI = \int_{\Omega_e} \tau (\vec{\nu} \cdot \nabla \mathcal{N}_i^\theta) (-k \Delta T^h) d\Omega$$

with

$$\Delta T^h = \Delta \sum_{i=1}^{m_T} \mathcal{N}_i^\theta(\vec{r}) T_i = \sum_{i=1}^{m_T} \Delta \mathcal{N}_i^\theta(\vec{r}) T_i = \Delta \vec{\mathcal{N}}^\theta \cdot \vec{T}$$

Remark. *If the basis functions are first-order ones, then $VI = 0$ since $\Delta \vec{\mathcal{N}}^\theta = \vec{0}$.*

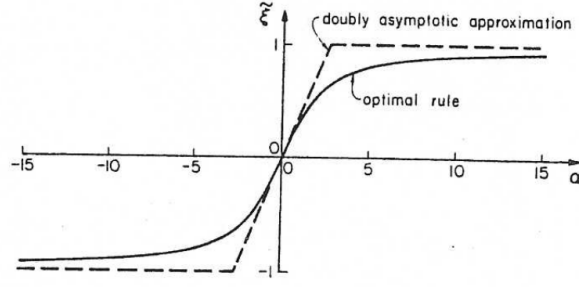
The SUPG-stabilised pure advection equation is extensively tested in the Stone ??.

About the choice of the parameter τ

- The approach used in step-63 and subsequently in ASPECT is the one in John & Knobloch [654, 715] which is also to be found in early 80's papers [154, 605]:

$$\tau_1 = \frac{h}{2|\vec{\nu}|p} \left(\coth(\text{Pe}) - \frac{1}{\text{Pe}} \right)$$

where τ is computed on each cell, h being its diameter in the direction of $\vec{\nu}$, and p is the order of the approximation (i.e. the maximum degree of polynomials). Note that Hughes & Brooks [605] replace the costly 'coth' term by an asymptotic curve:



Taken from [605]. Here α is the Peclet number and the vertical axis is the term $\coth(\text{Pe}) - 1/\text{Pe}$.

- Codina (2000) (see Eq.39 of [270]) defines τ as follows:

$$\tau_2 = \left(\frac{2|\vec{v}|}{h} + \frac{4\kappa}{h^2} + \sigma \right)^{-1}$$

where σ is a reaction term (which we neglect in what follows). This equation can be re-written:

$$\tau_2 = \frac{h}{2|\vec{v}|} \left(1 + \frac{1}{\text{Pe}} \right)^{-1}$$

- Shakib *et al.* (1991) (see Eq. 3.59 of [1152]) propose another formula:

$$\tau_3 = \frac{h}{2|\vec{v}|} \left(1 + \frac{9}{\text{Pe}^2} \right)^{-1}$$

- Braun [135] uses

$$\tau_4 = \frac{h}{|\vec{v}|\sqrt{15}} = \frac{h}{2|\vec{v}|} \frac{2}{\sqrt{15}}$$

This formula is independent of Pe and is also in Bochev *et al.* (2004) [104]. It is also used in Hughes & Brooks (1982) [605] for pure advection problems in 1D and this value is attributed to Raymond & Garder (1976) [1049].

- Following [1250] (see also Appendix A of Thieulot (2011) [1258]):

$$\tau_5 = \left(\frac{2|\vec{v}|}{h} + \frac{1}{\theta\delta t} + \frac{\kappa}{h^2} \right)^{-1}$$

Crank-Nicolson: $\theta = 1/2$, CFL condition yields $\delta t = C \frac{h}{|\vec{v}|}$ so

$$\tau_5 = \left(\frac{2|\vec{v}|}{h} + \frac{2}{\delta t} + \frac{\kappa}{h^2} \right)^{-1} = \left(\frac{2|\vec{v}|}{h} + \frac{2v}{Ch} + \frac{\kappa}{h^2} \right)^{-1} = \frac{h}{2|\vec{v}|} \left(1 + \frac{1}{C} + \frac{1}{4\text{Pe}} \right)^{-1}$$

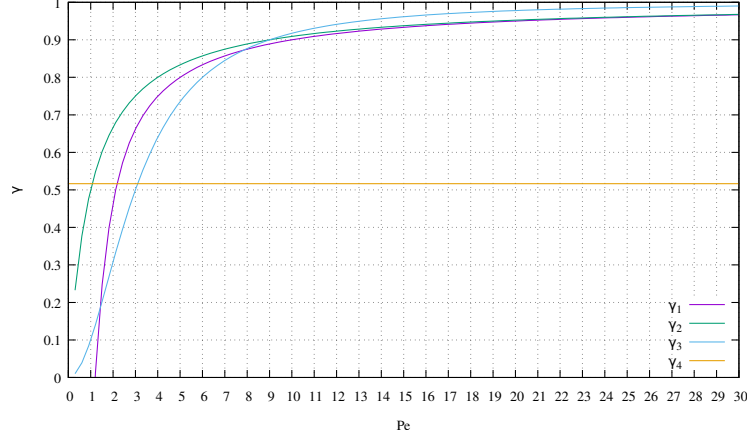
Note that the velocity in the CFL criterion is the maximum velocity in the domain, while in all other expressions above for τ_i it is the velocity in the cell/element. By doing so, the relationship for τ_5 is actually only valid for the cell with the highest velocity. Also, this has been used in the context of SUPG methods for the Navier-Stokes equations, not advection-diffusion equations! This approach is then not considered further.

- Franca *et al.* (2004) [413] use the following formula which originates from Franca *et al.* (1992) [412]: if $0 \leq \text{Pe} < 1$ then $\tau = \frac{h}{2|\vec{v}|} \text{Pe}$ and if $\text{Pe} \geq 1$ then $\tau = \frac{h}{2|\vec{v}|}$. Note however that their definition of the Peclet number includes a scalar parameter m which is the minimum between $1/3$ and $2C_k$ where the calculation of this parameter is discussed in Remarks 4 and 5 of [412]. The authors conclude that for linear quadrilaterals the value $1/3$ should be used while for biquadratic elements $1/12$ should be preferred. The same approach is found in Brezzi *et al.* (1992) [150].

- Knobloch [715] discusses other choices of τ .

Quoting Donea & Huerta: "It is obvious that τ must vanish when the mesh is refined (no stabilisation is necessary for a fine enough mesh)" and "Numerical experiments seem to indicate that for finite elements of order p the value of the stabilisation parameter should be approximately τ/p ."

Let us define $\gamma = \frac{\tau}{h/2|\vec{v}|}$ (dimensionless quantity) and plot this quantity against the Peclet number:



In the case when the equation to be stabilised is a pure advection equation, then $Pe \rightarrow \infty$ so

$$\gamma_1 \rightarrow 1 \quad \gamma_2 \rightarrow 1 \quad \gamma_3 \rightarrow 1 \quad \gamma_4 \simeq 0.52 \quad \gamma_5 \rightarrow (1 + C^{-1})^{-1}$$

Some remarks about Appendix A of Thieulot (2011) [1258]

In the DOUAR paper [136] or the FANTOM paper [1258], The advection matrix is simply modified and computed as follows:

$$(\mathbf{K}_a^e)_{SUPG} = \int_{x_k}^{x_{k+1}} (\vec{\mathcal{N}}^\star)^T \rho C_p \vec{v} \cdot \mathbf{B} dx \quad \text{with} \quad \vec{\mathcal{N}}^\star = \vec{\mathcal{N}} + \tau \vec{v} \cdot \mathbf{B}$$

Note that we can also write

$$(\mathbf{K}_a^e)_{SUPG} = \int_{x_k}^{x_{k+1}} \vec{\mathcal{N}}^T \rho C_p \vec{v} \cdot \mathbf{B} dx + \int_{x_k}^{x_{k+1}} \tau (\vec{v} \cdot \mathbf{B})^T \rho C_p (\vec{v} \cdot \mathbf{B}) dx$$

and we see that the SUPG method introduces an additional term that is akin to a diffusion term in the direction of the flow. This can be seen by looking at the advection matrix a regular grid of 1D elements of size h :

$$(\mathbf{K}_a^e)_{SUPG} = \mathbf{K}_a^e + \rho C_p \frac{\tau u^2}{h} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

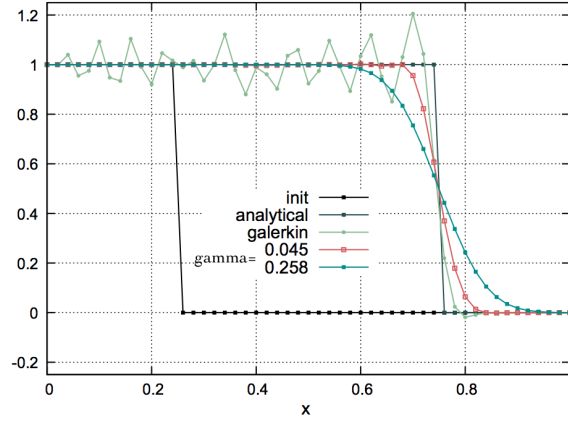
The additional matrix has the same structure as the 1D diffusion matrix in 6.1.

The parameter τ is chosen as follows [1250]:

$$\tau = \left(\frac{1}{\tau_1} + \frac{1}{\tau_2} + \frac{1}{\tau_3} \right)^{-1} \quad \tau_1 = \frac{h}{2|\vec{v}|}, \quad \tau_2 = \theta \delta t, \quad \tau_3 = \frac{h^2 \rho C_p}{k} \quad (6.60)$$

where h is a measure of the element size and θ is related to the time discretisation scheme ($\theta = 1/2$ corresponds to a mid-point implicit scheme), and we can define $\gamma = \tau|\vec{v}|/h$ (see Appendix A of [1258]).

A typical test case for testing an advection scheme is the step advection benchmark (see for instance Donea & Huerta (2003) [341]). At $t = 0$, a field $T(x)$ is prescribed in a 1D domain of unit length. For $x \leq 1/4$ we have $T(x) = 1$ and $T(x) = 0$ everywhere else as shown on the following figure:



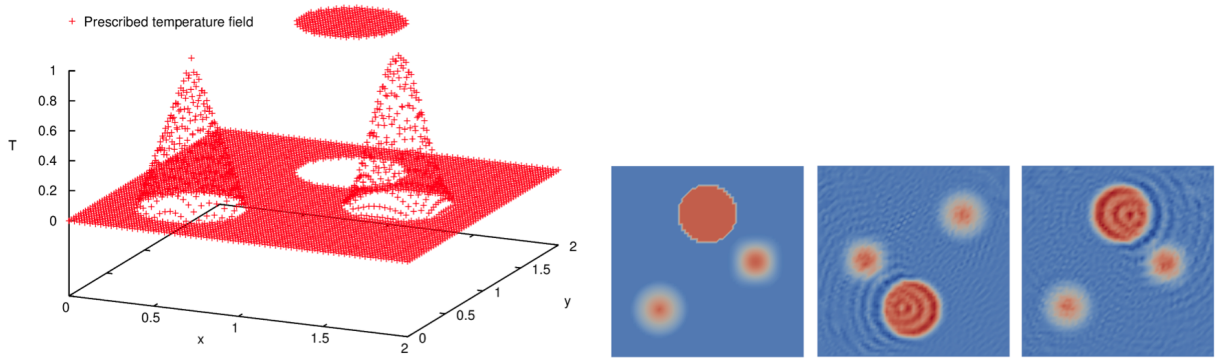
Taken and modified from Thieulot (2011) [1258]

The prescribed velocity is $\mathbf{v} = 1$, 50 elements are used ($h = 0.02$) and 250 time steps are carried out with $\delta t = 0.1h/\mathbf{v} = 0.002$ (CFL number of 0.1). Then it follows that $\tau_3 = \infty$ (no diffusion, i.e. $k = 0$) and

$$\tau = \left(\frac{2}{0.02} + \frac{2}{0.002} \right)^{-1} = \left(\frac{1}{0.01} + \frac{1}{0.001} \right)^{-1} = \left(\frac{1}{0.01} \left(1 + \frac{1}{10} \right) \right)^{-1} = 0.01 \left(\frac{11}{10} \right)^{-1} \simeq 0.009091$$

which yields $\gamma = 0.00909/0.02 = 0.4545\dots$ Using this value leads to a desired removal of the oscillations through a small amount of numerical diffusion. Braun [135] argues for a constant $\gamma = 1/\sqrt{15} = 0.258$ (citing Hughes & Brooks (1982) [605]), which effect is also shown in the figure above. This value is arguably too large and introduces an undesirable diffusion. Note that this same value is to be found in Bochev *et al.* (2004) [104]. The authors then state that "for 1D pure advection problems, this choice maximizes the phase accuracy in the semidiscrete equation" and cite Raymond & Gardner (1976) [1049] as source.

Another classic example of advection testing is a 2D problem where (for example) a cylinder, a Gaussian and a cone are prescribed and advected with a velocity field (see for instance [341]).



After a 2π rotation and in the absence of stabilisation we see that the temperature field showcases clearly visible ripples.

Remark. Note that ASPECT originally did not rely on the SUPG formulation to stabilise the advection(-diffusion) equations[732]. It instead relied on the Entropy Viscosity formulation [502, 501]. It is only during the 6th Hackathon in May 2019 that the SUPG was introduced on the code. Note that the ASPECT implementation is based on the deal.II step 63¹¹.

¹¹https://www.dealii.org/developer/doxygen/deal.II/step_63.html

Chapter 7

Solving the Stokes equations with the FEM

chapter_fem2.tex

7.1 A quick tour of similar literature

- *Treatise on Geophysics*, Volume 7, Edited by D. Bercovici and G. Schubert: "Numerical Methods for Mantle Convection", by S.J. Zhong, D.A. Yuen, L.N. Moresi and M.G. Knepley. Note that it is a revision of the previous edition chapter by S.J. Zhong, D.A. Yuen and L.N. Moresi, Volume 7, pp. 227-252, 2007.
- *Computational Science I*, Lecture Notes for CAAM 519, M.G. Knepley, 2017. <https://cse.buffalo.edu/~knepley/classes/caam519/>
- *Numerical Modeling of Earth Systems - An introduction to computational methods with focus on solid Earth applications of continuum mechanics*, Th.W. Becker and B.J.P. Kaus, 2018. <http://www-udc.ig.utexas.edu/external/becker/Geodynamics557.pdf>
- *Myths and Methods in Modeling*, M. Spiegelman, 2000. https://earth.usc.edu/~becker/teaching/557/reading/spiegelman_mmm.pdf

In the case of an incompressible flow, we have seen that the continuity (mass conservation) equation takes the simple form $\vec{\nabla} \cdot \vec{v} = 0$. In other words flow takes place under the constraint that the divergence of its velocity field is exactly zero everywhere (solenoidal constraint), i.e. it is divergence free.

We see that the pressure in the momentum equation is then a degree of freedom which is needed to satisfy the incompressibility constraint (and it is not related to any constitutive equation) (see for example Donea and Huerta [341]). In other words the pressure is acting as a Lagrange multiplier of the incompressibility constraint.

Various approaches have been proposed in the literature to deal with the incompressibility constraint but we will only focus on the penalty method (section 7.4) and the so-called mixed finite element method 7.5.

7.2 Strong and weak forms

strongweak.tex

As we have seen in Section 6.1 the strong form consists of the governing equation and the boundary conditions, i.e. the mass, momentum and energy conservation equations supplemented with Dirichlet and/or Neumann boundary conditions on (parts of) the boundary. Ultimately we have two main unknowns that we wish to solve for: velocity (a vector) and pressure (a scalar).

To develop the finite element formulation, the partial differential equations must be restated in an integral form called the weak form. In essence the PDEs are first multiplied by an arbitrary function and integrated over the domain.

7.3 Which velocity-pressure pair for Stokes?

The success of a mixed finite element formulation crucially depends on a proper choice of the local interpolations of the velocity and the pressure.

7.3.1 The compatibility condition (or LBB condition, or inf-sup condition)

lbb.tex

WARNING: I am not comfortable (yet) writing about this topic. What follows is a rough attempt at making sense of it.

The Ladyžhenskaya-Babuška-Brezzi (LBB¹) condition is a sufficient condition for a saddle point problem to have a unique solution. For saddle point problems coming from the Stokes equations, many discretizations (i.e. choices for the velocity and pressure polynomial spaces) are unstable, giving rise to artifacts such as spurious oscillations. The LBB condition gives criteria for when a discretization of a saddle point problem is stable. It also assures convergence at the optimal rate.

Bochev & Gunzburger [103] state: “The terminology ‘LBB’ originates from the facts that this condition was first explicitly discussed in the finite element setting for saddle point problems by Brezzi² [149] and that it is a special case of the general weak-coercivity condition first discussed for finite element methods by Ivo Babuška³ [37] and that, in the continuous setting of the Stokes equation, this condition was first proved to hold by Olga Ladyzhenskaya⁴; see [740].”

Unfortunately, to quote Donea & Huerta [341]: “In the finite element context, it is by no means easy to prove whether or not a given velocity-pressure pair satisfies the LBB compatibility condition.” Elman *et al.* state: “[...] Choosing spaces for which the discrete inf-sup condition holds and is a delicate matter, and seemingly natural choices of velocity and pressure approximation do not work. [...] In general, care must be taken to make the velocity space rich enough compared to the pressure space.” By rich enough the authors essentially mean that the order of the polynomials used to represent velocity must be higher than the one used for pressure.

The LBB condition, or inf-sup condition can be proven in different ways, and standard techniques have been designed as listed in Boffi *et al.* (2008) [108].

Elman *et al.* [371] state that “The inf-sup condition is a sufficient condition for the pressure to be unique up to constant in the case of an enclosed flow.” This can also be proven for other boundary conditions. This approach, based on the macro-element technique [1207] is explored in Appendix K.

It can be shown that, provided the kernel (null space) of matrix \mathbb{G} is zero, the Stokes matrix is non-singular, that is $\vec{\mathcal{V}}$ and $\vec{\mathcal{P}}$ are uniquely defined, and the Schur complement matrix \mathbb{S} is positive definite. Simply put, taking $\vec{\mathcal{V}} = \vec{0}$ in the discretised Stokes system without body forces yields $\mathbb{G} \cdot \vec{\mathcal{P}} = \vec{0}$ and implies that any pressure solution is only unique up to the null space of the matrix \mathbb{G} .

We know that the Schur complement matrix \mathbb{S} is positive definite if and only if all of its eigenvalues are positive. One could then (numerically) compute the eigenvalues of \mathbb{S} and check that these are indeed strictly positive to show that \mathbb{S} is positive definite but that would prove very costly.

Another way is to see that \mathbb{S} is positive definite only if $\ker(\mathbb{G}) = \{0\}$. Again to quote Donea & Huerta [341]: “If this is the case, the partitioned Stokes matrix is non-singular and delivers uniquely defined velocity and pressure fields. If this is not the case, a stable and convergent velocity field might be obtained, but the pressure field is likely to present spurious and oscillatory results.” Note

¹https://en.wikipedia.org/wiki/Ladyzhenskaya-Babuska-Brezzi_condition

²https://en.wikipedia.org/wiki/Franco_Brezzi

³https://en.wikipedia.org/wiki/Ivo_Babuska

⁴https://en.wikipedia.org/wiki/Olga_Ladyzhenskaya

that in the case of the $\mathbf{Q}_1 \times P_0$ element it has been shown that the multiple families of checkboard pressure modes actually lie in the kernel of \mathbb{G} . [1108, 1109]

We can look at this in a different manner, as explained in Elman, Silvester, and Wathen [371]: the unique solvability of the matrix system

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (7.1)$$

is determined by looking at the homogeneous system

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{0} \\ \vec{0} \end{pmatrix} \quad (7.2)$$

or,

$$\begin{aligned} \mathbb{K} \cdot \vec{\mathcal{V}} + \mathbb{G} \cdot \vec{\mathcal{P}} &= \vec{0} \\ \mathbb{G}^T \cdot \vec{\mathcal{V}} &= \vec{0} \end{aligned} \quad (7.3)$$

To start, premultiply the first equation by $\vec{\mathcal{V}}^T$ and the second by $\vec{\mathcal{P}}^T$. The second yields $\vec{\mathcal{P}}^T \cdot \mathbb{G}^T \cdot \vec{\mathcal{V}} = (\vec{\mathcal{V}}^T \cdot \mathbb{G} \cdot \vec{\mathcal{P}})^T = \vec{0}$ which is present in the first equation so that it simplifies to $\vec{\mathcal{V}}^T \cdot \mathbb{K} \cdot \vec{\mathcal{V}} = \vec{0}$. Since \mathbb{K} is positive definite, it follows that $\vec{\mathcal{V}} = \vec{0}$, implying unique solvability with respect to the velocity.

On the other hand, unique solvability with respect to the pressure is problematic. Substituting $\vec{\mathcal{V}} = \vec{0}$ in the system above gives $\mathbb{G} \cdot \vec{\mathcal{P}} = \vec{0}$, and implies that any pressure solution is only unique up to the nullspace of the matrix \mathbb{G} . The bottom line is that if Eq. (7.1) is to properly represent a continuous Stokes problem, then the mixed approximation spaces need to be chosen carefully. Specifically, we have to ensure that $\text{null}(\mathbb{G}) = \{1\}$ in the case of enclosed flow, and that $\text{null}(\mathbb{G}) = \{0\}$, otherwise.

Gresho and Sani [488] state: “LBB stable elements assure the existence of a unique solution to Stokes flow and assure convergence at optimal rate. [...] LBB-unstable elements may not converge, and if they do, they may not do so at the optimal rate.”

7.3.2 Families

The family of **Taylor-Hood** finite element spaces on triangular/tetrahedral grids is given by $\mathbf{P}_k \times P_{k-1}$ with $k \geq 2$, and on quadrilateral/hexahedral grids by $\mathbf{Q}_k \times Q_{k-1}$ with $k \geq 2$. This means that the pressure is then approximated by continuous functions.

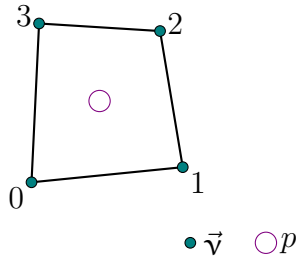
These finite elements are very popular, in particular the pairs for $k = 2$, i.e. $\mathbf{Q}_2 \times Q_1$ and $\mathbf{P}_2 \times P_1$. The reason why $k \geq 2$ comes from the fact that the $\mathbf{Q}_1 \times Q_0$ (often referred to as $\mathbf{Q}_1 \times P_0$) and $\mathbf{P}_1 \times P_0$ are not stable elements (they are not inf-sup stable), as shown in John [650, p64] and [650, p67].

Remark. Note that a similar element to $\mathbf{Q}_2 \times Q_1$ has been proposed and used successfully in Taylor and Hood [1240] (1973) and Hood and Taylor [589] (1974): it is denoted by $\mathbf{Q}_2^{(8)} \times Q_1$ since the center node (x^2y^2) and its associated degrees of freedom have been removed. It has also been proved to be LBB stable. These are also called **Serendipity** elements.

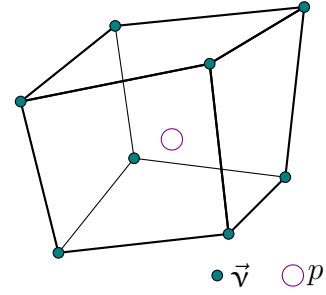
7.3.3 The bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$)

pair_q1p0.tex

4 vel. nodes, 1 press. node



8 vel. nodes, 1 press. node




However simple it may look, the element is one of the hardest elements to analyze and many questions are still open about its properties. The element does not satisfy the inf-sup condition [604, p211]. In Gresho & Sani [488] it is labeled as follows: “slightly unstable but highly usable”.

The $\mathbf{Q}_1 \times P_0$ mixed approximation is the lowest order conforming approximation method defined on a rectangular grid. It also happens to be the most famous example of an unstable mixed approximation method [371, p235]. Boland and Nicolaides [111] (1984) and Boland and Nicolaides [112] (1985) show that it is not stable.

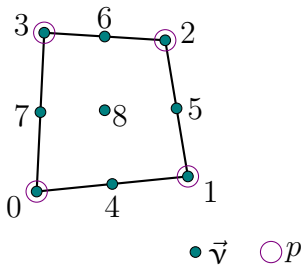
This element is discussed in Fortin (1981) [401], Fortin & Fortin (1985) [403] and in Pitkäranta & Saarinen (1985) [1001] in the context of multigrid use.

This element is plagued by so-called pressure checkerboard modes which have been thoroughly analysed, see for example Griffiths and Silvester [495] (1994), Chen, Pan, and Chang [221] (1995), Sani, Gresho, Lee, and Griffiths [1108] and Sani, Gresho, Lee, Griffiths, and Engelman [1109] (1981). These can be filtered out, see for example Chen, Pan, and Chang [221] (1995) or Lee, Gresho, and Sani [759] (1997), and explained in Section 9.7.

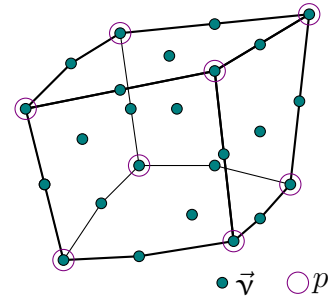
 **Relevant Literature:** Fortin and Boivin [402] (1990), [491] (1985), LeTallec and Ruas [774] (1986), Oden and Jacquotte [950] (1984).

7.3.4 The bi/tri-quadratic velocity - bi/tri-linear pressure element ($Q_2 \times Q_1$)

9 vel. nodes, 4 press. nodes



27 vel. nodes, 8 press. nodes



It belongs to the Taylor-Hood family of elements and satisfies the inf-sup (LBB) condition [604, p215]. Gresho & Sani [488, p554] write that in their opinion $\text{div}(\vec{v}) = 0$ is not strong enough. This element, implemented in penalised form, is discussed in Bercovier & Engelman (1979) [79] and the follow-up paper [80].

It is the default of the ASPECT code (see Appendix ??). It is implemented in [STONE](#) 18,21,48,91,120,...

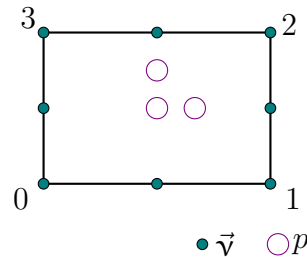
7.3.5 The bi/tri-quadratic velocity - discontinuous linear pressure element ($Q_2 \times P_{-1}$)

pair_q2pm1.tex

According to Boffi, Brezzi, and Fortin [108] “This element was apparently discovered around a blackboard at the Banff Conference on Finite Elements in Flow Problems (1979)”.

(tikz_p2pm1.tex)

9 vel. nodes, 3 press. nodes



This element is crowned “probably the most accurate 2D element” in Gresho and Sani [488].

It is characterised by piecewise bi/triquadratic velocities, and piecewise linear discontinuous polynomial pressure. The element satisfies the inf-sup condition, see p. 211 of Hughes [604], or p. 138 of Elman, Silvester, and Wathen [371]. It is used in Vosse, Steenhoven, Segal, and Janssen [1330] (1989) for steady laminar flow in a curved tube.

When using this element one must be aware of the fact that there are two possible choices for the definitions of the pressure space (mapped and un-mapped) as explained in Boffi and Gastaldi [107] (2002). See [STONE 76,120](#) for their implementation. Boffi, Brezzi, and Fortin [108] state: “On a general quadrilateral mesh, the [pressure] space can be defined in two different ways: either [it] consists of (discontinuous) piecewise linear functions, or it is built by considering three linear shape functions on the reference unit square and mapping them to the general elements like it is usually done for continuous finite elements. [...] We shall refer to the first possibility as unmapped pressure approach and to the second one as mapped pressure approach.” Furthermore they state “So far, we have shown that either the unmapped and the mapped pressure approach gives rise to a stable $Q_2 \times P_{-1}$ scheme. However, as a consequence of the results proved in Arnold, Boffi, and Falk [27] (2002), we have that the mapped pressure approach cannot achieve optimal approximation order. Namely, the unmapped pressure space provides a second-order convergence in L_2 , while the mapped one achieves only $\mathcal{O}(h)$ in the same norm.” See also discussion about mapped/unmapped in Boffi, Brezzi, and Fortin [109].

This element is mentioned in Kaus [679] (2010) and Pelletier, Fortin, and Camarero [984] (1989) and it is used in Frehner [417] (2014) to study 3D fold growth rates (see online supplementary material) and in Schmalholz [1118] (2008).

Note that the serendipity version of this pair, i.e. $Q_2^{(20)} \times P_{-1}$ is also LBB stable as shown in p180 of Reddy [1051].

7.3.6 The biquadratic velocity - discontinuous bilinear pressure element ($Q_2 \times Q_{-1}$)

This element is shown in Table 3.13-2 of Gresho & Sani’s book [488], and discussed in Section 3.13.6b of the book too. It is *not* LBB stable and has one checkerboard pressure mode.

It is used (alongside many other element pairs) in Christon, Gresho, and Sutton [256] (2002) in the context of a flow benchmark in a 2D box. The authors conclude that “[...] the Q2-Q-1 element

fared slightly better than the Q2-P-1 . Most surprising, though, were the good results obtained with the 'old' Taylor–Hood element, Q2-Q1 .”

It is also used in Gresho and Sutton [489] (2002) on a similar benchmark setup (8:1 thermal cavity problem) along with $\mathbf{Q}_1 \times P_0$, $\mathbf{Q}_2 \times P_{-1}$ and $\mathbf{Q}_2 \times Q_1$. The authors state that Q2Q-1 has div- stability problems but “produces excellent results and is still useful in general.” They also state “If the pesky-mode instability could be eciently dealt with, then the Q2xQ-1 element should be employed over the Q2xP-1 -especially in 3D (we believe).” Authors mention that it was also used in Vahl Davis and Jones [1300] and that it “performed EXTREMELY WELL.”

7.3.7 The stabilised bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$ -stab)

pair-q1p0stab.tex

Much has been written about the $Q_1 \times P_0$ element and the fact that it is not LBB-stable and that the pressure field contains a chequerboard mode that needs to be filtered out. It was the principal element used in computational geodynamics in codes such as Sopale, Citcom, Phantom, ... before being superseded by LBB-stable elements such as $Q_2 \times Q_1$ or $Q_2 \times P_{-1}$ [1260].

Many techniques have been proposed to stabilise this element but I here focus on those which keep the number of degrees of freedom unchanged, i.e. a matrix \mathbb{C} is added to the Stokes matrix:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix}$$

More specifically I will focus on the pressure jump methods.

Note that in 3D the physical dimension of the \mathbb{C} matrix is that of h^{dim}/η (i.e. $M^{-1}L^4T$) where h is the element size and η a viscosity. The Schur complement $\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$ has obviously the same dimensions with $[\mathbb{G}] = L^2$ and $[\mathbb{K}] = MT^{-1}$.

As explained in Silvester & Kechkar [1167]: “The system [without the \mathbb{C} matrix] is not strictly positive definite because of the zero coefficients on the diagonal. This fact makes pivoting necessary when solving [the system] by direct methods and limits the applicability of almost all iterative solution techniques. What is also well-known is that for certain combinations of the approximation velocity and pressure spaces, the uniqueness of the discrete solution may not be guaranteed. This is due to the occurrence of spurious pressure modes in the pressure approximation space.”

The stability of mixed finite element methods boils down to properties of the null space of the matrix \mathbb{G} . An approximation is unstable if $\mathbb{G} \cdot \vec{p} = \vec{0}$ where \vec{p} corresponds to some spurious pressure mode different from the constant value pressure. Note that if $\mathbb{G} \cdot \vec{p} = \vec{0}$, then $(\vec{0}, \vec{p})^T$ is a null vector of the homogeneous system.

The basic idea behind stabilization is to relax the incompressibility constraint in a special way so that this vector is no longer a null vector of the resulting coefficient matrix, and the discrete solutions satisfy rigorous error bounds. In other words the idea consists in regularising the system by replacing the zero block by an appropriate positive semi-definite matrix $-\mathbb{C}$ [1167].

We will here look at so-called local and global jump methods (and their various flavours), the macro-element method, as well as the penalty method (which is not really a stabilisation method, as we will see).

- global jump stabilisation: Hughes, Franca, and Balestra [607] (1987), Norburn and Silvester [945] (1998), Douglas and Wang [344] (1989), Christon and Cook [255] (2001), Cao [207] (2003), Eguchi [362] (2003), Y. Cao [1376] (2006)

- local jump stabilisation: Silvester and Kechkar [1167] (1990), Kechkar and Silvester [686] (1992), Vincent and Boyer [1326] (1992), Cao [207] (2003), Qin and Zhang [1025] (2007), Christon [254] (2002), Christon and Cook [255] (2001), Liao and Silvester [787] (2013), Y. Cao [1376] (2006)
- stabilisation through macro-elements: Fortin and Boivin [402] (1990), LeTallec and Ruas [774] (1986), LeTallec [773] (1981)

In effect, these jump stabilisation techniques provide an a-priori filter for the weakly unstable pressure modes associated with the $Q_1 \times P_0$ element.

Consistency: in Barth, Bochev, Gunzburger, and Shadid [50]: “we should define what we mean by a **consistent method**; perhaps a more apt terminology would be variationally consistent. In standard usage, consistency of numerical schemes for partial differential equations requires that the pointwise truncation error vanish as the grid size goes to zero; i.e., if one substitutes a smooth solution of the partial differential equation into the numerical scheme, then the residual is at least $o(h)$, where h denotes the grid size. Finite element schemes are not, in general, consistent in this sense. However, for standard finite element methods, sufficiently smooth exact solutions of the partial differential equations exactly satisfy the variational equation that defines the discrete finite element equations. This is what we mean by a consistent finite element scheme. This allows us to differentiate between the methods we consider in this paper and methods which are not consistent in this latter sense. For example, penalty methods for the Stokes problem are not consistent finite element methods since substitution of an exact solution into the discrete equations leaves a residual that is proportional to the penalty parameter. Thus, we consider only methods that do not suffer from this type of variational inconsistency.”

As explained in Elman book: “to ensure consistency we require $\vec{1} \in \text{null}(\mathbb{C})$ (this precludes the use of inconsistent ‘penalty methods’) and we require $\vec{p}^T \cdot \mathbb{C} \cdot \vec{p} > 0$ for all spurious pressure modes $\vec{p} \neq \vec{1}$ in $\text{null}(\mathbb{G})$.”

define jump operator !! looking at sike90, I realise that I don't understand how the(ir) jump operator works in practice. Eq on page 78?

The driving question behind all this, besides my wanting to understand these stabilisation schemes better, is the fact that 1) none of the existing publishing literature seems to address the problem of large and/or sharp viscosity contrasts/variations; 2) almost all papers deal with regular meshes and rectangular elements.

Penalty The conventional way of computing a regularisation matrix \mathbb{C} is to use a penalty formulation. In the framework presented in Silvester and Kechkar [1167] (1990), the standard penalty method corresponds to the specific choice of

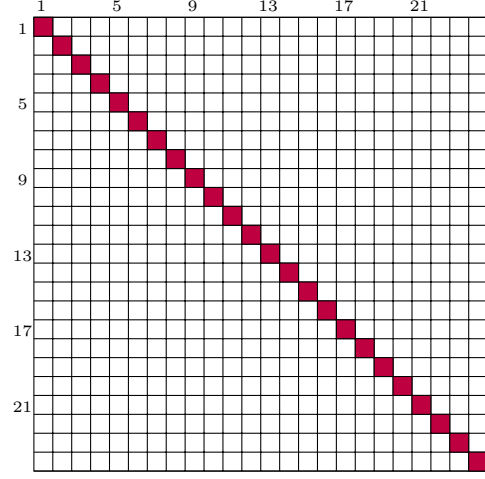
$$C(q^h, p^h) = \epsilon \int_{\Omega} q^h p^h dV = \epsilon \mathbb{M}_p \quad (7.4)$$

with $\epsilon > 0$ and \mathbb{M}_p is the pressure mass matrix. For a regular grid of squares with size h , it follows that

$$C_{ij} = 0 \quad \text{if } i \neq j \quad \text{and} \quad C_{ii} = \epsilon \int_{\Omega_i} dV = h^2 \epsilon$$

so that the stabilisation matrix is diagonal:

| | | | | | |
|----|----|----|----|----|----|
| 19 | 20 | 21 | 22 | 23 | 24 |
| 13 | 14 | 15 | 16 | 17 | 18 |
| 7 | 8 | 9 | 10 | 11 | 12 |
| 1 | 2 | 3 | 4 | 5 | 6 |



Two-dimensional grid composed of 6×4 elements on the left and the resulting sparsity pattern of the \mathbb{C} matrix on the right.

In Zienkiewicz & Taylor (Section 4.8.2) the matrix is also written as $\mathbb{C} = \epsilon h^2 \mathbf{1}$ where $\mathbf{1}$ is the unit matrix.

It is stressed here that the penalty technique does not stabilise an unstable mixed method [1167]. A small penalty parameter means that the original problem is solved quite accurately.

See Cuvelier, Segal, and Steenhoven [298] (1986) for some more details about the penalty method. The approach above is similar to the one presented in Section ???. The only difference is that instead of replacing the pressure in the momentum equation by $p = \lambda \vec{\nabla} \cdot \vec{v}$ we keep both velocity and pressure as unknowns and we take $\epsilon = \lambda^{-1}$. Since normally $\lambda \gg \eta$ then $\epsilon = \lambda^{-1}$ must be small. As explained in Silvester & Kechkar [1167], “despite its theoretical attraction, the penalty technique breaks down in practice because of its sensitivity to the particular choice of penalty parameter”.

The Stokes system for a single element then writes

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\epsilon \mathbb{M}_p \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix}$$

The second line yields

$$\mathbb{G}^T \cdot \vec{v} - \epsilon \mathbb{M}_p \cdot \vec{p} = \vec{h}$$

or,

$$\vec{p} = \frac{1}{\epsilon} \mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{v} - \vec{h})$$

which can be re-introduced in the first line:

$$\mathbb{K} \cdot \vec{v} + \mathbb{G} \cdot \frac{1}{\epsilon} \mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{v} - \vec{h}) = \vec{f}$$

or,

$$\left(\mathbb{K} + \frac{1}{\epsilon} \mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T \right) \cdot \vec{v} = \vec{f} + \frac{1}{\epsilon} \mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \vec{h}$$

This elimination could be carried out element by element so that one only solves for the velocity degrees of freedom.

- reduced integration? sike90 does not say anything about it.
- condition number explodes since $\frac{1}{\epsilon} \gg \eta$

Dohrmann and Bochev [336] (2004) state: “Penalty methods are another category of non-residual based regularizations. They, however, differ from stabilized methods in the sense that application of a penalty does not circumvent the inf–sup condition and only serves to uncouple pressure from velocity. In this sense, penalty methods should be viewed as solution, rather than stabilization procedures for the mixed equations.”

Global jump This method is explained in Silvester & Kechkar (1990) and the authors state that it was introduced by Hughes & Franca [607] in which a general theoretical framework for analysing global stabilisation techniques is presented. Using this framework, optimum rates of convergence for the $Q_1 \times P_0$ method stabilised with global jumps are established.

The global jump stabilisation formulation introduces a pressure diffusion operator that perturbs the incompressibility constraint. The global jump formulation insures mass conservation in a global sense since the null space of the stabilising matrix constrains the constant-pressure vector. However, the global jump stabilisation smears the div-free constraint over a small region, i.e., the divergence is not zero at the element level [257].

Consider the stabilisation term

$$C(q^h, p^h) = \beta h \sum_{s=1}^{N_s} \int_{\partial\Omega_s} \llbracket q^h \rrbracket \llbracket p^h \rrbracket ds \quad (7.5)$$

in which h is the mesh parameter (defined locally), $\llbracket \cdot \rrbracket$ is the jump operator, and $\beta > 0$ is a stabilising parameter. The summation is over *all* interior inter-element edges.

To illustrate, consider element 9 in the mesh consisting of equally sized squares represented here:

(tikz_globaljump.tex)

| | | | | | |
|----|----|----|----|----|----|
| 19 | 20 | 21 | 22 | 23 | 24 |
| 13 | 14 | 15 | 16 | 17 | 18 |
| 7 | 8 | 9 | 10 | 11 | 12 |
| 1 | 2 | 3 | 4 | 5 | 6 |

Element 9 has four direct neighbours: 3, 8, 10, and 15. The stabilisation term for this element involves the sum over its four neighbours:

$$\begin{aligned} \beta h \sum_{s=1}^4 \int_{\partial\Omega_s} \llbracket q^h \rrbracket \llbracket p^h \rrbracket ds &= \beta h^2 [(p_9 - p_3) + (p_9 - p_8) + (p_9 - p_{10}) + (p_9 - p_{15})] \\ &= \beta h^2 (4p_9 - p_3 - p_8 - p_{10} - p_{15}) \end{aligned}$$

The integral along each edge is simply the pressure difference across the edge multiplied by the edge surface/length which happens to be constant in this case. This means that the in the matrix \mathbb{C} , there will be entries on the 9th line at columns 3, 8, 10, and 15.

Be careful, let us now turn to element 6: it has 2 neighbours (5 and 12), so that the stabilisation term for this element involves the sum over its two neighbours:

$$\beta h^2 [(p_6 - p_5) + (p_6 - p_{12})] = \beta h^2 (2p_6 - p_5 - p_{12})$$

And looking now at element 23: it has three neighbours (17, 22, and 24), so that the stabilisation term for this element involves the sum over its three neighbours:

$$\beta h^2 [(p_{23} - p_{17}) + (p_{23} - p_{22}) + (p_{23} - p_{24})] = \beta h^2 (3p_{23} - p_{17} - p_{22} - p_{24})$$

The resulting assembled \mathbb{C} matrix is shown here:

Eguchi (2003) adds a linear form to the global stab \mathbb{C} to suppress the pressure nullspace.

Local jump According to Silvester and Kechkar [1167], the deficiencies of the global jump method can be overcome by a straightforward modification. Assume that the elements can now be assembled into N_m disjoint macro-elements of 2×2 elements, as shown in grey on the following figure:

(tikz_localjump.tex)

| | | | | | |
|----|----|----|----|----|----|
| 19 | 20 | 21 | 22 | 23 | 24 |
| 13 | 14 | 15 | 16 | 17 | 18 |
| 7 | 8 | 9 | 10 | 11 | 12 |
| 1 | 2 | 3 | 4 | 5 | 6 |

Consider now the bilinear form given by

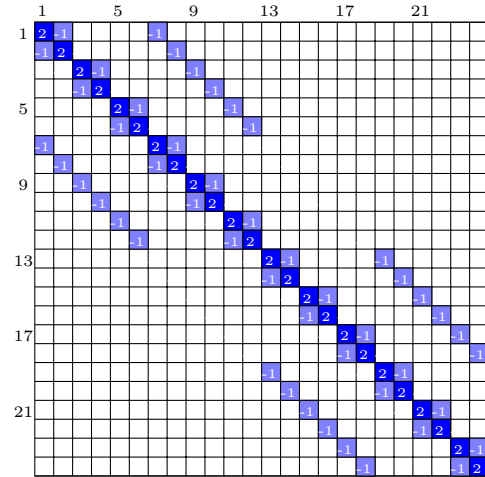
$$C(q^h, p^h) = \beta h \sum_{m=1}^{N_m} \sum_{i=1}^4 \int_{\partial\Omega_m} \llbracket q^h \rrbracket \llbracket p^h \rrbracket ds$$

where the first summation is over all 2×2 macroelements, and the second summation runs over all inter element edges strictly within each macroelement.

The form of the stabilisation matrix \mathbb{C} is similar to that above except that there is now a local basis.

For instance, considering again element 9, it now belongs to the second macro-element and therefore only 'sees' neighbours 3 and 10. The resulting \mathbb{C} matrix is shown on the figure here after and its structure is obviously different than in the global stabilisation case, albeit also pentadiagonal.

| | | | | | |
|----|----|----|----|----|----|
| 19 | 20 | 21 | 22 | 23 | 24 |
| 13 | 14 | 15 | 16 | 17 | 18 |
| 7 | 8 | 9 | 10 | 11 | 12 |
| 1 | 2 | 3 | 4 | 5 | 6 |



Remark. Obviously, one could re-number the elements so that matrix \mathbb{C} is block diagonal: macroelement 1 would contain elements 1,2,3,4, macroelement 2 would contain elements 5,6,7,8, etc ...

According to Silvester and Kechkar [1167] or Chibani and Kechkar [233], the advantages of this local method over the global jump formulation are:

1. implementation is more straightforward because for assembly purposes each 2×2 block of elements can be treated as a single macroelement ⁷

⁷same here, not sure what they mean by this

2. mass is conserved locally (over a macroelement), using the global jump formulation mass is only conserved globally
3. robustness is improved in the sense that the discrete velocity solution is less sensitive to the magnitude of β , the influence of the stabilisation matrix being localised (will need to be tested numerically!).

Remark. *The globally stabilised formulation corresponds to the extreme case of a local stabilisation based on a single macro-element [488].*

One of the features of the local stabilisation is that if the discrete incompressibility constraints are added together then the jump terms sum to zero in each macro element. Indeed, let us consider the following macro element:

(tikz_macro.tex)



The corresponding matrix (making abstraction of the β term) writes:

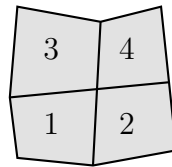
$$h^2 \begin{pmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & 0 & -1 \\ -1 & 0 & 2 & -1 \\ 0 & -1 & -1 & 2 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$$

and the row/column sum of its entries is always null. Also, \mathbb{C} is obviously positive semi-definite [1167].

Gresho & Sani [488] state: “ This is crucially important to the success of the method since it implies that the local incompressibility of the $Q_1 \times P_0$ method is retained after stabilisation (albeit over macro-elements). It also suggests that a good strategy when constructing the partition is to form macro-elements containing as few elements as possible. Once a suitable macro-element partitioning has been formed, the local stabilisation matrices can be calculated by running through the component elements, summing jump contributions corresponding to the internal edges.”

If one now considers the following irregular macro-element,

(tikz_macro2.tex)



the corresponding matrix is given by⁸

$$\tilde{h} \begin{pmatrix} h_{12} + h_{13} & -h_{12} & -h_{13} & 0 \\ -h_{12} & h_{12} + h_{24} & 0 & -h_{24} \\ -h_{13} & 0 & h_{13} + h_{34} & -h_{34} \\ 0 & -h_{24} & -h_{34} & h_{24} + h_{34} \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$$

where h_{ij} is the length/surface of the edge between elements i and j . The reference length \tilde{h} may be computed by simply defining it to be the average diameter of the constituent elements.

⁸I suspect it should involve the normal vectors to the edges ...?

Remark. In three dimensions, the $2 \times 2 \times 2$ block is the obvious starting point for stabilising $Q_1 \times P_0$ [488].

Perhaps the most serious potential drawback of the local framework is that stability is only guaranteed if the stabilisation parameter β is bigger than some critical value β_0 , which needs to be estimated.

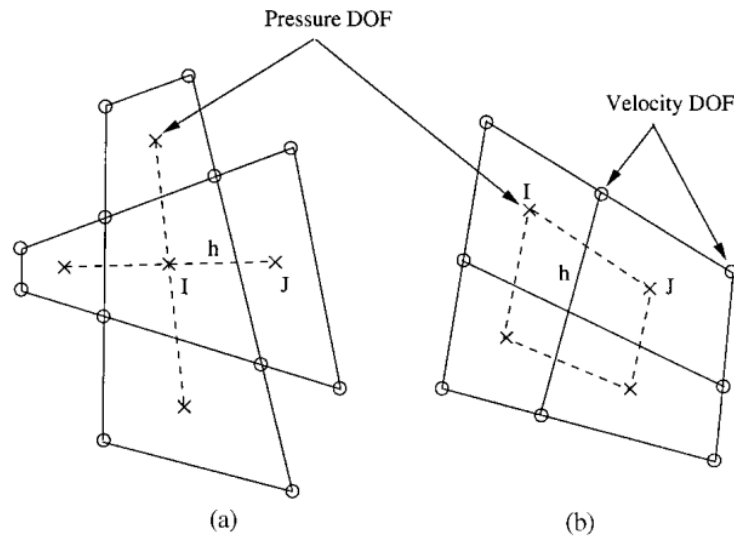
It can be estimated that $\beta = 1/4$ in 2D and $\beta = 1/6$ in 3D (see Gresho & Sani [488, p636] for a detailed derivation, see also Vincent and Boyer [1326] (1992)).

Silvester & Kechkar state: “The advantages of the stabilisation procedures over the penalty method are especially relevant to the discretisation of 3D incompressible flow problems, since iterative solution methods have to be used. Similar stabilisation techniques to those described here are applicable to the three-dimensional version of the $Q_1 \times P_0$ mixed method”.

The conclusion in [1167] is as follows: The local jump formulation proves to be an efficient method for a priori filtering of spurious pressure modes. It cleanly stabilises the $Q_1 \times P_0$ mixed method without compromising its simplicity and resulting efficiency; in particular, it is very robust with respect to the magnitude of the stabilisation parameter.

It is reported in Gresho & Sani [488] that when using an iterative solver, iteration counts are only independent of the grid in the stabilised cases: using the raw $Q_1 \times P_0$ method the iteration counts significantly increase with decreasing h . Note that the deterioration of the condition number of the matrix with decreasing h is worse in 3D than in 2D (but bear in mind that one almost always use higher resolutions in 2D than in 3D, so it does not help). 3D is also discussed in [chsu97].

A way to look at the global vs. local stabilisation schemes is presented on the following figure from Christon (2002) [254]:



Element configuration for pressure stabilization: (a) global jump; (b) local jump.

Remark. The locally stabilised $Q1P0$ and $P1P0$ elements have been analysed in Kechkar & Silvester [686]. Penalty, global and local approaches are mentioned in Vincent & Boyer (1992) but only the local jump stabilisation is used.

Silvester (1994) investigates the value for β for $Q1P0$ and arrives at $0.0615 \leq \beta \leq 0.25$. Chang & Sugiyama (1997) report that “a value between 0.01 and 0.1 appears to work well for most applications conducted so far”. In [371] the authors state that $\beta = \frac{1}{4}$ is an idea value which “ensures stability independently of the rectangle aspect ratio”.

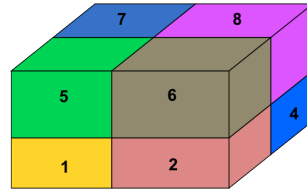
Norburn & Silvester (1998) state: “Although the ‘optimal parameter’ (the value of β which minimises the discretisation error on a given mesh) is impossible to determine a priori, good parameter choices (usually over-estimating the optimal parameter) can be found by minimising the condition

number of the pressure Schur complement matrix. [...] The motivation for choosing the parameter value β which minimises the condition number (ratio of largest to smallest eigenvalue) of the pressure Schur complement is that this quantity roughly determines the rate of convergence of Uzawa type iteration methods.” They conclude that $\beta \simeq 0.1$ is appropriate for P1P0.

Rather interestingly we see that choosing the ‘best’ β is primarily based in the literature on the Schur complement condition number, and less on the accuracy of the solutions (although it is sometimes assessed, see Fig. 4 of Norburn and Silvester [945] (1998)).

Finally it is worth mentioning a recent paper by Chibani and Kechkar [233] (2020) who present modified local jump stabilisation schemes which effectively only take 2 or 1 one pressure jump into account instead of 4 per macro-element.

In three dimensions the matrix \mathbb{S} is obtained by assembling the submatrices for the macroelements which, in general, are made up of 8 (i.e., $2 \times 2 \times 2$) adjacent elements, as shown in the following sketch [chsu97]:



Numbering of hexahedrons in a macroelement.

Element No. 3 is behind No. 1 and below No. 7.

For a macroelement of 8 elements as ordered in the above sketch, the submatrix is defined below:

$$\mathbb{S} = \begin{pmatrix} h_{12} + h_{13} + h_{15} & -h_{12} & -h_{13} & 0 & -h_{15} & 0 & 0 & 0 \\ -h_{21} & h_{21} + h_{24} + h_{26} & 0 & -h_{24} & 0 & -h_{26} & 0 & 0 \\ -h_{31} & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ -h_{51} & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

in which h_{ij} is the length scale for elements ‘i’ and ‘j’ in the macroelement. The matrix \mathbb{S} is symmetric, since $h_{ij} = h_{ji}$.

Two approaches are possible for computing the length scale h_{ij} :

- the square root of the interior inter-element surface area between elements ‘i’ and ‘j’
- the quotient of the interior inter-element surface area divided by the cube root of the average element volume of the macroelement.

Macro-element One source for this stabilisation approach is Section 5.3.2 of the book by Elman, Silvester and Wathen [371]. The \mathbb{C} matrix for the 2D macroelement is shown to be:

$$\mathbb{C} = \beta \begin{pmatrix} 1 & -1 & 1 & -1 \\ -1 & 1 & -1 & 1 \\ 1 & -1 & 1 & -1 \\ -1 & 1 & -1 & 1 \end{pmatrix}$$

and the authors suggest $\beta = \frac{1}{4}h_x h_y$.

Also see [402] (1990) .

Numerical scaling of \mathbb{C} Since the matrix \mathbb{K} contains the viscosity, it is to be expected that the magnitude of the entries in matrix \mathbb{C} must somehow follows the values in the $\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$ term. Indeed the Schur complement is $\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}$ so if the entries in \mathbb{C} are wrongly scaled it will either have no effect at all or will alter the solution too much. This is indeed what is advocated in Christon (2002) [254]:

The global jump stabilization formulation attempts to control the jump in pressure across element boundaries, and results in a PPE that is perturbed by a ‘pressure-diffusion’ operator. The off-diagonal entries in the global jump stabilization matrix are defined as

$$S_{IJ} = \beta \frac{[C^T M_L^{-1} C]_{IJ}}{\Gamma_{IJ}} \int_{\Gamma_{IJ}} [\![\psi_I]\!] [\![\psi_J]\!] d\Gamma, \quad (23)$$

where I and J identify adjacent elements that share a common face as shown in Figure 1(a). Here, Γ_{IJ} represents the shared inter-element boundary, $[\![\cdot]\!]$ is the jump operator, and β is a non-dimensional scaling parameter. For the $Q_1 Q_0$ element, the pressure approximation is piecewise constant with $\psi_I = 1$ inside Ω_e and zero outside. In two dimensions, Γ_{IJ} represents the length of the element edge shared by element I and J , and in three dimensions, it represents the area of

In the paper S stands for \mathbb{C} , and C stands for \mathbb{G} . Also, because Christon is solving the N-S equations, and because of his algorithmic choice to do so, there is a velocity mass matrix where our \mathbb{K} block resides. As such his implementation showcases the lumped mass matrix in the equation above rather than the lumped \mathbb{K} matrix.

In the paper Christon states that the PPE term $C^T M_L^{-1} C$ is symmetric. This is obviously true since it is a scalar for the $Q_1 \times P_0$ element. Am I missing something here?

Using our notations, the off-diagonal entries of the stabilisation matrix \mathbb{C} for the global approach becomes:

$$\mathbb{C}_{ef} = -\beta (\mathbb{G}^T \cdot \mathbb{K}_L^{-1} \cdot \mathbb{G})_e \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} [\![\psi_e]\!] [\![\psi_f]\!] d\Gamma \quad \text{for } e \neq f \quad (7.6)$$

where e and f identify adjacent elements that share a common face, Γ_{ef} represents the shared inter-element boundary (a length in 2D, a surface in 3D), β is a non-dimensional scaling parameter, and \mathbb{K}_L is the row-wise lumped \mathbb{K} matrix⁹. For the $Q_1 \times P_0$ element, the pressure approximation is piecewise constant with $\psi_i = 1$ inside the element and zero outside.

The inclusion of the $\mathbb{G}^T \cdot \mathbb{K}_L^{-1} \cdot \mathbb{G}$ term in the stabilisation yields proper dimensionality of the stabilization matrix (the integrand is dimensionless so that the dimensions of \mathbb{C} are those of the elemental Schur complement block), accounts for scaling due to irregular elements, and still preserves the symmetry [254].

Finally, the diagonal element of C for element e is computed as

$$\mathbb{C}_{ee} = \beta (\mathbb{G}^T \cdot \mathbb{K}_L^{-1} \cdot \mathbb{G})_e \sum_{f \neq e} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} [\![\psi_e]\!] [\![\psi_f]\!] d\Gamma$$

We find that the sum of the terms on the row corresponding to e is indeed zero, consistency is ensured even for non-rectangular elements.

There is however a major problem with this approach: even when the viscosity is constant in the domain, Eq. (7.6) does not yield a symmetric matrix if elements are not identical in shape. Since $\mathbb{C}_{ef} \propto (\mathbb{G}^T \cdot \mathbb{K}_L^{-1} \cdot \mathbb{G})_e$ then $\mathbb{C}_{ef} \neq \mathbb{C}_{fe}$! I then suspect that Christon’s notation $[C^T M_L^{-1} C]_{IJ}$ indicates that some care must be taken so as to ensure $S_{IJ} = S_{JI}$ but it is not further specified. We will then have to figure this out.

Let us consider the case of square elements of size $h_x = h_y = h$. Then the \mathbb{K}_e and \mathbb{G}_e matrices

⁹The quantity $\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$ for $Q_1 \times P_0$ elements is a scalar, which is rather convenient as it gives in a simple way the scaling for the stabilisation term. However the inverse of \mathbb{K} is costly so there is a cheaper alternative which consists in lumping it so that it becomes diagonal and its inverse is then trivial.

are given by:

$$\mathbb{K}_e = \begin{pmatrix} 1 & 0.25 & -0.5 & -0.25 & -0.5 & -0.25 & 0 & 0.25 \\ 0.25 & 1 & 0.25 & 0 & -0.25 & -0.5 & -0.25 & -0.5 \\ -0.5 & 0.25 & 1 & -0.25 & 0 & -0.25 & -0.5 & 0.25 \\ -0.25 & 0 & -0.25 & 1 & 0.25 & -0.5 & 0.25 & -0.5 \\ -0.5 & -0.25 & 0 & 0.25 & 1 & 0.25 & -0.5 & -0.25 \\ -0.25 & -0.5 & -0.25 & -0.5 & 0.25 & 1 & 0.25 & 0 \\ 0 & -0.25 & -0.5 & 0.25 & -0.5 & 0.25 & 1 & -0.25 \\ 0.25 & -0.5 & 0.25 & -0.5 & -0.25 & 0 & -0.25 & 1 \end{pmatrix} \quad \mathbb{G}_e = h \begin{pmatrix} +1/2 \\ +1/2 \\ -1/2 \\ +1/2 \\ -1/2 \\ -1/2 \\ +1/2 \\ -1/2 \end{pmatrix}$$

so that the Schur complement is

$$\mathbb{S}_e = \mathbb{G}_e^T \cdot \tilde{\mathbb{K}}_e^{-1} \cdot \mathbb{G}_e = \frac{2}{3}h^2$$

In that case we almost recover the expression of for example the macro-element.

Dealing with viscosity contrasts/large variations This topic is almost never discussed as many papers consider the standard Stokes equations with $\eta = 1$ (and also regular meshes made of identical elements). This is not problematic in engineering where often the fluid in question has a constant viscosity (or the equations have been rendered dimensionless). However, in geodynamical applications we know that the viscosity field can showcase very sharp gradients (sinking/rising objects, shear bands, free surface, ...). In what follows we assume that each element e has an effective viscosity η_e .

The scaling of the \mathbb{C} matrix in the previous section is not formulated when viscosity contrasts from one element to the other are present. For example scaling the row entries of the \mathbb{C} matrix by the element viscosity still yields a structurally symmetric matrix, but not a numerically symmetric one which is problematic since we have seen that \mathbb{C} must be semi-positive definite. Some form of viscosity averaging must then take place between adjacent elements so that the contribution from element e to f is exactly the same as f to e .

In order for the stabilisation to remain consistent it must satisfy $\mathbb{C} \cdot \vec{1} = 0$, i.e. it should have zero effect on a constant pressure field, which then forces the sum of the entries for each row (or column) to be null. This requirement makes the above viscosity averaging idea very difficult in practice in the global case (satisfying both symmetry and consistency).

Let us start by defining the elemental Schur complement

$$\tilde{\mathbb{S}}_e = |(\mathbb{G}^T \cdot \mathbb{K}_L^{-1} \cdot \mathbb{G})_e|$$

and then

$$\tilde{\mathbb{S}}_{ef} = \phi(\tilde{\mathbb{S}}_e, \tilde{\mathbb{S}}_f)$$

where ϕ is a function to be specified later such that $\phi(x, y) = \phi(y, x)$. The local and global jump stabilisation can then be formulated as follows:

$$\mathbb{C}_{ef} = -\beta \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma \quad \text{for } e \neq f \quad (7.7)$$

supplemented by

$$\mathbb{C}_{ee} = \beta \sum_{f \neq e} \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma$$

In this case the matrix \mathbb{C} is symmetric and consistent!

Recap : For each element compute its (scalar) Schur complement

$$\tilde{\mathbb{S}}_e = |(\mathbb{G}_e^T \cdot \tilde{\mathbb{K}}_e^{-1} \cdot \mathbb{G}_e)|$$

We here assume that elements can have different viscosities and/or shape so that \mathbb{S}_e varies from element to element.

- Global jump Assuming elements e and f share an edge, build the \mathbb{C} matrix as follows:

$$\mathbb{C}_{ef} = -\beta \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma \quad \text{for } e \neq f \quad (7.8)$$

supplemented by

$$\mathbb{C}_{ee} = \beta \sum_{f \neq e} \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma$$

- Local jump Let \mathcal{M}_e be the macroelement that element e is in. Assuming elements e and f share an edge, build the \mathbb{C} matrix as follows:

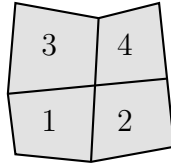
$$\mathbb{C}_{ef} = -\beta \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma \quad \text{for } e \neq f \quad \text{and} \quad f \in \mathcal{M}_e \quad (7.9)$$

supplemented by

$$\mathbb{C}_{ee} = \beta \sum_{f \neq e, f \in \mathcal{M}_e} \tilde{\mathbb{S}}_{ef} \frac{1}{\Gamma_{ef}} \int_{\Gamma_{ef}} \llbracket \psi_e \rrbracket \llbracket \psi_f \rrbracket d\Gamma$$

- Macroelement stab Let \mathcal{M}_e be the macroelement that element e is in.

(tikz_macro2.tex)

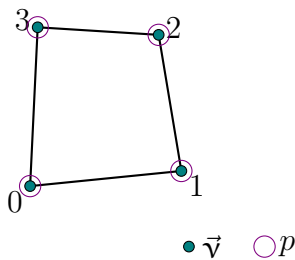


This is less straightforward than the local jump since (for example) elements 1 and 4 do not have an edge in common so that the jump operator cannot be used. One could think of assigning all four elements a single effective viscosity but elements shapes/sizes can differ and the matrix is then not necessarily consistent. One should probably go back to the derivations in Elman, Silvester, and Wathen [371] and see whether a more generic form of the macroelement stabilisation matrix \mathbb{C} could be derived?

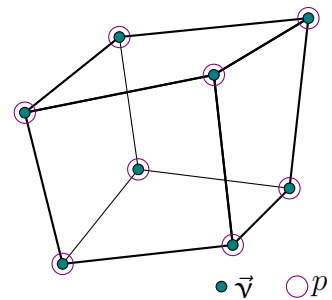
7.3.8 The stabilised bi/tri-linear velocity - bi/tri-linear pressure element ($Q_1 \times Q_1$ -stab)

pair_q1q1stab.tex

4 vel. nodes, 4 press. nodes




8 vel. nodes, 8 press. nodes



The $\mathbf{Q}_1 \times Q_1$ element is not LBB-stable but it can be stabilised. Despite some applications in geodynamics (it is used in Burstedde, Ghattas, Stadler, Tu, and Wilcox [188] (2009) and Burstedde et al. [189] (2013)), it is not appropriate for buoyancy-driven flows, as shown in Thieulot and Bangerth [1260].

See Norburn and Silvester [944] (2001) for a fourier analysis of the normal and stabilised (a la Hughes, Franca, and Balestra [606] (1986)) $\mathbf{Q}_1 \times Q_1$ element. Stabilisation is worked out in Dohrmann and Bochev [336] (2004), Bochev, Dohrmann, and Gunzburger [101] (2006), and Bochev and Dohrmann [102] (2006).

- $\mathbf{Q}_1 \times P_0$ -stab. Pro: stabilisation can be switched off; Con: stabilisation for deformed elements? problem near boundaries: incomplete stencil? choice of parameter β .
- $\mathbf{Q}_1 \times Q_1$ -stab. Pro: easier to implement than $\mathbf{Q}_1 \times P_0$ -stab, stabilisation local to element, easier when elements are not rectangular, no free parameter; Con: stabilisation cannot be switched off.

 **Relevant Literature:** Schneider, Raithby, and Yovanovich [1130], Tezduyar, Mittal, Ray, and Shih [1249], Tezduyar [1247], Gresho, Chan, Christon, and Hindmarsch [485], Idelsohn, Storti, and Nigro [619], Knobloch and Tobiska [714], Franca, Oliveira, and Sarkis [411], and Li, He, and Chen [781]. See Braack and Lube [127] for a review of local projection stabilisation for incompressible flow problems.

This unstable pair is also used in ice sheet modelling Helanow and Ahlkrone [561], Zhang, Ju, Gunzburger, Ringler, and Price [1405], Zwinger, Greve, Gagliardini, Shiraiwa, and Lyly [1444]. A $P_1 \times P_1$ version of it is used in Karabelas, Haase, Plank, and Augustin [670] (2020).

7.3.9 The Rannacher-Turek element - rotated $Q_1 \times P_0$

This element is the natural quadrilateral analogue of the well-known triangular P_1^{nc} Stokes element of Crouzeix-Raviart [290]. This element is sometimes called $\mathbf{Q}_1^{rot} \times Q_0$ or the Rannacher-Turek element [650, Section 3.6.5] (see also Appendix B.4, example B.53 of John [650]). This rectangular nonconforming [289] element is termed the rotated \mathbf{Q}_1 element because of the fact that $r^2 - s^2$ can be generated from rs (occurring in the bilinear Q_1 element) by a rotation of 45° [224, p93]. The velocity approximation is achieved by rotated dim-linear functions that have continuous degrees of freedom on the faces of the mesh cells as we have seen in Section 5.3.15. This element was introduced in Rannacher & Turek (1992) [1045] has been proven to satisfy the inf-sup condition. It has been studied comprehensively in Schieweck (1997)¹⁰, [1162] and in Turek [1293, 1291]. Superconvergence properties have also been reported [875, 874]. It has been used in 2D [824] and 3D [710, 445] and forms the basis of the FeatFlow software¹¹. It is used in the PhD thesis of Gastaldo [441] and Ouazzi [965]. It has been successfully coupled to multigrid solvers [226, 1295]. This element has been compared to the stabilised $\mathbf{Q}_1 \times P_0$ element [787]. It is mentioned in [527]

It essentially comes in two flavours, the Middle Point (MP) and the Mid Value (MV) one.

Remark. John [650] explains that: *"For the point-value-oriented non-conforming finite element spaces (MP), the value of the Dirichlet boundary condition in the barycenter of the faces at the boundary is taken. Using the mean-value-oriented spaces (MV), one computes the integrals of the boundary condition on these faces and normalizes with the area of the faces to set the boundary values. In the case of homogeneous Dirichlet boundary conditions, the boundary values computed in both ways are zero."*

¹⁰Habilitation thesis in German

¹¹<http://www.featflow.de/en/index.html>

Remark. John also makes a very important point: "There are also unmapped (non-parametric) versions of these finite element spaces, which define the polynomials directly on the mesh cell K . It is shown in Rannacher and Turek (1992) [1045] that these versions are inf-sup stable on more general meshes than the mapped (parametric) version of the $\mathbf{Q}_1^{\text{rot}} \times Q_0$ finite element, e.g., on strongly nonuniform meshes. Considering all four types of $\mathbf{Q}_1^{\text{rot}} \times Q_0$ finite elements, the optimal order of convergence on perturbed meshes is achieved only by the mean-value-oriented version of the unmapped $\mathbf{Q}_1^{\text{rot}} \times Q_0$ finite element.

Mahmood *et al.* [824] mention a very important fact: "The chosen nonconforming element requires additional stabilization for handling the deformation tensor formulation due to missing Korn's inequality [592, 713, 144]. To this end we employ the standard edge oriented stabilization [1295, 1294] in our simulations." This is a rather unfortunate fact that although LBB stable this element needs an additional term in the weak form (see Turek *et al.* (2002) [1295]) so as to suppress parasitic velocity modes when the div-grad formulation of the Stokes equation is used (as opposed to the Laplace formulation – see [341, Section 6.5.2]).

This element is used in Hansbo, Larson, and Larson [528] (2001) in the context of near incompressible elasticity. It is mentioned that it does not fulfill the discrete Korn's inequality. It is then stabilised in a discontinuous Galerkin framework.



Relevant Literature Sheen [1155] (2020), **chen92** (1992) , **chen93** (1993)

7.3.10 The $P_1 \times P_0$ pair

example 3.70 in John [650] (book),

Elman Silvester Wathen say (5.3.3) that "it can be readily stabilized using the pressure jump stabilization together with an appropriate macroelement subdivision." See Norburn and Silvester [945] (1998) for globally and locally stabilised versions.

Qin and Zhang [1026] (2007) states: "the element is unstable for any mesh since the dimension of the discrete velocity space is always less than that of the pressure space (with Dirichlet boundary condition)." The authors explain a filter algorithm to make the element usable.

Arnold [28] (1993) states: "Unfortunately, this simplest possible Stokes element is notoriously unstable. On any tri- angulation with at least three vertices on the boundary the dimension of the pressure space exceeds that of the velocity space [...] and the finite dimensional problem is singular. Moreover, while the discrete velocity field u_h is uniquely determined (as it is for any conforming method for the Stokes problem), for this choice of elements u_h belongs to the space of divergence-free fields piecewise linear fields, and on many meshes, for example on a uniform diagonal mesh of the square [...], this space is known to reduce to zero. So even after accounting for the indeterminacy of the pressure we have no convergence."

Example 3.2 of Boffi, Brezzi, and Fortin [108] (2008) explains neatly the locking phenomenon and how to circumvent it via a so-called cross-grid macroelement. See also Hong, Kim, and Lee [588] (2003).

In his lecture notes¹², Guermond states "A simple alternative to the $Q_1 \times P_0$ element consists of using the $P_1 \times P_0$ element. Let \mathcal{T}_h be a mesh of D composed of affine simplices, and approximate the velocity with continuous piecewise linear polynomials and the pressure with (discontinuous) piecewise constants. Since the velocity is piecewise linear, its divergence is constant on each simplex. As a result, testing the divergence of the velocity with piecewise constants enforces the divergence to be zero everywhere. That is to say, the $P_1 \times P_0$ finite element yields a velocity approximation that is exactly divergence-free [...]. Unfortunately, this pair does not satisfy the inf-sup condition."

I have not found a paper yet which showcases its accuracy on a manufactured solution and compares it to other element pairs.

¹²<https://www.math.tamu.edu/~guermond/>

7.3.11 The $P_2 \times P_0$ pair

[108], [cakp18] stable (Kanschat book)

compared with P1NC-P0 and BR element in Carstensen, Köhler, Peterseim, and Schedensack [212] (2015).

7.3.12 The $Q_2 \times Q_0$ pair

Quadratic velocities, constant pressure. The element satisfies the inf-sup condition, but the constant pressure assumption may require fine discretisation. source?

I have implemented it in `STONE` ?? using the penalty method.

7.3.13 The $P_1 \times P_1$ -stabilised pair

Like its quadrilateral counterpart $Q_1 \times Q_1$, the $P_1 \times P_1$ pair is not stable and needs to be stabilised [945, 1224]. TerraNeo code uses stabilised with PSPG [56].

Norburn and Silvester [945]

7.3.14 The $P_1^+ \times P_1$ (MINI) pair in 2D & 3D

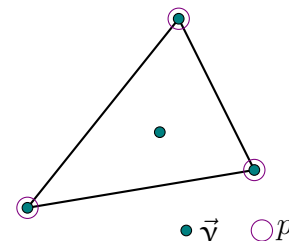
pair-mini.tex

The MINI element was first introduced in Arnold, Brezzi, and Fortin [26] (1984). It is also discussed in Section 3.6.1 of John [650] (2016) and in Section 6.1 of Boffi, Brezzi, and Fortin [108] (2008). It is thoroughly studied in Cioncolini and Boffi [258] (2019).

(tikz-mini.tex)

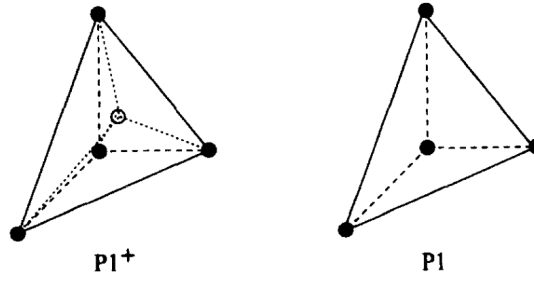
As explained in Braess [128], since the support of the cubic bubble is restricted to the element, the associated variable (dofs living on the bubble) can be eliminated from the resulting system of linear equations by static condensation. Also, the MINI element is cheaper than the Taylor-Hood element but it is commonly accepted that it yields a poorer approximation of the pressure.

4 vel. nodes, 3 press. nodes



Remark. Note that Franca and Oliveira [410] (2003) propose an equal-order-linear-continuous velocity-pressure variables which is enriched with velocity and pressure bubble functions to model the Stokes problem. They show by static condensation that these bubble functions give rise to a stabilized method involving least-squares forms of the momentum and of the continuity equations. In some cases their approach recovers the MINI element. Also check Ganesan, Matthies, and Tobiska [434] (2008).

The 3D MINI element is not very common but it is used for instance in Pichelin and Coupez [997] (1998) or Tommasi, Knoll, Vauchez, Signorelli, Thoraval, and Logé [1271] (2009). It is also said to be LBB stable in Reddy [1051, p180]. It is used in [427] phd thesis in the context of microstructures deformation modeling, which itself cites Cao, Montmitonnet, and Bouchard [206] (2013).



Velocity and pressure nodes for the 3D MINI element, taken from [997]

Note that this element is used in Braess & Wriggers (2000) [130] in the context of Arbitrary Lagrangian Eulerian finite element analysis of free surface flows, and also in Zlotnik, Diez, Fernandez, and Verges [1440] (2007) for subduction with X-FEM technique. . It is also mentioned in Nafa and Thatcher [920] (1993).


The 2D element is implemented in [STONE](#) ??.

7.3.15 The $P_2 \times P_1$ pair

pair_p2p1.tex

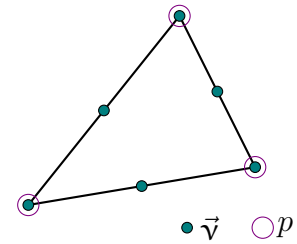
From Segal [1147]: “Taylor-Hood elements [1240] are characterized by the fact that the pressure is continuous in the region Ω . A typical example is the quadratic triangle ($P_2 \times P_1$ element). In this element the velocity is approximated by a quadratic polynomial and the pressure by a linear polynomial. One can easily verify that both approximations are continuous over the element boundaries.”

It can be shown, Segal (1979), that this element is admissible if at least 3 elements are used. The quadrilateral counterpart of this triangle is the $Q_2 \times Q_1$ element. Reddy and Gartling [1051, p179] also report this element to be LBB stable. It is also mentioned in Nafa and Thatcher [920].

 **Relevant Literature:** Schubert & Anderson [1141], Leng *et al.* [771], Cuffaro *et al.* [293]

(tikz_p2p1.tex)

6 vel. nodes, 3 press. nodes



7.3.16 The $P_2^+ \times P_{-1}$ pair (Crouzeix-Raviart)

pair_crouzeixraviart.tex

Since the $P_2 \times P_{-1}$ pair is not LBB stable [1051, p179], (see also table 3.13-1 of Gresho and Sani [488]) it is enhanced by a cubic bubble and is therefore called $P_2^+ \times P_{-1}$.

This element was first introduced in [290]. It is the element used in the MILAMIN code [299]. It is a seven-node triangle with quadratic velocity shape functions enhanced by a cubic bubble function and discontinuous linear interpolation for the pressure field [298]. This element is LBB stable and no additional stabilization techniques are required[371]. The '+' in its name stands for the bubble while the '-' stands for the discontinuous character of the pressure field: once again, it is P_1 over the element, but discontinuous across element edges.

Remark. Cuvelier *et al.* , 1986 [298] recommend a 6-point or 7-point quadrature rule for this element.

Remark. Segal [1147] explains for output purposes (printing, plotting etc.) the discontinuous pressures are averaged in vertices for all the adjoining elements. See also Fig. 7.3 of [298].


Remark. The simplest Crouzeix-Raviart element is the non-conforming linear triangle with constant pressure [298], see Section 7.3.27.

It is worth noting that this element has more degrees of freedom than the Taylor-Hood element for the same order of accuracy. However, since the bubble can be eliminated, one can design a modified version of this element.

Check Cuvelier book chapter 8 for modified element

Remark. I have once asked the (main) author of MILAMIN why he chose this element, for example over the $P_2 \times P_1$. His answer is as follows: "Elements with continuous pressure are incapable of converging in the L_{inf} norm for mechanical problems exhibiting pressure jumps such as the inclusion-host setup. During my MSc and PhD I was focusing on sharp heterogeneities, so this is why I decided to choose $P_2^+ \times P_{-1}$. You will see that it is also easy to invert the pressure mass matrix for such elements, which is really useful (both for the augmentation and preconditioning)."

This element is used by Poliakov and Podlachikov [1008] to study the deformation of the surface above a rising diapir. Note that they actually use a "13 point integration formula (Hughes 1987) for calculation of the stiffness matrix was used in order to conserve detailed information from the marker field in the coarse FEM mesh". It is also used in [24] in the context of a new free-surface stabilization scheme. It is the element used in LaCoDe [322]. It is mentioned in Section 6.2 in Boffi, Brezzi, and Fortin [108] (2008). It is compared to the $P_2 \times P_1$ element for the Navier-Stokes equations in Krahel and Bänsch [729] (2005).

 Relevant Literature: Hansbo and Larson [529]

7.3.17 The $P_2^+ \times P_1$ pair

This element pair is not to be mistaken for the Crouzeix-Raviart. Both share the same P_2^+ space for the velocity but this element has a continuous linear pressure. It is mentioned in Table 3.13-1 of Gresho and Sani [488]: "LBB stable. Second order. cubic bubble. good element". It is also mentioned in Soulaimani, Fortin, Ouellet, Dhett, and Bertrand [1181] (1987).

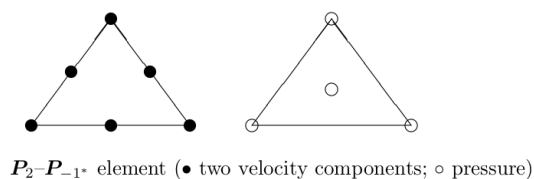
Implemented in [STONE](#) 120.

7.3.18 The $P_2 \times (P_1 + P_0)$ pair

pair_p2p1p0.tex

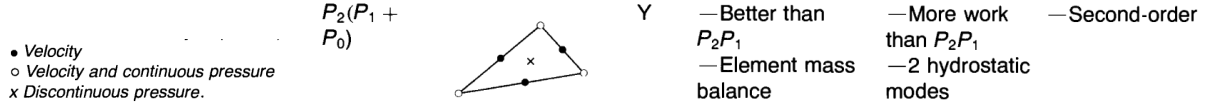
This element pair is discussed in 5.3.3 of Elman, Silvester and Wathen:

"[Another] possibility is to construct a hybrid pressure approximation by combining the continuous linear pressure approximation with the discontinuous constant pressure approximation. The resulting mixed method is referred to as the $P_2 - P_{-1*}$ approximation and enjoys the best of both worlds; it has locally incompressibility, and yet it does not have its accuracy compromised by the lower order pressure [of $P_2 \times P_0$]. Perhaps surprisingly, this element is also uniformly stable."



Taken from Elman, Silvester, and Wathen [371].

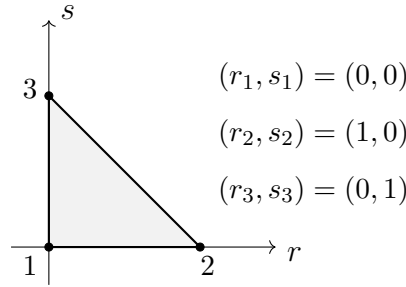
In Gresho & Sani table 3.13-1: “LBB stable yes. Better than P2P1, element mass balance. more work than P2P1. 2 hydrostatic modes. Second order.”



Taken from Gresho and Sani [488]. Unfortunately they do not provide a source for its origin, for the LBB-stability proof, or any source at all, actually.

Looking at the figure above, it is clear that the P_1 space is to be understood as a continuous pressure space, with an additional constant bubble.

For the continuous P_1 space, we have the following reference element



and the basis functions are simply

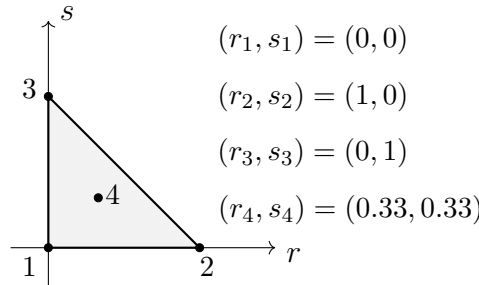
$$\mathcal{N}_1(r, s) = 1 - r - s \quad (7.10)$$

$$\mathcal{N}_2(r, s) = r \quad (7.11)$$

$$\mathcal{N}_3(r, s) = s \quad (7.12)$$

with the interpolation requirement $\mathcal{N}_i(r_j, s_j) = \delta_{ij}$ fulfilled, as well as $\sum_i \mathcal{N}_i = 1$.

Now, following the figure by Gresho and Sani, I build the reference element for the $P_1 + P_0$ space:



$P_1 + P_0$ means that the pressure inside the element is given by

$$p^h(r, s) = a\mathcal{N}_1(r, s) + b\mathcal{N}_2(r, s) + c\mathcal{N}_3(r, s) + d\mathcal{N}_4(r, s)$$

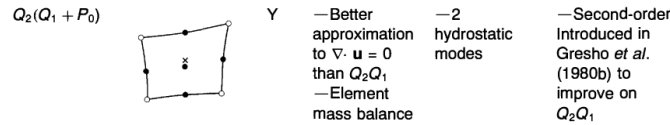
Note that it is then impossible to find a, b, c, d such that the interpolation requirement $\mathcal{N}_i(r_j, s_j) = \delta_{ij}$ is fulfilled. In other words, the element is not interpolatory, i.e., there is no δ_{ij} property.

With regards to the 'element mass balance', W.B. states : “the mass conservation requires that the function that is constant 1 on one cell and zero on all other cells is part of the function space. That is indeed true – it’s the \mathcal{N}_4 function. Indeed, that’s the purpose of the enrichment with the P_0 part. It is not necessary that *all* shape functions are discontinuous.”

Boffi, Cavallini, Gardini, and Gastaldi [110] (2012) state: “[...] the pressure space Q_h is defined as the sum of two finite element spaces, namely $P_k + P_0$ ($k \geq d-1$) [...] for the enhanced Hood–Taylor [...]. However, it can be easily observed that the sum is not direct, since globally constant functions can be represented exactly by means of piecewise P_0 or continuous P_k ($k \geq 1$) elements. Concerning the implementation of the method, we avoid the computation of the basis functions of such a finite element by testing the discrete problem (2.3) with the basis functions of the two subspaces separately. By the above discussion it turns out that the resulting matrix is rank-deficient, with kernel of dimension 1.”

7.3.19 The $Q_2 \times (Q_1 + Q_0)$ pair

It is a rather peculiar element pair (triplet?). The velocity space is the standard Q_2 space but the pressure space is the sum of two spaces, i.e. Q_1 and Q_0 . Please see Section 7.3.18 on the $P_2 \times (P_1 + P_0)$ element.



Taken from Gresho and Sani [488]’s book.

It is implemented in [STONE](#) 120.

7.3.20 The $P_3 \times P_2$ pair

pair_p3p2.tex

$P_3 \times P_2$ mentioned in Stenberg [1207]. The P_3 basis functions are presented in Section 5.3.12 and the P_2 basis functions in Section 5.3.10. See [STONE](#) ??.

7.3.21 The Raviart-Thomas family

- Raviart Thomas 0 RT0 [1048] ? mentioned/defined/drawn in 4.2.2 of Kanschat book. Also exist for quads see 4.2.37 Hanert, Legat, and Deleersnijder [526]: “ $P_1^\perp \times P_0$ symbol denotes an element with normal velocity nodes in the middle of each edge of the triangulation [...]. This element, also called low order Raviart–Thomas element (Raviart and Thomas, 1977), is based on flux conservation on elements edges and the resulting scheme is very close to a finite volume scheme.”

Mentioned in John [650], appendix B.3, example B.45: “the normal component of \mathbf{v} on each face is a constant. The normal component of functions from RT0 is continuous across faces of the mesh cells.”

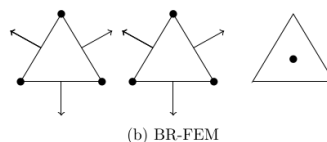
Check Brezzi and Fortin [148]

Mentioned in Chen [223] (1993).

7.3.22 The Bernaudi-Raugel pair

pair_bernardi-raugel.tex

In Carstensen, Köhler, Peterseim, and Schedensack [212] (2015) we find: “The BR-FEM after Bernardi and Raugel [83] is a modification of the $P_2 \times P_0$ FEM. It is sometimes also called reduced $P_2 \times P_0$ FEM”. They also state that this element also exists in 3D.



It is also mentioned in Boffi, Brezzi, and Fortin [109] although it seems it is there called the SMALL element (p474).

In Lederer: ”Consider the case $d = 2$. [...] we only need to control the normal velocity at the edge, i.e. adding the edge bubble for both components of the velocity seems to be sub optimal (with

respect to computational costs and the expected approximation properties). The idea now is to only add the normal edge bubble.”

According to John, Linke, Merdon, Neilan, and Rebholz [655] (2017) (example 6.3), “the velocity space in the Bernardi-Raugel element consists of P_1 functions which are enriched with edge bubble functions”. The authors also speak of ‘reconstructing the test functions’ and state: “the results of the method with reconstruction are generally more accurate. In summary, the use of an appropriately reconstructed test function in the Bernardi–Raugel pair of spaces led to a clear improvement of the accuracy of the computed results compared with the standard method.”

7.3.23 The Scott-Vogelius pair

It originates in Scott and Vogelius [1146] (1985).

Example 3.73 ($P_k/P_{k-1}^{\text{disc}}$, $k \geq 2$, on a Special Macro Cell) Consider the pair of spaces $P_k/P_{k-1}^{\text{disc}}$, $k \geq 2$, on the grid shown in Fig. 3.5, in particular in the macro cell which is surrounded boldly. The diagonals of the macro cell should be parallel to the coordinate axes. The pair $P_k/P_{k-1}^{\text{disc}}$, $k \geq 2$, is called Scott–Vogelius finite element, see Scott and Vogelius (1985).

Let $\mathbf{v}^h \in P_k$, then $\nabla \cdot \mathbf{v}^h \in P_{k-1}^{\text{disc}} = Q^h$. Hence, Remark 3.56 gives that discretely divergence-free finite element velocities from $P_k/P_{k-1}^{\text{disc}}$ are even weakly divergence-free, which is a desirable property. However, it turns out that $P_k/P_{k-1}^{\text{disc}}$ does not fulfill the discrete inf-sup condition on the mesh presented in Fig. 3.5.

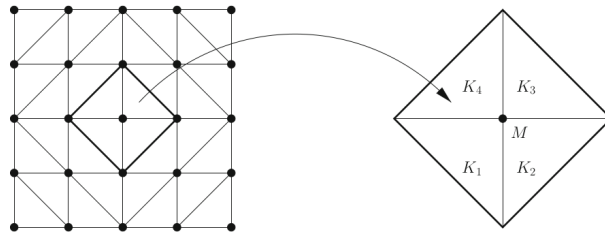


Fig. 3.5 Grid for $P_k/P_{k-1}^{\text{disc}}$, $k \geq 2$, finite element

Taken from John [650, p70].

See also John, Linke, Merdon, Neilan, and Rebholz [655] (2017).

Chen [ref?!] says: (P^k, P^{k-1}) : stable if $k \geq 4$ in R^2 and for meshes without singular-vertex. Exact divergence free. Not easy to code due to the high degree.

7.3.24 The BDM (Brezzi-Douglas-Marini) pair

BDM (Brezzi-Douglas-Marini) element mentioned in Kanschat book, section 4.2.14. Also exist for quads see section 4.2.39. Mentioned in Chen [223] (1993), Also check Brezzi and Fortin [148]

7.3.25 The DSSY pair

This element is often referred to as the ‘DSSY’ element because of the four authors of the original paper: Douglas, Santos, sheen and Ye (1999) [345].

The non-conforming finite element space Q_l is defined based on the reference square element on $[-1, 1]^2$:

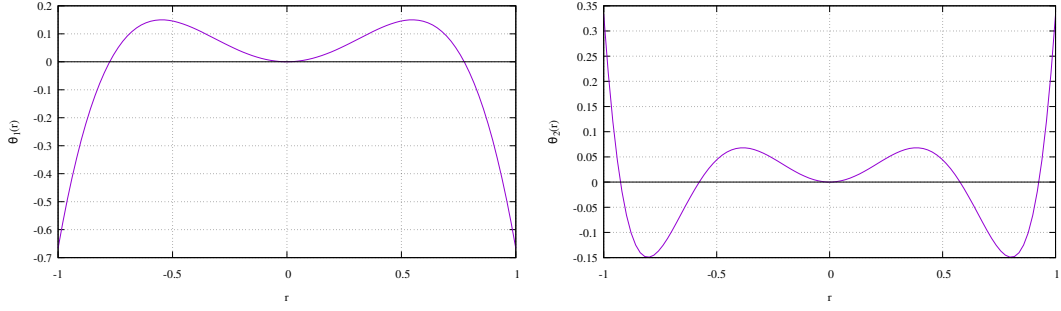
$$Q_l = \text{Span} \{1, r, s, \theta_l(r) - \theta_l(s)\} \quad l = 1, \text{ or } 2$$

with

$$\begin{aligned}\theta_1(r) &= r^2 - \frac{5}{3}r^4 \\ \theta'_1(r) &= 2r - \frac{20}{3}r^3 \\ \theta_2(r) &= r^2 - \frac{25}{6}r^4 + \frac{7}{2}r^6\end{aligned}\tag{7.13}$$

$$\theta'_2(r) = 2r - \frac{50}{3}r^3 + 21r^5\tag{7.14}$$

The dimension of Q_l is four and the θ_l functions look like:

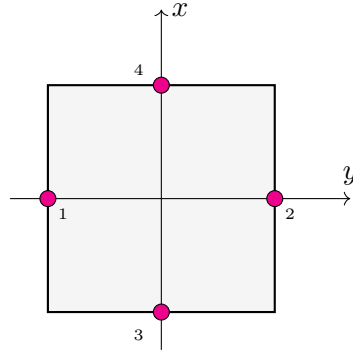


We have:

- $\theta_1(r = -1) = \theta_1(r = +1) = -\frac{2}{3}$, $\theta_1(r = 0) = 0$
- $\theta_2(r = -1) = \theta_2(r = +1) = \frac{1}{3}$, $\theta_2(r = 0) = 0$

The nodes are situated at the mid-edges of the quadrilateral:

(tikz_dssy2D.tex)



The basis function corresponding to the node (1, 0) is given by

$$\begin{aligned}\mathcal{N}_1(r, s)^{(l)} &= \frac{1}{4} - \frac{1}{2}r + \frac{\theta_l(r) - \theta_l(s)}{4\theta_l(1)} \\ \mathcal{N}_2(r, s)^{(l)} &= \frac{1}{4} + \frac{1}{2}r + \frac{\theta_l(r) - \theta_l(s)}{4\theta_l(1)} \\ \mathcal{N}_3(r, s)^{(l)} &= \frac{1}{4} - \frac{1}{2}s - \frac{\theta_l(r) - \theta_l(s)}{4\theta_l(1)} \\ \mathcal{N}_4(r, s)^{(l)} &= \frac{1}{4} + \frac{1}{2}s - \frac{\theta_l(r) - \theta_l(s)}{4\theta_l(1)}\end{aligned}\tag{7.15}$$

We can easily verify that $\sum_i \mathcal{N}_i(r, s, t) = 1$ and that $\mathcal{N}_i(\vec{r}_j) = \delta_{ij}$:

$$\begin{aligned}
\mathcal{N}_1^{(l)}(r_1, s_1) &= \frac{1}{4} - \frac{1}{2}(-1) + \frac{\theta_l(-1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} + \frac{1}{2} + \frac{\theta_l(-1)}{4\theta_l(1)} = \frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1 \\
\mathcal{N}_1^{(l)}(r_2, s_2) &= \frac{1}{4} - \frac{1}{2}(+1) + \frac{\theta_l(+1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{2} + \frac{\theta_l(+1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{2} + \frac{1}{4} = 0 \\
\mathcal{N}_1^{(l)}(r_3, s_3) &= \frac{1}{4} - \frac{1}{2}(0) + \frac{\theta_l(0) - \theta_l(-1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_1^{(l)}(r_4, s_4) &= \frac{1}{4} - \frac{1}{2}(0) + \frac{\theta_l(0) - \theta_l(+1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_2^{(l)}(r_1, s_1) &= \frac{1}{4} + \frac{1}{2}(-1) + \frac{\theta_l(-1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{2} + \frac{1}{4} = 0 \\
\mathcal{N}_2^{(l)}(r_2, s_2) &= \frac{1}{4} + \frac{1}{2}(+1) + \frac{\theta_l(+1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1 \\
\mathcal{N}_2^{(l)}(r_3, s_3) &= \frac{1}{4} + \frac{1}{2}(0) + \frac{\theta_l(0) - \theta_l(-1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_2^{(l)}(r_4, s_4) &= \frac{1}{4} + \frac{1}{2}(0) + \frac{\theta_l(0) - \theta_l(+1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_3^{(l)}(r_1, s_1) &= \frac{1}{4} - \frac{1}{2}(0) - \frac{\theta_l(-1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_3^{(l)}(r_2, s_2) &= \frac{1}{4} - \frac{1}{2}(0) - \frac{\theta_l(+1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_3^{(l)}(r_3, s_3) &= \frac{1}{4} - \frac{1}{2}(-1) - \frac{\theta_l(0) - \theta_l(-1)}{4\theta_l(1)} = \frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1 \\
\mathcal{N}_3^{(l)}(r_4, s_4) &= \frac{1}{4} - \frac{1}{2}(+1) - \frac{\theta_l(0) - \theta_l(+1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{2} + \frac{1}{4} = 0 \\
\mathcal{N}_4^{(l)}(r_1, s_1) &= \frac{1}{4} + \frac{1}{2}(0) - \frac{\theta_l(-1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_4^{(l)}(r_2, s_2) &= \frac{1}{4} + \frac{1}{2}(0) - \frac{\theta_l(+1) - \theta_l(0)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{4} = 0 \\
\mathcal{N}_4^{(l)}(r_3, s_3) &= \frac{1}{4} + \frac{1}{2}(-1) - \frac{\theta_l(0) - \theta_l(-1)}{4\theta_l(1)} = \frac{1}{4} - \frac{1}{2} + \frac{1}{4} = 0 \\
\mathcal{N}_4^{(l)}(r_4, s_4) &= \frac{1}{4} + \frac{1}{2}(1) - \frac{\theta_l(0) - \theta_l(1)}{4\theta_l(1)} = \frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1
\end{aligned}$$

The basis functions can also be explicitly written for θ_1 as in Cai *et al.* [203]:


$$\begin{aligned}
\mathcal{N}_1(r, s)^{(l)} &= \frac{1}{4} - \frac{1}{2}r - \frac{3}{8} \left[\left(r^2 - \frac{5}{3}r^4 \right) - \left(s^2 - \frac{5}{3}s^4 \right) \right] \\
\mathcal{N}_2(r, s)^{(l)} &= \frac{1}{4} + \frac{1}{2}r - \frac{3}{8} \left[\left(r^2 - \frac{5}{3}r^4 \right) - \left(s^2 - \frac{5}{3}s^4 \right) \right] \\
\mathcal{N}_3(r, s)^{(l)} &= \frac{1}{4} - \frac{1}{2}s + \frac{3}{8} \left[\left(r^2 - \frac{5}{3}r^4 \right) - \left(s^2 - \frac{5}{3}s^4 \right) \right] \\
\mathcal{N}_4(r, s)^{(l)} &= \frac{1}{4} + \frac{1}{2}s + \frac{3}{8} \left[\left(r^2 - \frac{5}{3}r^4 \right) - \left(s^2 - \frac{5}{3}s^4 \right) \right]
\end{aligned} \tag{7.16}$$

The derivatives of the basis functions are as follows:

$$\begin{aligned}
\partial_r \mathcal{N}_1(r, s)^{(l)} &= -\frac{1}{2} + \frac{\theta'_l(r)}{4\theta_l(1)} \\
\partial_r \mathcal{N}_2(r, s)^{(l)} &= +\frac{1}{2} + \frac{\theta'_l(r)}{4\theta_l(1)} \\
\partial_r \mathcal{N}_3(r, s)^{(l)} &= -\frac{\theta'_l(r)}{4\theta_l(1)} \\
\partial_r \mathcal{N}_4(r, s)^{(l)} &= -\frac{\theta'_l(r)}{4\theta_l(1)}
\end{aligned} \tag{7.17}$$

$$\begin{aligned}
\partial_s \mathcal{N}_1(r, s)^{(l)} &= -\frac{\theta'_l(s)}{4\theta_l(1)} \\
\partial_s \mathcal{N}_2(r, s)^{(l)} &= -\frac{\theta'_l(s)}{4\theta_l(1)} \\
\partial_s \mathcal{N}_3(r, s)^{(l)} &= -\frac{1}{2} + \frac{\theta'_l(s)}{4\theta_l(1)} \\
\partial_s \mathcal{N}_4(r, s)^{(l)} &= +\frac{1}{2} + \frac{\theta'_l(s)}{4\theta_l(1)}
\end{aligned} \tag{7.18}$$

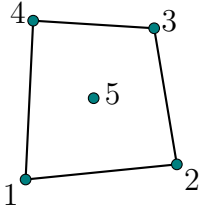
Note that a correction was issued in Cai, Douglas Jr, Santos, Sheen, and Ye [202] (2000) if a true quadrilateral (i.e., one having two opposite, nonparallel edges) is included in the partition. The authors state that in the case of rectangles the original method is fine.

 **Relevant Literature:** Park & Sheen (2003) [974], Jeon *et al.* (2013) [645], Park, Sheen & Shin (2013) [975], Bangerth *et al.* (2017) [43], Sheen (2020) [1155]

7.3.26 The Han pair

It is based on Han [524] (also mentioned in Sheen (2020) [1155]). The nodes are at the same location as for the RT element above, but there is an additional bubble function in the middle:

(tikz_han.tex)



Inside the reference element we assume that a field f can be represented by

$$f^h(r, s) = a + br + cs + d \underbrace{\frac{5r^4 - 3r^2}{2}}_{\phi(r)} + e \underbrace{\frac{5s^4 - 3s^2}{2}}_{\phi(s)}$$

We then must have

$$\begin{aligned}
f_1 &= f^h(r = 1, s = 0) &= a + b + d \\
f_2 &= f^h(r = 0, s = 1) &= a + c + e \\
f_3 &= f^h(r = -1, s = 0) &= a - b + d \\
f_4 &= f^h(r = 0, s = -1) &= a - c + e \\
f_5 &= f^h(r = 0, s = 0) &= a
\end{aligned}$$

and we easily get

$$a = f_5 \quad f_1 - f_3 = 2b \quad f_2 - f_4 = 2c$$

followed by

$$d = f_1 - a - b = f_1 - f_5 - \frac{1}{2}(f_1 - f_3) = \frac{f_1 - 2f_5 + f_3}{2}$$

and

$$e = f_2 - a - c = f_2 - f_5 - \frac{1}{2}(f_2 - f_4) = \frac{f_2 - 2f_5 + f_4}{2}$$

Finally:

$$f(r, s) = f_5 + \frac{1}{2}(f_1 - f_3)r + \frac{1}{2}(f_2 - f_4)s + \frac{f_1 - 2f_5 + f_3}{2}\phi(r) + \frac{f_2 - 2f_5 + f_4}{2}\phi(s)$$

i.e.

$$f(r, s) = \left(\frac{r + \phi(r)}{2}\right) f_1 + \left(\frac{s + \phi(s)}{2}\right) f_2 + \left(-\frac{r - \phi(r)}{2}\right) f_3 + \left(-\frac{s - \phi(s)}{2}\right) f_4 + (1 - \phi(r) - \phi(s)) f_5$$

which has us define

$$\begin{aligned} \mathcal{N}_1(r, s) &= \frac{r + \phi(r)}{2} \\ \mathcal{N}_2(r, s) &= \frac{s + \phi(s)}{2} \\ \mathcal{N}_3(r, s) &= -\frac{r - \phi(r)}{2} \\ \mathcal{N}_4(r, s) &= -\frac{s - \phi(s)}{2} \\ \mathcal{N}_5(r, s) &= 1 - \phi(r) - \phi(s) \end{aligned}$$

We have of course the following properties $\sum_{i=1}^5 \mathcal{N}_i(r, s) = 1$ and $\mathcal{N}_i(r_j, s_j) = \delta_{ij}$, $i, j \in 1, 5$. The partial derivatives of the basis functions are as follows

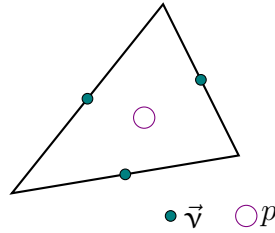
$$\begin{aligned} \partial_r \mathcal{N}_1(r, s) &= \frac{1 + \phi'(r)}{2} \\ \partial_r \mathcal{N}_2(r, s) &= 0 \\ \partial_r \mathcal{N}_3(r, s) &= -\frac{1 - \phi'(r)}{2} \\ \partial_r \mathcal{N}_4(r, s) &= 0 \\ \partial_r \mathcal{N}_5(r, s) &= -\phi'(r) \\ \partial_s \mathcal{N}_1(r, s) &= 0 \\ \partial_s \mathcal{N}_2(r, s) &= \frac{1 + \phi'(s)}{2} \\ \partial_s \mathcal{N}_3(r, s) &= 0 \\ \partial_s \mathcal{N}_4(r, s) &= -\frac{1 - \phi'(s)}{2} \\ \partial_s \mathcal{N}_5(r, s) &= -\phi'(s) \end{aligned}$$

This element is implemented in the `stone_han.py` file in [STONE](#) 77 and also in [STONE](#) 120.

7.3.27 The Divergence-free nonconforming $P_1^{NC} \times P_0$ pair

It belongs to the Crouzeix-Raviart family. The midside nodes are used as degrees of freedom for the velocities. It is mentioned in Section 6.3 of Boffi, Brezzi, and Fortin [108] (2008): “[...] It is exactly divergence free. Another important feature of this element is that it can be seen as a “mass conservation” scheme. The present element has been generalized to second order in Fortin and Soulie [404] (1983). It must also be said that coerciveness may be a problem for the $P_1^{NC} \times P_0$ element, as it does not satisfy the discrete version of Korn’s inequality. This issue has been deeply investigated and clearly illustrated in Arnold [28] (1993).”

(tikz_p1ncp0.tex)

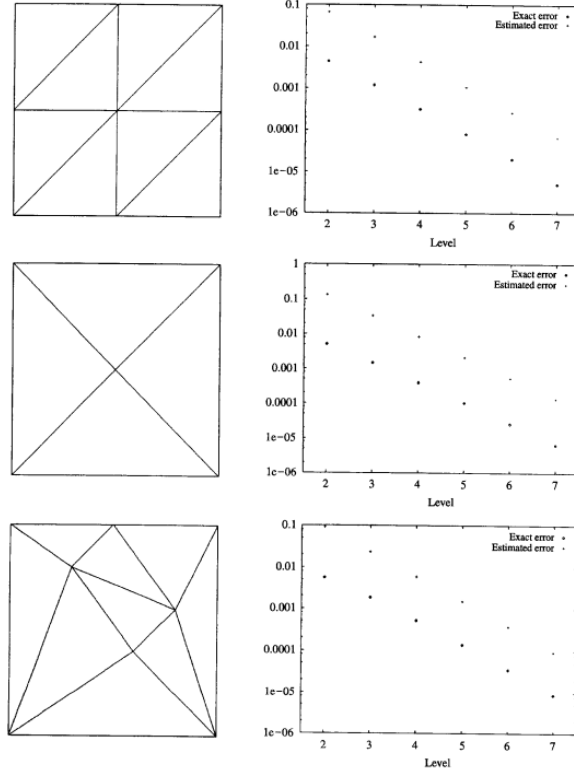


At page 170 of [128] it is stated that “an analogous quadrilateral element was developed and studied by Rannacher and Turek [1045] (1992)”.

In Boffi, Brezzi, and Fortin [109] we find: “We consider the classical (almost¹³) stable nonconforming triangular element introduced in Crouzeix and Raviart [290], in which mid-side nodes are used as degrees of freedom for the velocities. This generates a piecewise linear nonconforming approximation; pressures are taken constant on each element. It is also possible to build a three-dimensional version of this element, using mid-face nodes as degrees of freedom.” Also: “It must also be recalled that coercivity is a problem for the $P_1^{NC} \times P_0$ element. The trouble is that the bilinear form (8.2.1) is not coercive on the nonconforming space V_h and we do not have the discrete version of Korn’s inequality.”

It is also mentioned in John [650], appendix B.3, example B.43, in 2D and 3D, in Brezzi and Fortin [148] (example 8.1), and studied extensively in John [649] (1998).

¹³What does that mean?!



Taken from John [649].

In John, Linke, Merdon, Neilan, and Rebholz [655] (2017) the authors show results obtained with this element (fig 6) but also explain that these are obtained with so-called reconstructed test functions.

7.3.28 The Chen nonconforming $Q_1 \times Q_0$ pair (?)

pair_chen.tex

What follows is tentative!

This space is proposed in Chen [225] (1993), albeit not in the context of the Stokes equations. It is based on the mid-point variant of the RT basis functions,

$$\begin{aligned}\mathcal{N}_1(r, s) &= \frac{1}{4}(1 - 2s - (r^2 - s^2)) \\ \mathcal{N}_2(r, s) &= \frac{1}{4}(1 + 2r + (r^2 - s^2)) \\ \mathcal{N}_3(r, s) &= \frac{1}{4}(1 + 2s - (r^2 - s^2)) \\ \mathcal{N}_4(r, s) &= \frac{1}{4}(1 - 2r + (r^2 - s^2))\end{aligned}$$

to which a P_2 bubble is added

$$\phi(r, s) = 1 - \frac{3}{4}(r^2 + s^2)$$

Note that this function is zero at locations $\pm 1/\sqrt{3}$ on all four edges and exactly 1 in the middle.

A field f is represented inside the element by

$$f^h(r, s) = a\mathcal{N}_1(r, s) + b\mathcal{N}_2(r, s) + c\mathcal{N}_3(r, s) + d\mathcal{N}_4(r, s) + e\phi(r, s)$$

We immediately see that this space is not interpolatory, i.e. the basis function $\phi(r, s)$ cannot be 1 in the middle and 0 at the other four nodes.

Chen [224] also extends this to 3D in the paper.

This space is used for velocity and a Q_0 space is used for pressure in [STONE](#) 120 (only because the basis functions above are based on the Rannacher-Turek ones).

7.3.29 Other FE element pairs

- $Q_2 \times Q_2$: This element is never used, probably because a) it is unstable, b) it is very costly. There is one reference to it in Hughes, Franca, and Balestra [606] (1986).
- $Q_1 \times P_{-1}$ Bilinear velocities, piecewise linear discontinuous polynomial pressure.
- See Fortin [401] for various stable low order elements other than the enriched $Q_1^+ \times P_0$
- $Q_1 \times Q_1$ + nonconforming null edge average [411]
- check Dhatt and Hubert [332] (1986) many flavours of triangles and quads.
- Bercovier-Pironneau element pair, or P_1isoP_2 . See Boffi, Cavallini, Gardini, and Gastaldi [110] (2012).

7.3.30 A note about incompressibility and standard mixed methods

What follows is nicely explained and demonstrated in John *et al.* [655]. In their example 1.1 they look at the velocity error of benchmark VJ2 (see Section 12.1.9) which analytical solution is a zero velocity field. They show that for the MINI, Taylor-Hood and Crouzeix-Raviart triangular elements the velocity error grows with the magnitude of the rhs. They also make this statement: “there are important applications, e.g., natural convection problems, where the pressure is larger than the velocity by orders of magnitude. In such situations, one cannot expect to compute accurate velocity fields with classical mixed methods, at least for low order methods.”

7.4 The penalty approach for viscous flow

penalty.tex

In order to impose the incompressibility constraint, two widely used procedures are available, namely the Lagrange multiplier method and the penalty method [53, 604]. The latter allows for the elimination of the pressure variable from the momentum equation (resulting in a reduction of the matrix size).

Mathematical details on the origin and validity of the penalty approach applied to the Stokes problem can for instance be found in Cuvelier *et al.* [298], Reddy [1050] or Gunzburger [507].

The penalty formulation of the mass conservation equation is based on a relaxation of the incompressibility constraint and writes

$$\vec{\nabla} \cdot \vec{v} + \frac{p}{\lambda} = 0 \quad (7.19)$$

where λ is the penalty parameter, that can be interpreted (and has the same dimension) as a bulk viscosity. It is equivalent to say that the material is weakly compressible. It can be shown that if one chooses λ to be a sufficiently large number, the continuity equation $\vec{\nabla} \cdot \vec{v} = 0$ will be approximately satisfied in the finite element solution. The value of λ is often recommended to be 6 to 7 orders of magnitude larger than the shear viscosity [341, 608].

Equation (7.19) can be used to eliminate the pressure in the momentum equation so that the mass and momentum conservation equations fuse to become :

$$\vec{\nabla} \cdot (2\eta\dot{\epsilon}(\vec{v})) + \lambda\vec{\nabla}(\vec{\nabla} \cdot \vec{v}) + \rho\vec{g} = \vec{0} \quad (7.20)$$

Malkus & Hughes (1978) [830] have established the equivalence for incompressible problems between the reduced integration of the penalty term and a mixed Finite Element approach if the pressure nodes coincide with the integration points of the reduced rule.

In the end, the elimination of the pressure unknown in the Stokes equations replaces the original saddle-point Stokes problem [72] by an elliptical problem, which leads to a symmetric positive definite (SPD) FEM matrix. This is the major benefit of the penalized approach over the full indefinite solver with the velocity-pressure variables. Indeed, the SPD character of the matrix lends itself to efficient solving strategies and is less memory-demanding since it is sufficient to store only the upper half of the matrix including the diagonal [473].

The penalty approach for example is used in the SOPALE , DOUAR , ConMan, FANTOM and ELEFANT geodynamical codes.

Remark. *FEM codes relying on the penalty approach all rely on direct solvers, because as explained in Brezzi & Fortin [148]: "Using a penalty method is, for instance, almost impossible if an iterative method is used for the solution of the linear system, iterative methods being in general quite sensitive to the condition number of the matrix at hand."*

Since the penalty formulation is only valid for incompressible flows, then $\dot{\epsilon}(\vec{v}) = \dot{\epsilon}^d(\vec{v})$ so that the d superscript is omitted in what follows. We here focus on Cartesian coordinates only and because

the stress tensor is symmetric one can also rewrite it the following vector format:

$$\begin{aligned}
\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} &= \begin{pmatrix} -p \\ -p \\ -p \\ 0 \\ 0 \\ 0 \end{pmatrix} + 2\eta \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yz} \end{pmatrix} \\
&= \lambda \begin{pmatrix} \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ 0 \\ 0 \\ 0 \end{pmatrix} + 2\eta \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yz} \end{pmatrix} \\
&= \left[\lambda \underbrace{\begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}}_K + \eta \underbrace{\begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_C \right] \cdot \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix}
\end{aligned}$$

Remember that

$$\frac{\partial u^h}{\partial x} = \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial x} u_i \quad \frac{\partial v^h}{\partial y} = \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial y} v_i \quad \frac{\partial w^h}{\partial z} = \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial z} w_i$$

and

$$\begin{aligned}
\frac{\partial u^h}{\partial y} + \frac{\partial v^h}{\partial x} &= \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial y} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial x} v_i \\
\frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial x} &= \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial z} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial x} w_i \\
\frac{\partial v^h}{\partial z} + \frac{\partial w^h}{\partial y} &= \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial z} v_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i}{\partial y} w_i
\end{aligned}$$

so that, since in $m_v = 8$ in 3D:

$$\begin{pmatrix} \frac{\partial u^h}{\partial x} \\ \frac{\partial w^h}{\partial y} \\ \frac{\partial w^h}{\partial z} \\ \frac{\partial u^h}{\partial y} + \frac{\partial v^h}{\partial x} \\ \frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial x} \\ \frac{\partial v^h}{\partial z} + \frac{\partial w^h}{\partial y} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \mathcal{N}_1}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_2}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_3}{\partial x} & 0 & 0 & \dots & \frac{\partial \mathcal{N}_8}{\partial x} & 0 & 0 \\ 0 & \frac{\partial \mathcal{N}_1}{\partial y} & 0 & 0 & \frac{\partial \mathcal{N}_2}{\partial y} & 0 & 0 & \frac{\partial \mathcal{N}_3}{\partial y} & 0 & \dots & 0 & \frac{\partial \mathcal{N}_8}{\partial y} & 0 \\ 0 & 0 & \frac{\partial \mathcal{N}_1}{\partial z} & 0 & 0 & \frac{\partial \mathcal{N}_2}{\partial z} & 0 & 0 & \frac{\partial \mathcal{N}_3}{\partial z} & \dots & 0 & 0 & \frac{\partial \mathcal{N}_8}{\partial z} \\ \frac{\partial \mathcal{N}_1}{\partial y} & \frac{\partial \mathcal{N}_1}{\partial x} & 0 & \frac{\partial \mathcal{N}_2}{\partial y} & \frac{\partial \mathcal{N}_2}{\partial x} & 0 & \frac{\partial \mathcal{N}_3}{\partial y} & \frac{\partial \mathcal{N}_3}{\partial x} & 0 & \dots & \frac{\partial \mathcal{N}_8}{\partial y} & \frac{\partial \mathcal{N}_8}{\partial x} & 0 \\ \frac{\partial \mathcal{N}_1}{\partial z} & 0 & \frac{\partial \mathcal{N}_1}{\partial x} & \frac{\partial \mathcal{N}_2}{\partial z} & 0 & \frac{\partial \mathcal{N}_2}{\partial x} & \frac{\partial \mathcal{N}_3}{\partial z} & 0 & \frac{\partial \mathcal{N}_3}{\partial x} & \dots & \frac{\partial \mathcal{N}_8}{\partial z} & 0 & \frac{\partial \mathcal{N}_8}{\partial x} \\ 0 & \frac{\partial \mathcal{N}_1}{\partial z} & \frac{\partial \mathcal{N}_1}{\partial y} & 0 & \frac{\partial \mathcal{N}_2}{\partial z} & \frac{\partial \mathcal{N}_2}{\partial y} & 0 & \frac{\partial \mathcal{N}_3}{\partial z} & \frac{\partial \mathcal{N}_3}{\partial y} & \dots & 0 & \frac{\partial \mathcal{N}_8}{\partial z} & \frac{\partial \mathcal{N}_8}{\partial y} \end{pmatrix}}_{\mathbf{B}(6 \times 24)} \cdot \underbrace{\begin{pmatrix} u1 \\ v1 \\ w1 \\ u2 \\ v2 \\ w2 \\ u3 \\ v3 \\ w3 \\ \dots \\ u8 \\ v8 \\ w8 \end{pmatrix}}_{\vec{V}(24 \times 1)}.$$

Finally,

$$\vec{\sigma} = \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} = (\lambda \mathbf{K} + \eta \mathbf{C}) \cdot \mathbf{B} \cdot \vec{V}$$

We will now establish the weak form of the momentum conservation equation. We start again from

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{b} = \vec{0}$$

For the \mathcal{N}_i 's 'regular enough', we can write:

$$\int_{\Omega_e} \mathcal{N}_i \vec{\nabla} \cdot \boldsymbol{\sigma} dV + \int_{\Omega_e} \mathcal{N}_i \vec{b} dV = 0$$

We can integrate by parts and drop the surface term¹⁴:

$$\int_{\Omega_e} \vec{\nabla} \mathcal{N}_i \cdot \boldsymbol{\sigma} dV = \int_{\Omega_e} \mathcal{N}_i \vec{b} dV$$

or,

$$\int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial y} & \frac{\partial \mathcal{N}_i}{\partial z} & 0 \\ 0 & \frac{\partial \mathcal{N}_i}{\partial y} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & 0 & \frac{\partial \mathcal{N}_i}{\partial z} \\ 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial z} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & \frac{\partial \mathcal{N}_i}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \mathcal{N}_i \vec{b} dV$$

¹⁴We will come back to this at a later stage

Let $i = 1, 2, 3, 4, \dots, 8$ and stack the resulting eight equations on top of one another.

$$\begin{aligned}
& \int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial y} & \frac{\partial \mathcal{N}_i}{\partial z} & 0 \\ 0 & \frac{\partial \mathcal{N}_i}{\partial y} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & 0 & \frac{\partial \mathcal{N}_i}{\partial z} \\ 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial z} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & \frac{\partial \mathcal{N}_i}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \mathcal{N}_i \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} dV \\
& \int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial y} & \frac{\partial \mathcal{N}_i}{\partial z} & 0 \\ 0 & \frac{\partial \mathcal{N}_i}{\partial y} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & 0 & \frac{\partial \mathcal{N}_i}{\partial z} \\ 0 & 0 & \frac{\partial \mathcal{N}_i}{\partial z} & 0 & \frac{\partial \mathcal{N}_i}{\partial x} & \frac{\partial \mathcal{N}_i}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \mathcal{N}_2 \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} dV \\
& \dots \\
& \int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_8}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_8}{\partial y} & \frac{\partial \mathcal{N}_8}{\partial z} & 0 \\ 0 & \frac{\partial \mathcal{N}_8}{\partial y} & 0 & \frac{\partial \mathcal{N}_8}{\partial x} & 0 & \frac{\partial \mathcal{N}_8}{\partial z} \\ 0 & 0 & \frac{\partial \mathcal{N}_8}{\partial z} & 0 & \frac{\partial \mathcal{N}_8}{\partial x} & \frac{\partial \mathcal{N}_8}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \mathcal{N}_8 \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} dV \quad (7.21)
\end{aligned}$$

We easily recognize \mathbf{B}^T inside the integrals! Let us define

$$\vec{\mathcal{N}}_b^T = (\mathcal{N}_1 b_x, \mathcal{N}_1 b_y, \mathcal{N}_1 b_z, \dots, \mathcal{N}_8 b_x, \mathcal{N}_8 b_y, \mathcal{N}_8 b_z)$$

then we can write

$$\int_{\Omega_e} \mathbf{B}^T \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \vec{\mathcal{N}}_b dV$$

and finally:

$$\int_{\Omega_e} \mathbf{B}^T \cdot [\lambda \mathbf{K} + \eta \mathbf{C}] \cdot \mathbf{B} \cdot \vec{V} dV = \int_{\Omega_e} \vec{\mathcal{N}}_b dV$$

Since \vec{V} is the vector of unknowns (i.e. the velocities at the corners), it does not depend on the x , y or z coordinates so it can be taken outside of the integral (remember $dV = dx dy dz$ here):

$$\underbrace{\left(\int_{\Omega_e} \mathbf{B}^T \cdot [\lambda \mathbf{K} + \eta \mathbf{C}] \cdot \mathbf{B} dV \right)}_{\mathbf{A}_{el}(24 \times 24)} \cdot \underbrace{\vec{V}}_{(24 \times 1)} = \underbrace{\int_{\Omega_e} \vec{\mathcal{N}}_b dV}_{\vec{B}_{el}(24 \times 1)}$$

or,

$$\left[\underbrace{\left(\int_{\Omega_e} \lambda \mathbf{B}^T \cdot \mathbf{K} \cdot \mathbf{B} dV \right)}_{\mathbf{A}_{el}^\lambda(24 \times 24)} + \underbrace{\left(\int_{\Omega_e} \eta \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} dV \right)}_{\mathbf{A}_{el}^\eta(24 \times 24)} \right] \cdot \underbrace{\vec{V}}_{(24 \times 1)} = \underbrace{\int_{\Omega_e} \vec{\mathcal{N}}_b dV}_{\vec{B}_{el}(24 \times 1)}$$

Once the elemental matrix and rhs have been computed for an element its contribution is added to the global FEM matrix. The linear system is then solved and the velocity field at all nodes is obtained. From this velocity field the elemental pressure can be recovered by means of Eq. (7.19):

$$p = -\lambda \vec{\nabla} \cdot \vec{v}.$$

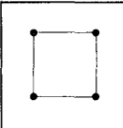
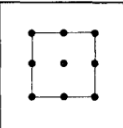
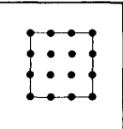
In two dimensions the equations above are very similar. Let us assume that the flow is taking place in the xy -plane and that the domain is infinite in the z -direction. Then $w = 0$ and $\partial_z \rightarrow 0$. From the 6 terms of the strain rate tensor only three remain: $\dot{\epsilon}_{xx}$, $\dot{\epsilon}_{yy}$ and $\dot{\epsilon}_{xy}$. Then, since $m_v = 4$ (each element has 4 velocity nodes), we have

$$\mathbf{K} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \mathbf{C} = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Conversely the matrix \mathbf{B} has size 3×8 , etc ... Note that we will come across matrix \mathbf{C} again when we solve the (non penalty-formulated) Stokes equations in the following sections.

As stated before the implementation is rather straightforward since only one FE matrix must be computed and assembled. However, there is one specific point which needs to be addressed: reduced integration.


To quote Hughes *et al.* (1979) [608]: *When a quadrature rule of lower order than the “standard” one is employed, this is called reduced integration. If all terms employ the same reduced integration, this is called uniform reduced integration; if reduced integration is used on some terms while standard integration is used on others, this is called selective reduced integration.*

| | | | |
|-----------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------|
| |  |  |  |
| shape functions | bilinear | biquadratic | bicubic |
| λ -term | 1 point | 2×2 | 3×3 |
| μ -term | 2×2 | 3×3 | 4×4 |

Selective Gauss-Legendre integration rules for 2-dimensional isoparametric Lagrange elements. Taken from Hughes *et al.* (1979) [608].

In the context of penalty-based codes, it has been shown [830] that it is crucial to resort to a selective reduced integration approach. The viscosity term \mathbf{A}_{el}^η is integrated on 2^{ndim} points but the penalty term \mathbf{A}_{el}^λ must be integrated on a single quadrature point. Finally, when the pressure is computed from the velocity via Eq. (7.19), the divergence term must also be computed at a single point in the middle of the element.

See [STONE 01](#) for a concrete example of a 2D penalty-based Stokes solver.

 **Relevant Literature:** Oden *et al.* [951], Dhatt & Hubert[332].

7.5 The mixed FEM for viscous flow

mixed.tex

7.5.1 In three dimensions

The FEM formulation of the Stokes equation is quite complex so we simplify things as much as possible for now by assuming the flow to be incompressible, isoviscous and isothermal.

The methodology to derive the discretised equations of the mixed system is quite similar to the one we have used in the case of the penalty formulation. The big difference comes from the fact that we are now solving for both velocity and pressure at the same time, and that we therefore must solve the mass and momentum conservation equations together. As before, velocity inside an element is given by

$$\vec{v}^h(\vec{r}) = \sum_{i=1}^{m_v} \mathcal{N}_i^v(\vec{r}) \vec{v}_i \quad (7.22)$$

where \mathcal{N}_i^v are the polynomial basis functions for the velocity, and the summation runs over the m_v velocity nodes composing the element. A similar expression is used for pressure:

$$p^h(\vec{r}) = \sum_{i=1}^{m_p} \mathcal{N}_i^p(\vec{r}) p_i \quad (7.23)$$

Note that the velocity is a vector while pressure (and temperature) is a scalar. There are then $ndof_v = ndim$ velocity degrees of freedom per node and $ndof_p = 1$ pressure degrees of freedom. It is also very important to remember that the numbers of velocity nodes and pressure nodes for a given element are more often than not different and that velocity and pressure nodes need not be colocated. Indeed, unless so-called 'stabilised elements' are used, we have $m_v > m_p$, which means that the polynomial order of the velocity field is higher than the polynomial order of the pressure field (usually by value 1).

Other notations will be sometimes used for Eqs. (7.22) and (7.23):

$$u^h(\vec{r}) = \vec{\mathcal{N}}^v \cdot \vec{u} \quad v^h(\vec{r}) = \vec{\mathcal{N}}^v \cdot \vec{v} \quad w^h(\vec{r}) = \vec{\mathcal{N}}^v \cdot \vec{w} \quad p^h(\vec{r}) = \vec{\mathcal{N}}^p \cdot \vec{p} \quad (7.24)$$

where $\vec{v} = (u, v, w)$ and $\vec{\mathcal{N}}^v$ is the vector containing all basis functions evaluated at location \vec{r} :

$$\vec{\mathcal{N}}^v = (\mathcal{N}_1^v(\vec{r}), \mathcal{N}_2^v(\vec{r}), \mathcal{N}_3^v(\vec{r}), \dots, \mathcal{N}_{m_v}^v(\vec{r})) \quad (7.25)$$

$$\vec{\mathcal{N}}^p = (\mathcal{N}_1^p(\vec{r}), \mathcal{N}_2^p(\vec{r}), \mathcal{N}_3^p(\vec{r}), \dots, \mathcal{N}_{m_p}^p(\vec{r})) \quad (7.26)$$

and with

$$\vec{u} = (u_1, u_2, u_3, \dots, u_{m_v}) \quad (7.27)$$

$$\vec{v} = (v_1, v_2, v_3, \dots, v_{m_v}) \quad (7.28)$$

$$\vec{w} = (w_1, w_2, w_3, \dots, w_{m_v}) \quad (7.29)$$

$$\vec{p} = (p_1, p_2, p_3, \dots, p_{m_p}) \quad (7.30)$$

We will now establish the weak form of the momentum conservation equation. We start again from

$$\vec{\nabla} \cdot \vec{\sigma} + \vec{b} = \vec{0} \quad (7.31)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (7.32)$$

For the \mathcal{N}_i^γ 's and \mathcal{N}_i^p 'regular enough', we can write:

$$\int_{\Omega_e} \mathcal{N}_i^\gamma \vec{\nabla} \cdot \boldsymbol{\sigma} dV + \int_{\Omega_e} \mathcal{N}_i^\gamma \vec{b} dV = \vec{0} \quad (7.33)$$

$$\int_{\Omega_e} \mathcal{N}_i^p \vec{\nabla} \cdot \vec{v} dV = 0 \quad (7.34)$$

We can integrate by parts and drop the surface term¹⁵:

$$\int_{\Omega_e} \vec{\nabla} \mathcal{N}_i^\gamma \cdot \boldsymbol{\sigma} dV = \int_{\Omega_e} \mathcal{N}_i^\gamma \vec{b} dV \quad (7.35)$$

$$\int_{\Omega_e} \mathcal{N}_i^p \vec{\nabla} \cdot \vec{v} dV = 0 \quad (7.36)$$

or,

$$\int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_i^\gamma}{\partial x} & 0 & 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial y} & \frac{\partial \mathcal{N}_i^\gamma}{\partial z} & 0 \\ 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial y} & 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial x} & 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial z} \\ 0 & 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial z} & 0 & \frac{\partial \mathcal{N}_i^\gamma}{\partial x} & \frac{\partial \mathcal{N}_i^\gamma}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega = \int_{\Omega_e} \mathcal{N}_i^\gamma \vec{b} dV \quad (7.37)$$

The above equation can ultimately be written:

$$\int_{\Omega_e} \mathbf{B}^T \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} dV = \int_{\Omega_e} \vec{N}_b dV \quad (7.38)$$

We have previously established that the strain rate vector $\vec{\dot{\epsilon}}$ is:

$$\vec{\dot{\epsilon}} = \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix} = \begin{pmatrix} \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i \\ \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i \\ \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial y} u_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial x} v_i \right) \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial z} u_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial x} w_i \right) \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial z} v_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial y} w_i \right) \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & 0 & 0 & \dots & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial x} & 0 & 0 \\ 0 & \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & 0 & \dots & 0 & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial y} & 0 \\ 0 & 0 & \frac{\partial \mathcal{N}_1^\gamma}{\partial z} & \dots & 0 & 0 & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial z} \\ \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & 0 & \dots & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial y} & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial x} & 0 \\ \frac{\partial \mathcal{N}_1^\gamma}{\partial z} & 0 & \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & \dots & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial z} & 0 & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial x} \\ 0 & \frac{\partial \mathcal{N}_1^\gamma}{\partial z} & \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & \dots & 0 & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial z} & \frac{\partial \mathcal{N}_{m_v}^\gamma}{\partial y} \end{pmatrix}}_{\mathbf{B}} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ w_1 \\ u_2 \\ v_2 \\ w_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \\ w_{m_v} \end{pmatrix}}_{\vec{\mathcal{V}}} \quad (7.39)$$

or, $\vec{\dot{\epsilon}} = \mathbf{B} \cdot \vec{\mathcal{V}}$ where \mathbf{B} is the gradient matrix and $\vec{\mathcal{V}}$ is the vector of all velocity degrees of freedom for the element. The matrix \mathbf{B} is then of size $6 \times (m_v \cdot ndof_v)$ and the vector $\vec{\mathcal{V}}$ is $m_v \cdot ndof_v$ long.

¹⁵We will come back to this at a later stage

we have

$$\sigma_{xx} = -p + 2\eta\dot{\epsilon}_{xx}^d \quad (7.40)$$

$$\sigma_{yy} = -p + 2\eta\dot{\epsilon}_{yy}^d \quad (7.41)$$

$$\sigma_{zz} = -p + 2\eta\dot{\epsilon}_{zz}^d \quad (7.42)$$

$$\sigma_{xy} = 2\eta\dot{\epsilon}_{xy}^d \quad (7.43)$$

$$\sigma_{xz} = 2\eta\dot{\epsilon}_{xz}^d \quad (7.44)$$

$$\sigma_{yz} = 2\eta\dot{\epsilon}_{yz}^d \quad (7.45)$$

Since we here only consider incompressible flow, we have $\dot{\epsilon}^d = \dot{\epsilon}$ so

$$\vec{\sigma} = - \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} p + \mathbf{C} \cdot \vec{\epsilon} = - \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \vec{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (7.46)$$

with

$$\mathbf{C} = \eta \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \vec{\epsilon} = \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} \quad (7.47)$$

Let us define matrix \mathcal{N}^p of size $6 \times m_p$:

$$\mathcal{N}^p = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \vec{N}^p = \begin{pmatrix} \vec{N}^p \\ \vec{N}^p \\ \vec{N}^p \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (7.48)$$

so that

$$\vec{\sigma} = -\mathcal{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (7.49)$$

finally

$$\int_{\Omega_e} \mathbf{B}^T \cdot [-\mathcal{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V}] d\Omega = \int_{\Omega_e} \mathcal{N}_b d\Omega \quad (7.50)$$

or,

$$\underbrace{\left(- \int_{\Omega_e} \mathbf{B}^T \cdot \mathcal{N}^p d\Omega \right)}_{\mathbb{G}} \cdot \vec{P} + \underbrace{\left(\int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} d\Omega \right)}_{\mathbb{K}} \cdot \vec{V} = \underbrace{\int_{\Omega_e} \mathcal{N}_b d\Omega}_{\vec{f}} \quad (7.51)$$

where the matrix \mathbb{K} is of size $(m_v \cdot ndof_v \times m_v \cdot ndof_v)$, and matrix \mathbb{G} is of size $(m_v \cdot ndof_v \times m_p \cdot ndof_p)$.

Turning now to the mass conservation equation:

$$\begin{aligned}
\vec{0} &= \int_{\Omega_e} \vec{\mathcal{N}}^p \vec{\nabla} \cdot \vec{v} d\Omega \\
&= \int_{\Omega_e} \vec{\mathcal{N}}^p \sum_{i=1}^{m_v} \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \right) d\Omega \\
&= \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_1^p \left(\sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \right) \\ \mathcal{N}_2^p \left(\sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \right) \\ \mathcal{N}_3^p \left(\sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \right) \\ \vdots \\ \mathcal{N}_{m_p}^p \left(\sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \right) \end{pmatrix} dV \\
&= \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_1^p & \mathcal{N}_1^p & \mathcal{N}_1^p & 0 & 0 & 0 \\ \mathcal{N}_2^p & \mathcal{N}_2^p & \mathcal{N}_2^p & 0 & 0 & 0 \\ \mathcal{N}_3^p & \mathcal{N}_3^p & \mathcal{N}_3^p & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathcal{N}_{m_p}^p & \mathcal{N}_{m_p}^p & \mathcal{N}_{m_p}^p & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial x} u_i \\ \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial y} v_i \\ \sum_i \frac{\partial \mathcal{N}_i^\gamma}{\partial z} w_i \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial y} u_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial x} v_i \right) \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial z} u_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial x} w_i \right) \\ \sum_i \left(\frac{\partial \mathcal{N}_i^\gamma}{\partial z} v_i + \frac{\partial \mathcal{N}_i^\gamma}{\partial y} w_i \right) \end{pmatrix} dV \\
&= \int_{\Omega_e} \underbrace{\begin{pmatrix} \mathcal{N}_1^p & \mathcal{N}_1^p & \mathcal{N}_1^p & 0 & 0 & 0 \\ \mathcal{N}_2^p & \mathcal{N}_2^p & \mathcal{N}_2^p & 0 & 0 & 0 \\ \mathcal{N}_3^p & \mathcal{N}_3^p & \mathcal{N}_3^p & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathcal{N}_{m_p}^p & \mathcal{N}_{m_p}^p & \mathcal{N}_{m_p}^p & 0 & 0 & 0 \end{pmatrix}}_{(\mathcal{N}^p)^T} \cdot \vec{\varepsilon} dV \\
&= \left(\int (\mathcal{N}^p)^T \cdot \mathbf{B} dV \right) \cdot \vec{V} \\
&= -\mathbb{G}_e^T \cdot \vec{V}
\end{aligned} \tag{7.52}$$

Note that it is common to actually start from $-\vec{\nabla} \cdot \vec{v} = 0$ (see Eq.(3) in [848]) so as to arrive at $\mathbb{G}_e^T \cdot \vec{V} = \vec{0}$

Ultimately we obtain the following system for each element:

$$\begin{pmatrix} \mathbb{K}_e & \mathbb{G}_e \\ -\mathbb{G}_e^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{V} \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_e \\ 0 \end{pmatrix}$$

Such a matrix is then generated for each element and then must be assembled into the global F.E. matrix. Note that in this case the elemental Stokes matrix is antisymmetric. One can also define

the following symmetric modified Stokes matrix:

$$\begin{pmatrix} \mathbb{K}_e & \mathbb{G}_e \\ \mathbb{G}_e^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{V} \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_e \\ 0 \end{pmatrix} \quad (7.53)$$

This matrix is symmetric, but indefinite. It is non-singular if $\ker(\mathbb{G}^T) = 0$, which is the case if the compatibility condition holds.

CHECK: Matrix \mathbb{K} is the viscosity matrix. Its size is $(ndof_v * N_v) \times (ndof_v * N_v)$ where $ndof_v$ is the number of velocity degrees of freedom per node (typically 1, 2 or 3) and N_v is the number of velocity nodes. The size of matrix \mathbb{G} is $(ndof_v * N_v) \times (ndof_p * N_p)$ where $ndof_p (= 1)$ is the number of velocity degrees of freedom per node and N_p is the number of pressure nodes. Conversely, the size of matrix \mathbb{G}^T is $(ndof_p * N_p) \times (ndof_v * N_v)$. The size of the global FE matrix is $N = ndof_v * N_v + ndof_p * N_p$. Note that matrix \mathbb{K} is analogous to a discrete Laplacian operator, matrix \mathbb{G} to a discrete gradient operator, and matrix \mathbb{G}^T to a discrete divergence operator.

On the physical dimensions of the Stokes matrix blocks

We start from the Stokes equations:

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}) + \rho \vec{g} = \vec{0} \quad (7.54)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (7.55)$$

We have $[p] = ML^{-1}T^{-2}$, $[\vec{\nabla}] = L^{-1}$, so the dimensions of the terms in the first equation are: $ML^{-2}T^{-2}$. The blocks \mathbb{K} and \mathbb{G} stem from the weak form which is obtained by multiplying the strong form equations by the (dimensionless) basis functions and integrating over the 3D domain, so that it follows that

$$[\mathbb{K} \cdot \vec{\mathcal{V}}] = [\mathbb{G} \cdot \vec{\mathcal{P}}] = [f] = (ML^{-2}T^{-2}) \cdot L^3 = MLT^{-2}$$

We can then easily deduce:

$$[\mathbb{K}] = MT^{-1} \quad [\mathbb{G}] = L^2$$

Turning to the mass conservation equation, we have $[\vec{\nabla} \cdot \vec{v}] = L^{-1}LT^{-1} = T^{-1}$, which yields the discretised weak form $\mathbb{G} \cdot \vec{\mathcal{V}} = 0$ so that $[\mathbb{G} \cdot \vec{\mathcal{V}}] = L^3T^{-1}$ and we of course recover $[\mathbb{G}] = L^2$.

If we wanted both equations to have the same dimensions, we would need to multiply the second one by a characteristic quantity which dimension is $ML^{-2}T^{-1}$, i.e. for example η/L (since $[\eta] = ML^{-1}T^{-1}$). This is indeed what we end up doing in practice, see Section 7.5.4.

On elemental level mass balance

Note that in what is above no assumption has been made about whether the pressure basis functions are continuous or discontinuous from one element to another.

Indeed, as mentioned in Gresho & Sani [488], since the weak formulation of the momentum equation involves integration by parts of $\vec{\nabla} p$, the resulting weak form contains no derivatives of pressure. This introduces the possibility of approximating it by functions (piecewise polynomials, of course) that are not C^0 -continuous, and indeed this has been done and is quite popular/useful (e.g. P_0 or P_{-1}).

It is then worth noting that *only* discontinuous pressure elements assure an element-level mass balance [488]: if for instance N_i^p is piecewise-constant on element e (of value 1), the elemental weak form of the mass conservation equation is

$$\int_{\Omega_e} N_i^p \vec{\nabla} \cdot \vec{v} = \int_{\Omega_e} \vec{\nabla} \cdot \vec{v} = \int_{\Gamma_e} \vec{n} \cdot \vec{v} = 0$$

One potentially unwelcome consequence of using discontinuous pressure elements is that they do not possess uniquely defined pressure on the element boundaries; they are dual valued there, and often multi-valued at certain velocity nodes.

On the \mathbf{C} matrix

The relationship between deviatoric stress and deviatoric strain rate tensor is

$$\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}^d \quad (7.56)$$

$$= 2\eta \left(\dot{\boldsymbol{\epsilon}} - \frac{1}{3}(\vec{\nabla} \cdot \vec{\mathbf{v}})\mathbf{1} \right) \quad (7.57)$$

$$= 2\eta \left[\begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} & \dot{\epsilon}_{yz} \\ \dot{\epsilon}_{zx} & \dot{\epsilon}_{zy} & \dot{\epsilon}_{zz} \end{pmatrix} - \frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz}) \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right] \quad (7.58)$$

$$= \frac{2}{3}\eta \begin{pmatrix} 2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} - \dot{\epsilon}_{zz} & 3\dot{\epsilon}_{xy} & 3\dot{\epsilon}_{xz} \\ 3\dot{\epsilon}_{yx} & -\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy} - \dot{\epsilon}_{yy} & 3\dot{\epsilon}_{yz} \\ 3\dot{\epsilon}_{zx} & 3\dot{\epsilon}_{zy} & -\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} + 2\dot{\epsilon}_{zz} \end{pmatrix} \quad (7.59)$$

so that

$$\vec{\tau} = \frac{2}{3}\eta \begin{pmatrix} 2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} - \dot{\epsilon}_{zz} \\ -\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy} - \dot{\epsilon}_{yy} \\ -\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} + 2\dot{\epsilon}_{zz} \\ 3\dot{\epsilon}_{xy} \\ 3\dot{\epsilon}_{xz} \\ 3\dot{\epsilon}_{yz} \end{pmatrix} = \frac{\eta}{3} \underbrace{\begin{pmatrix} 4 & -2 & -2 & 0 & 0 & 0 \\ -2 & 4 & -2 & 0 & 0 & 0 \\ -2 & -2 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{pmatrix}}_{\mathbf{C}^d} \cdot \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} = \mathbf{C}^d \cdot \vec{\epsilon} \quad (7.60)$$

which is identical to the one in the Appendix A of Schmalholz (2008) [1118]. In two dimensions, we have

$$\vec{\tau} = \frac{1}{3}\eta \underbrace{\begin{pmatrix} 4 & -2 & 0 \\ -2 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix}}_{\mathbf{C}^d} \cdot \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ 2\dot{\epsilon}_{xy} \end{pmatrix}$$

see for instance Andres-Martinez *et al.* (2015) [24].

In the case where we assume incompressible flow from the beginning, i.e. $\dot{\boldsymbol{\epsilon}} = \dot{\boldsymbol{\epsilon}}^d$, then

$$\vec{\tau} = \eta \underbrace{\begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{C}} \cdot \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} = \mathbf{C} \cdot \vec{\epsilon} \quad (7.61)$$

A slightly different formulation

The momentum conservation equation can be written as follows:

$$\vec{\nabla} \cdot (2\eta\dot{\boldsymbol{\epsilon}}(\vec{\mathbf{v}})) - \vec{\nabla} p + \vec{b} = \vec{0}$$

When the viscosity η is constant and the flow is incompressible this equation becomes

$$\eta \Delta \vec{v} - \vec{\nabla} p + \vec{b} = \vec{0}$$

In this case the matrix \mathbf{B} takes a different form (See Donea & Huerta [341, Eq. 6.24]) and one should be aware that this can have consequences for the Neumann boundary conditions.

In Burstedde *et al.* (2009) [188] the authors state that when the Laplacian formulation is used it has the computational advantage that the velocity components are coupled only through the incompressibility condition. While the two formulations are equivalent only for constant viscosity, they state that they employ the Laplacian approach formulation as a preconditioner for the viscous term.

Concretely, we apply the same method as above, i.e. we reorganise the terms of the velocity gradient tensor in a vector:

$$\begin{aligned} \vec{\nabla} \vec{v} &\rightarrow \begin{pmatrix} \partial_x u \\ \partial_y u \\ \partial_z u \\ \partial_x v \\ \partial_y v \\ \partial_z v \\ \partial_x w \\ \partial_y w \\ \partial_z w \end{pmatrix} = \begin{pmatrix} \sum_i \partial_x \mathcal{N}_i u_i \\ \sum_i \partial_y \mathcal{N}_i u_i \\ \sum_i \partial_z \mathcal{N}_i u_i \\ \sum_i \partial_x \mathcal{N}_i v_i \\ \sum_i \partial_y \mathcal{N}_i v_i \\ \sum_i \partial_z \mathcal{N}_i v_i \\ \sum_i \partial_x \mathcal{N}_i w_i \\ \sum_i \partial_y \mathcal{N}_i w_i \\ \sum_i \partial_z \mathcal{N}_i w_i \end{pmatrix} \\ &= \underbrace{\begin{pmatrix} \partial_x \mathcal{N}_1^\gamma & 0 & 0 & \partial_x \mathcal{N}_2^\gamma & 0 & 0 & \dots & \partial_x \mathcal{N}_{m_\gamma}^\gamma & 0 & 0 \\ \partial_y \mathcal{N}_1^\gamma & 0 & 0 & \partial_y \mathcal{N}_2^\gamma & 0 & 0 & \dots & \partial_y \mathcal{N}_{m_\gamma}^\gamma & 0 & 0 \\ \partial_z \mathcal{N}_1^\gamma & 0 & 0 & \partial_z \mathcal{N}_2^\gamma & 0 & 0 & \dots & \partial_z \mathcal{N}_{m_\gamma}^\gamma & 0 & 0 \\ 0 & \partial_x \mathcal{N}_1^\gamma & 0 & 0 & \partial_x \mathcal{N}_2^\gamma & 0 & \dots & 0 & \partial_x \mathcal{N}_{m_\gamma}^\gamma & 0 \\ 0 & \partial_y \mathcal{N}_1^\gamma & 0 & 0 & \partial_y \mathcal{N}_2^\gamma & 0 & \dots & 0 & \partial_y \mathcal{N}_{m_\gamma}^\gamma & 0 \\ 0 & \partial_z \mathcal{N}_1^\gamma & 0 & 0 & \partial_z \mathcal{N}_2^\gamma & 0 & \dots & 0 & \partial_z \mathcal{N}_{m_\gamma}^\gamma & 0 \\ 0 & 0 & \partial_x \mathcal{N}_1^\gamma & 0 & 0 & \partial_x \mathcal{N}_2^\gamma & \dots & 0 & 0 & \partial_x \mathcal{N}_{m_\gamma}^\gamma \\ 0 & 0 & \partial_y \mathcal{N}_1^\gamma & 0 & 0 & \partial_y \mathcal{N}_2^\gamma & \dots & 0 & 0 & \partial_y \mathcal{N}_{m_\gamma}^\gamma \\ 0 & 0 & \partial_z \mathcal{N}_1^\gamma & 0 & 0 & \partial_z \mathcal{N}_2^\gamma & \dots & 0 & 0 & \partial_z \mathcal{N}_{m_\gamma}^\gamma \end{pmatrix}}_{\mathbf{B}} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ w_1 \\ u_2 \\ v_2 \\ w_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \\ w_{m_v} \end{pmatrix}}_{\vec{v}} \end{aligned}$$

and in two dimensions:

$$\vec{\nabla} \vec{v} \rightarrow \begin{pmatrix} \partial_x u \\ \partial_y u \\ \partial_x v \\ \partial_y v \end{pmatrix} = \begin{pmatrix} \sum_i \partial_x \mathcal{N}_i u_i \\ \sum_i \partial_y \mathcal{N}_i u_i \\ \sum_i \partial_x \mathcal{N}_i v_i \\ \sum_i \partial_y \mathcal{N}_i v_i \end{pmatrix} = \underbrace{\begin{pmatrix} \partial_x \mathcal{N}_1^\gamma & 0 & \partial_x \mathcal{N}_2^\gamma & 0 & \dots & \partial_x \mathcal{N}_i^\gamma m_\gamma & 0 \\ \partial_y \mathcal{N}_1^\gamma & 0 & \partial_y \mathcal{N}_2^\gamma & 0 & \dots & \partial_y \mathcal{N}_i^\gamma m_\gamma & 0 \\ 0 & \partial_x \mathcal{N}_1^\gamma & 0 & \partial_x \mathcal{N}_2^\gamma & \dots & 0 & \partial_x \mathcal{N}_i^\gamma m_\gamma \\ 0 & \partial_y \mathcal{N}_1^\gamma & 0 & \partial_y \mathcal{N}_2^\gamma & \dots & 0 & \partial_y \mathcal{N}_i^\gamma m_\gamma \end{pmatrix}}_{\mathbf{B}} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix}}_{\vec{v}}.$$

If such a formulation is used, it makes more sense to actually group the unknowns as follows:

$$\vec{v} = (u_1, \dots, u_{m_\gamma}, v_1, \dots, v_{m_\gamma}, w_1, \dots, w_{m_\gamma})$$

We start from

$$\eta \Delta \vec{v} - \vec{\nabla} p + \rho \vec{g} = \vec{0}$$

In 2D Cartesian coordinates this becomes:

$$\eta \Delta u - \partial_x p + \rho g_x = 0 \quad (7.62)$$

$$\eta \Delta v - \partial_y p + \rho g_y = 0 \quad (7.63)$$

or,

$$\eta \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) - \partial_x p + \rho g_x = 0 \quad (7.64)$$

$$\eta \left(\frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) - \partial_y p + \rho g_y = 0 \quad (7.65)$$

Assuming that we prescribe the normal velocity on all sides (i.e. no Neumann boundary conditions), we can establish the weak form of these equations:

$$\underbrace{\left(\int_{\Omega} \eta \left(\frac{\partial \vec{N}^v}{\partial x} \frac{\partial \vec{N}^v}{\partial x} + \frac{\partial \vec{N}^v}{\partial y} \frac{\partial \vec{N}^v}{\partial y} \right) dV \right)}_{\mathbb{N}} \cdot \vec{v}_x + \underbrace{\left(- \int_{\Omega} \frac{\partial \vec{N}^v}{\partial x} \vec{N}^p dV \right)}_{G_x} \cdot \vec{p} = \underbrace{\int_{\Omega} \vec{N}^v \rho g_x dV}_{\vec{f}_x} \quad (7.66)$$

$$\underbrace{\left(\int_{\Omega} \eta \left(\frac{\partial \vec{N}^v}{\partial x} \frac{\partial \vec{N}^v}{\partial x} + \frac{\partial \vec{N}^v}{\partial y} \frac{\partial \vec{N}^v}{\partial y} \right) dV \right)}_{\mathbb{N}} \cdot \vec{v}_y + \underbrace{\left(- \int_{\Omega} \frac{\partial \vec{N}^v}{\partial y} \vec{N}^p dV \right)}_{G_y} \cdot \vec{p} = \underbrace{\int_{\Omega} \vec{N}^v \rho g_y dV}_{\vec{f}_y} \quad (7.67)$$

Turning now to the continuity equation

$$-\vec{\nabla} \cdot \vec{v} = 0$$

or,

$$-\frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} = 0$$

Its weak form then is

$$\underbrace{\left(- \int_{\Omega} \vec{N}^p \frac{\partial \vec{N}^v}{\partial x} dV \right)}_{G_x^T} \cdot \vec{v}_x + \underbrace{\left(- \int_{\Omega} \vec{N}^p \frac{\partial \vec{N}^v}{\partial y} dV \right)}_{G_y^T} \cdot \vec{v}_y = 0$$

In the end:

$$\begin{pmatrix} \mathbb{N} & 0 & G_x \\ 0 & \mathbb{N} & G_y \\ G_x & G_y & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{v}_x \\ \vec{v}_y \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f}_x \\ \vec{f}_y \\ 0 \end{pmatrix}$$

This approach is implemented in [STONE](#) 48.

On the 'forgotten' surface terms

FINISH write

7.5.2 Revisiting the penalty method

We have just seen that the discretised Stokes equation yield the following saddle point system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{0} \end{pmatrix}$$

One can perturb the continuity equation by a term $\mathbb{C}_\epsilon = \epsilon \mathbb{M}_p$ where \mathbb{M}_p is the pressure mass matrix. This yields

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C}_\epsilon \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{0} \end{pmatrix}$$

or,

$$\vec{\mathcal{P}} = \frac{1}{\epsilon} \mathbb{M}_p^{-1} \cdot \mathbb{G}^T \cdot \vec{\mathcal{V}}$$

Substituting pressure in the first equation yields:

$$(\mathbb{K} + \frac{1}{\epsilon} \mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T) \cdot \vec{\mathcal{V}} = \vec{f} \quad (7.68)$$

If we want to solve these equations, it is necessary that the matrix \mathbb{M}_p^{-1} can be computed easily. This is for example the case if \mathbb{M}_p is a lumped mass matrix (often done for Taylor-Hood elements). When discontinuous pressure elements are used, \mathbb{M}_p is in a block diagonal matrix, i.e. a diagonal matrix consisting of small matrices as diagonal elements. One can easily verify that these small matrices have the size of the number of pressure unknowns per element. Note that this is all carried out at the elemental level.

7.5.3 A much more compact derivation of the Stokes matrix blocks

What follows is inspired by chapter 6 of Donea and Huerta [341]. One can easily show that the weak form of the Stokes system can be written

$$\int_{\Omega} \vec{\nabla} \vec{\omega} : \boldsymbol{\sigma} \, d\Omega = \int_{\Omega} \vec{\omega} \cdot \vec{b} \, d\Omega + \int_{\Gamma} \vec{\omega} \cdot \vec{t} \, d\Gamma \quad (7.69)$$

$$\int_{\Omega} q \vec{\nabla} \cdot \vec{v} \, d\Omega = 0 \quad (7.70)$$

where $\vec{\omega}$ and q are the velocity and pressure test functions respectively, and with

$$\vec{\nabla} \vec{\omega} : \boldsymbol{\sigma} = \sum_{i,j}^{ndim} \frac{\partial \omega_i}{\partial x_j} \sigma_{ij}$$

Assuming the Cauchy stress $\boldsymbol{\sigma}$ is given by the linear Stokes' law $\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau}$ with $\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}(\vec{v}) = \mathcal{C} : \vec{\nabla} \vec{v}$ where \mathcal{C} is a fourth-order tensor with $\mathcal{C}_{ijkl} = \eta(\delta_{ik}\delta_{jl} + \delta_{il}\delta_{jk})$. Then

$$\begin{aligned} \int_{\Omega} \vec{\nabla} \vec{\omega} : \boldsymbol{\sigma} \, d\Omega &= \int_{\Omega} \vec{\nabla} \vec{\omega} : (-p\mathbf{1} + \boldsymbol{\tau}) \, d\Omega \\ &= - \int_{\Omega} p \vec{\nabla} \vec{\omega} : \mathbf{1} \, d\Omega + \int_{\Omega} \vec{\nabla} \vec{\omega} : \boldsymbol{\tau} \, d\Omega \\ &= - \int_{\Omega} p \vec{\nabla} \cdot \vec{\omega} \, d\Omega + \int_{\Omega} \vec{\nabla} \vec{\omega} : \mathcal{C} : \vec{\nabla} \vec{v} \, d\Omega \end{aligned} \quad (7.71)$$

The weak form of the Stokes system now takes the form

$$-\int_{\Omega} p \vec{\nabla} \cdot \vec{\omega} \, d\Omega + \int_{\Omega} \vec{\nabla} \vec{\omega} : \mathcal{C} : \vec{\nabla} \vec{v} \, d\Omega = \int_{\Omega} \vec{\omega} \cdot \vec{b} \, d\Omega + \int_{\Gamma} \vec{\omega} \cdot \vec{t} \, d\Gamma \quad (7.72)$$

$$\int_{\Omega} q \vec{\nabla} \cdot \vec{v} \, d\Omega = 0 \quad (7.73)$$

Actually, the bilinear form with the two double dot products is not particularly convenient so it is always rewritten in terms of the strain rate vector

$$\vec{\varepsilon}(\vec{v}) = \begin{pmatrix} \dot{\varepsilon}_{xx}(\vec{v}) \\ \dot{\varepsilon}_{yy}(\vec{v}) \\ \dot{\varepsilon}_{zz}(\vec{v}) \\ 2\dot{\varepsilon}_{xy}(\vec{v}) \\ 2\dot{\varepsilon}_{xz}(\vec{v}) \\ 2\dot{\varepsilon}_{yz}(\vec{v}) \end{pmatrix}$$

as one can easily show that

$$\vec{\nabla} \vec{\omega} : \mathcal{C} : \vec{\nabla} \vec{v} = \vec{\varepsilon}(\vec{\omega})^T \cdot \mathbf{C}_{\eta} \cdot \vec{\varepsilon}(\vec{v})$$

with

$$\mathbf{C}_{\eta} = \eta \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

and since $\vec{\varepsilon}(\vec{v}) = \mathbf{B} \cdot \vec{v}$ and $\vec{\varepsilon}(\vec{\omega}) = \mathbf{B} \cdot \vec{\omega}$ then

$$\begin{aligned} \int_{\Omega} \vec{\nabla} \vec{\omega} : \mathcal{C} : \vec{\nabla} \vec{v} \, d\Omega &= \int_{\Omega} \vec{\varepsilon}(\vec{\omega})^T \cdot \mathbf{C}_{\eta} \cdot \vec{\varepsilon}(\vec{v}) \, d\Omega \\ &= \int_{\Omega} \vec{\omega}^T \cdot \mathbf{B}^T \cdot \mathbf{C}_{\eta} \cdot \mathbf{B} \cdot \vec{v} \, d\Omega \\ &= \vec{\omega}^T \cdot \int_{\Omega} \mathbf{B}^T \cdot \mathbf{C}_{\eta} \cdot \mathbf{B} \, d\Omega \cdot \vec{v} \\ &= \vec{\omega}^T \cdot \mathbb{K} \cdot \vec{v} \end{aligned} \quad (7.74)$$

Let us now turn to the mass conservation equation:

$$\int_{\Omega} q \vec{\nabla} \cdot \vec{v} \, d\Omega = \int_{\Omega} \vec{Q}^T \vec{N}^p \vec{\nabla} \cdot \vec{v} \, d\Omega = \vec{Q}^T \cdot \int_{\Omega} \vec{N}^p \vec{\nabla} \cdot \vec{v} \, d\Omega$$

We have $\vec{\nabla} \cdot \vec{v} = \text{Tr}[\dot{\varepsilon}(\vec{v})]$ but also

$$\begin{aligned} \vec{\nabla} \cdot \vec{v} &= \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} \dot{\varepsilon}_{xx}(\vec{v}) \\ \dot{\varepsilon}_{yy}(\vec{v}) \\ \dot{\varepsilon}_{zz}(\vec{v}) \\ 2\dot{\varepsilon}_{xy}(\vec{v}) \\ 2\dot{\varepsilon}_{xz}(\vec{v}) \\ 2\dot{\varepsilon}_{yz}(\vec{v}) \end{pmatrix} \\ &= \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \end{pmatrix} \cdot \vec{\varepsilon}(\vec{v}) \\ &= \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \end{pmatrix} \cdot \mathbf{B} \cdot \vec{v} \end{aligned} \quad (7.75)$$

so that

$$\begin{aligned}
\vec{\mathcal{N}}^p \vec{\nabla} \cdot \vec{v} &= \vec{\mathcal{N}}^p \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \end{pmatrix} \cdot \mathbf{B} \cdot \vec{v} \\
&= \begin{pmatrix} \vec{\mathcal{N}}^p & \vec{\mathcal{N}}^p & \vec{\mathcal{N}}^p & 0 & 0 & 0 \end{pmatrix} \cdot \mathbf{B} \cdot \vec{v} \\
&= \mathbf{N}^T \cdot \mathbf{B} \cdot \vec{v}
\end{aligned} \tag{7.76}$$

Finally

$$\begin{aligned}
\int_{\Omega} q \vec{\nabla} \cdot \vec{v} \, d\Omega &= \vec{Q}^T \cdot \int_{\Omega} \mathbf{N}^T \cdot \mathbf{B} \, d\Omega \cdot \vec{v} \\
&= -\vec{Q}^T \cdot \mathbb{G}^T \cdot \vec{v}
\end{aligned} \tag{7.77}$$

where

$$\mathbb{G} = - \int_{\Omega} \mathbf{B}^T \cdot \mathbf{N} \, d\Omega$$

Obviously the term $-\int_{\Omega} p \vec{\nabla} \vec{\omega} : \mathbf{1} \, d\Omega = -\int_{\Omega} p \vec{\nabla} \cdot \vec{\omega} \, d\Omega$ will take the form $\vec{\mathcal{P}}^T \cdot \mathbf{G}^T \cdot \vec{\mathcal{W}} = \vec{\mathcal{W}}^T \cdot \mathbf{G} \cdot \vec{\mathcal{P}}$ so that

$$-\int_{\Omega} p \vec{\nabla} \cdot \vec{\omega} \, d\Omega + \int_{\Omega} \vec{\nabla} \vec{\omega} : \mathcal{C} : \vec{\nabla} \vec{v} \, d\Omega = \vec{\mathcal{W}}^T \cdot \mathbf{G} \cdot \vec{\mathcal{P}} + \vec{\mathcal{W}}^T \cdot \mathbb{K} \cdot \vec{v} = \vec{\mathcal{W}}^T \cdot (\mathbf{G} \cdot \vec{\mathcal{P}} + \mathbb{K} \cdot \vec{v})$$

The rhs are handled as shown previously. Since these relationships must work for any test function then it means that what multiplies $\vec{\mathcal{W}}$ and \vec{Q} must be null and we recover Eq. (7.53).

Note that this approach is quite versatile since it does not require to specify the coordinate system. The vector $\vec{\varepsilon}$ will contain the components of the strain rate tensor and in the end the matrix \mathbf{B} will reflect the exact form of the strain rate tensor terms (see axisymmetric formulation).

7.5.4 Pressure scaling

pressure_scaling.tex

As nicely explained in the step 32 of deal.ii¹⁶, we often need to scale the \mathbb{G} block since it is many orders of magnitude smaller than \mathbb{K} (especially in geodynamics where viscosities are $\sim 10^{22}$), which introduces large inaccuracies in the solving process to the point that the solution is nonsensical. This scaling coefficient is η/L where η and L are representative viscosities and lengths. We start from

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix}$$

and introduce the scaling coefficient as follows (which in fact does not alter the solution at all):

$$\begin{pmatrix} \mathbb{K} & \frac{\eta}{L} \mathbb{G} \\ \frac{\eta}{L} \mathbb{G}^T & -\frac{\eta^2}{L^2} \mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \frac{L}{\eta} \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \frac{\eta}{L} \vec{h} \end{pmatrix}$$

We then end up with the modified Stokes system:

$$\begin{pmatrix} \mathbb{K} & \underline{\mathbb{G}} \\ \underline{\mathbb{G}}^T & \underline{\mathbb{C}} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \underline{\vec{p}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \underline{\vec{h}} \end{pmatrix}$$

where

$$\underline{\mathbb{G}} = \frac{\eta}{L} \mathbb{G} \quad \underline{\vec{p}} = \frac{L}{\eta} \vec{p} \quad \underline{\mathbb{C}} = \frac{\eta^2}{L^2} \mathbb{C} \quad \underline{\vec{h}} = \frac{\eta}{L} \vec{h}$$

After the solve phase, we recover the real pressure with $\vec{p} = \frac{\eta}{L} \underline{\vec{p}}$.

¹⁶https://www.dealii.org/9.0.0/doxygen/deal.II/step_32.html

7.5.5 Going from 3D to 2D

The world is three-dimensional. However, for many different reasons one may wish to solve problems which are two-dimensional.

Following ASPECT manual, we will think of two-dimensional models in the following way:

- We assume that the domain we want to solve on is a two-dimensional cross section (in the $x - y$ plane) that extends infinitely far in both negative and positive z direction.
- We assume that the velocity is zero in the z direction and that all variables have no variation in the z direction.

As a consequence, two-dimensional models are three-dimensional ones in which the z component of the velocity is zero and so are all z derivatives. This allows to reduce the momentum conservation equations from 3 equations to 2 equations. However, contrarily to what is often seen, the 3D definition of the deviatoric strain rate remains, i.e. in other words:

$$\dot{\epsilon}^d = \dot{\epsilon} - \frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} \quad (7.78)$$

and not $1/2$. In light of all this, the full strain rate tensor and the deviatoric strain rate tensor in 2D are given by:

$$\epsilon = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} & \dot{\epsilon}_{yz} \\ \dot{\epsilon}_{zx} & \dot{\epsilon}_{zy} & \dot{\epsilon}_{zz} \end{pmatrix} = \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) & 0 \\ \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) & \frac{\partial v}{\partial y} & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (7.79)$$

$$\dot{\epsilon}^d = \frac{1}{3} \begin{pmatrix} 2\frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} & \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) & 0 \\ \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) & -\frac{\partial u}{\partial x} + 2\frac{\partial v}{\partial y} & 0 \\ 0 & 0 & -\frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} \end{pmatrix} \quad (7.80)$$

Although the bottom right term may be surprising, it is of no consequence when this expression of the deviatoric strain rate is used in the Stokes equation:

$$\vec{\nabla} \cdot 2\eta \dot{\epsilon}^d =$$

FINISH!

In two dimensions the velocity is then $\vec{v} = (u, v)$ and the FEM building blocks and matrices are simply:

$$\vec{\epsilon} = \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ 2\dot{\epsilon}_{xy} \end{pmatrix} = \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial N_1^\vee}{\partial x} & 0 & \frac{\partial N_2^\vee}{\partial x} & 0 & \frac{\partial N_3^\vee}{\partial x} & 0 & \dots & \frac{\partial N_{m_v}^\vee}{\partial x} & 0 \\ 0 & \frac{\partial N_1^\vee}{\partial y} & 0 & \frac{\partial N_2^\vee}{\partial y} & 0 & \frac{\partial N_3^\vee}{\partial y} & \dots & 0 & \frac{\partial N_{m_v}^\vee}{\partial x} \\ \frac{\partial N_1^\vee}{\partial y} & \frac{\partial N_1^\vee}{\partial x} & \frac{\partial N_2^\vee}{\partial y} & \frac{\partial N_2^\vee}{\partial x} & \frac{\partial N_3^\vee}{\partial y} & \frac{\partial N_3^\vee}{\partial x} & \dots & \frac{\partial N_{m_v}^\vee}{\partial y} & \frac{\partial N_{m_v}^\vee}{\partial x} \end{pmatrix}}_B \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix}}_{\vec{v}} \quad (7.81)$$

we have

$$\sigma_{xx} = -p + 2\eta\dot{\epsilon}_{xx} \quad (7.82)$$

$$\sigma_{yy} = -p + 2\eta\dot{\epsilon}_{yy} \quad (7.83)$$

$$\sigma_{xy} = +2\eta\dot{\epsilon}_{xy} \quad (7.84)$$

so

$$\vec{\sigma} = - \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} p + \mathbf{C} \cdot \vec{\epsilon} = - \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \vec{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (7.85)$$

with

$$\mathbf{C} = \eta \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{or} \quad \mathbf{C} = \frac{\eta}{3} \begin{pmatrix} 4 & -2 & 0 \\ -2 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad (7.86)$$

check the right **C**

Finally the matrix \mathbf{N}^p is of size $3 \times m_p$:

$$\mathbf{N}^p = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \vec{N}^p = \begin{pmatrix} \vec{N}^p \\ \vec{N}^p \\ 0 \end{pmatrix} \quad (7.87)$$

7.5.6 The cylindrical axisymmetric case

mixed_axisymmetric.tex

In cylindrical coordinates the velocity gradient is given by

$$\vec{\nabla} \vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & \frac{1}{r} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r} & \frac{\partial v_r}{\partial z} \\ \frac{\partial v_\theta}{\partial r} & \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} + \frac{v_r}{r} & \frac{\partial v_\theta}{\partial z} \\ \frac{\partial v_z}{\partial r} & \frac{1}{r} \frac{\partial v_z}{\partial \theta} & \frac{\partial v_z}{\partial z} \end{pmatrix} \quad (7.88)$$

In the case of axisymmetry, and in this case symmetry about the z axis, there is invariance with respect to the rotation around the axis so stresses and other quantities are independent of the θ coordinate, or simply put $\partial_\theta \rightarrow 0$. The velocity gradient simplifies to:

$$\vec{\nabla} \vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & -\frac{v_\theta}{r} & \frac{\partial v_r}{\partial z} \\ \frac{\partial v_\theta}{\partial r} & \frac{v_r}{r} & \frac{\partial v_\theta}{\partial z} \\ \frac{\partial v_z}{\partial r} & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix} \quad (7.89)$$

Also, it follows logically that $v_\theta = 0$ so that ultimately:

$$\vec{\nabla} \vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & 0 & \frac{\partial v_r}{\partial z} \\ 0 & \frac{v_r}{r} & 0 \\ \frac{\partial v_z}{\partial r} & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix} \quad (7.90)$$

and the strain rate tensor is then given by

$$\dot{\epsilon}(\vec{v}) = \frac{1}{2} (\vec{\nabla} \vec{v} + \vec{\nabla} \vec{v}^T) = \begin{pmatrix} \frac{\partial v_r}{\partial r} & 0 & \frac{1}{2} (\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) \\ 0 & \frac{v_r}{r} & 0 \\ \frac{1}{2} (\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix} \quad (7.91)$$

The velocity divergence $\vec{\nabla} \cdot \vec{v}$ is simply the trace of $\dot{\epsilon}(\vec{v})$ so

$$\vec{\nabla} \cdot \vec{v} = \frac{\partial v_r}{\partial r} + \frac{v_r}{r} + \frac{\partial v_z}{\partial z}$$

The components of the $\vec{\epsilon}(\vec{v})$ vector are

$$\vec{\epsilon}(\vec{v}) = \begin{pmatrix} \dot{\epsilon}_{rr} \\ \dot{\epsilon}_{\theta\theta} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{r\theta} \\ 2\dot{\epsilon}_{rz} \\ 2\dot{\epsilon}_{\theta z} \end{pmatrix} = \begin{pmatrix} \frac{\partial v_r}{\partial r} \\ \frac{v_r}{r} \\ \frac{\partial v_z}{\partial z} \\ 0 \\ \frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z} \\ 0 \end{pmatrix}$$

We see that there are two zeroes and consequently we only keep the four non zero components:

$$\vec{\epsilon}(\vec{v}) = \begin{pmatrix} \frac{\partial v_r}{\partial r} \\ \frac{v_r}{r} \\ \frac{\partial v_z}{\partial z} \\ \frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z} \end{pmatrix}$$

Only displacements in the r and z directions remain (note that $\dot{\varepsilon}_{\theta\theta}$ is in fact equal to \mathbf{v}_r/r). In what follows I rename $u = \mathbf{v}_r$ and $w = \mathbf{v}_z$ to simplify notations. Then, inside an element we have

$$u^h(r, z) = \sum_{i=1}^{m_\nu} \mathcal{N}_i^\nu(r, z) u_i \quad (7.92)$$

$$w^h(r, z) = \sum_{i=1}^{m_\nu} \mathcal{N}_i^\nu(r, z) w_i \quad (7.93)$$

where \mathcal{N}_i^ν are the velocity basis functions attached to the m_ν nodes of the element. We compute the elements of the $\vec{\varepsilon}(\vec{\mathbf{v}})$ vector as follows:

$$\dot{\varepsilon}_{rr} = \frac{\partial u^h}{\partial r} = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, z) u_i \quad (7.94)$$

$$\dot{\varepsilon}_{\theta\theta} = \frac{u_r^h}{r} = \frac{1}{r} \sum_{i=1}^m \mathcal{N}_i(r, z) u_i \quad (7.95)$$

$$\dot{\varepsilon}_{zz} = \frac{\partial w^h}{\partial z} = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial z}(r, z) w_i \quad (7.96)$$

$$2\dot{\varepsilon}_{rz} = \frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial r} = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial z}(r, z) u_i + \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, z) w_i \quad (7.97)$$

and then

$$\vec{\varepsilon}^h = \begin{pmatrix} \frac{\partial u^h}{\partial r} \\ \frac{u^h}{r} \\ \frac{\partial w^h}{\partial z} \\ \frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial r} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \mathcal{N}_1}{\partial r} & 0 & \frac{\partial \mathcal{N}_2}{\partial r} & 0 & \dots & \dots & \frac{\partial \mathcal{N}_{m_\nu}}{\partial r} & 0 \\ \frac{\mathcal{N}_1}{r} & 0 & \frac{\mathcal{N}_2}{r} & 0 & \dots & \dots & \frac{\mathcal{N}_{m_\nu}}{r} & 0 \\ 0 & \frac{\partial \mathcal{N}_1}{\partial z} & 0 & \frac{\partial \mathcal{N}_2}{\partial z} & \dots & \dots & 0 & \frac{\partial \mathcal{N}_{m_\nu}}{\partial z} \\ \frac{\partial \mathcal{N}_1}{\partial z} & \frac{\partial \mathcal{N}_1}{\partial r} & \frac{\partial \mathcal{N}_2}{\partial z} & \frac{\partial \mathcal{N}_2}{\partial r} & \dots & \dots & \frac{\partial \mathcal{N}_{m_\nu}}{\partial z} & \frac{\partial \mathcal{N}_{m_\nu}}{\partial r} \end{pmatrix}}_{\mathbf{B}(4 \times 2m_\nu)} \cdot \underbrace{\begin{pmatrix} u_1 \\ w_1 \\ u_2 \\ w_2 \\ \vdots \\ u_{m_\nu} \\ w_{m_\nu} \end{pmatrix}}_{\vec{\mathbf{v}}(2m_\nu \times 1)} \quad (7.98)$$

or $\vec{\varepsilon}^h = \mathbf{B} \cdot \vec{\mathbf{v}}$, where $\vec{\mathbf{v}}$ is the vector of velocity dofs for an element.

Following the presentation of Section 7.5.3, we know that we will obtain

$$\underbrace{\left(- \int_{\Omega_e} \mathbf{B}^T \cdot \mathcal{N}^p dV \right)}_{\mathbb{G}_e} \cdot \vec{P} + \underbrace{\left(\int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} dV \right)}_{\mathbb{K}_e} \cdot \vec{\mathbf{v}} = \underbrace{\int_{\Omega_e} \vec{\mathcal{N}}_b dV}_{\vec{f}_e} \quad (7.99)$$

with

$$\mathbf{C}_\eta = \eta \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Note: We have in cylindrical coordinates $dV = r dr d\theta dz$. The integral over the θ coordinate yields a factor 2π so for instance

$$\mathbb{K}_e = 2\pi \iint_{\Omega_e} \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} r dr dz \quad (7.100)$$

Note the r term in the integrand. The integration can now be performed as simply as was the case in the plane strain problem.

Note that it is common to actually start from $-\vec{\nabla} \cdot \vec{v} = 0$ (see Eq. (3) in [848]) so as to arrive at $\mathbb{G}_e^T \cdot \vec{V} = \vec{0}$. Ultimately we obtain the following system for each element:

$$\begin{pmatrix} \mathbb{K}_e & \mathbb{G}_e \\ \mathbb{G}_e^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{V} \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_e \\ 0 \end{pmatrix}$$

Such a matrix is then generated for each element and then must be assembled into the global F.E. matrix.

Unfortunately there is not much teaching/practical material to be found in the literature with regards to axisymmetric flow. For example Donea and Huerta [341] do not even mention this problem. In Hughes [604] we find:

$$\boldsymbol{\sigma} = \begin{Bmatrix} \sigma_{11} \\ \sigma_{22} \\ \sigma_{12} \\ \sigma_{33} \end{Bmatrix} = \begin{Bmatrix} \sigma_{rr} \\ \sigma_{zz} \\ \sigma_{rz} \\ \sigma_{\theta\theta} \end{Bmatrix} \quad (2.12.3)$$

$$\boldsymbol{\epsilon} = \begin{Bmatrix} \epsilon_{11} \\ \epsilon_{22} \\ 2\epsilon_{12} \\ \epsilon_{33} \end{Bmatrix} = \begin{Bmatrix} \epsilon_{rr} \\ \epsilon_{zz} \\ 2\epsilon_{rz} \\ \epsilon_{\theta\theta} \end{Bmatrix} \quad (2.12.4)$$

The ordering emanates from the following generalization of Table 2.7.1:

| I J | i k | j l |
|------------|------------|------------|
| 1 | 1 | 1 |
| 3 | 1 | 2 |
| 3 | 2 | 1 |
| 2 | 2 | 2 |
| 4 | 3 | 3 |

The \mathbf{D} array takes on the following form:

$$\mathbf{D} = [D_{IJ}] = \underbrace{\begin{bmatrix} D_{33} & D_3 \\ D_3^T & D_{44} \end{bmatrix}}_{4 \times 4} \quad (2.12.5)$$

$$\mathbf{D}_{33} = \begin{bmatrix} D_{11} & D_{12} & D_{13} \\ & D_{22} & D_{23} \\ \text{symmetric} & & D_{33} \end{bmatrix} \quad (2.12.6)$$

$$\mathbf{D}_3 = \begin{Bmatrix} D_{14} \\ D_{24} \\ D_{34} \end{Bmatrix} \quad (2.12.7)$$

where $D_{IJ} = c_{ijkl}$, in which the indices are related by the table. The \mathbf{B}_e -matrix takes on the form

$$\mathbf{B}_e = \begin{bmatrix} N_{a,1} & 0 \\ 0 & N_{a,2} \\ N_{a,2} & N_{a,1} \\ \hline \frac{N_a}{r} & 0 \end{bmatrix} \quad (2.12.8)$$

Again, a factor of $2\pi r$ needs to be included in all integrands.

Remark

The *plane strain* case may be obtained from the axisymmetric formulation by

- i. Ignoring the $2\pi r$ factors; and
- ii. Ignoring the fourth row of \mathbf{B}_e and the fourth row and column of \mathbf{D} .

Taken from pages 101 of T.J.R. Hughes. *The Finite Element Method. Linear Static and Dynamic Finite Element Analysis*. Dover Publications, Inc., 2000. ISBN:

0-486-41181-8.

Also check page 469 of Gresho & Sani's book, and see their remark on axisymmetric case for the N-S equations on page 545.



Relevant Literature:

- Masahisa Tabata. “Finite element analysis of axisymmetric flow problems”. In: *Zeitschrift für angewandte Mathematik und Mechanik* 76 (1996), pp. 171–174
- Vitoriano Ruas. “Mixed finite element methods with discontinuous pressures for the axisymmetric Stokes problem”. In: *ZAMM-Journal of Applied Mathematics and Mechanics*/Zeitschrift für Angewandte Mathematik und Mechanik: Applied Mathematics and Mechanics 83.4 (2003), pp. 249–264. DOI: 10.1002/zamm.200310032
- Young-Ju Lee and Hengguang Li. “Axisymmetric Stokes equations in polygonal domains: regularity and finite element approximations”. In: *Computers & Mathematics with Applications* 64.11 (2012), pp. 3500–3521. DOI: 10.1016/j.camwa.2012.08.014
- TE Price. “Numerically exact integration of a family of axisymmetric finite elements”. In: *Communications in numerical methods in engineering* 19.4 (2003), pp. 253–261. DOI: 10.1002/cnm.583

7.6 Mappings & Jacobians

mappings.tex

The name *isoparametric* derives from the fact that the same ('iso') functions are used as basis functions and for the mapping to the reference element.

More generally, if n_e denotes the number of nodes of an element and n_g denotes the number of nodes describing the geometry of the element, then the element is termed *subparametric* when $n_g < n_e$ and *superparametric* when $n_g > n_e$.

7.6.1 General case

What follows is written for the 2d case but extending it to 3d is trivial.

Any variable defined on the element can be approximated using the basis functions:

$$f^h(r, s) = \sum_i \mathcal{N}_i(r, s) f_i. \quad (7.101)$$

If we treat the coordinate variables x and y themselves as functions, then the basis functions can be used to construct the mapping:

$$x = \sum_{i=1}^4 \mathcal{N}_i(r, s) x_i \quad (7.102)$$

$$y = \sum_{i=1}^4 \mathcal{N}_i(r, s) y_i \quad (7.103)$$

This is a relationship between the reduced coordinates r, s and the 'real' coordinates x, y .

Let us compute the space derivatives of these quantities:

$$\frac{\partial x}{\partial r} = \sum_i \frac{\partial \mathcal{N}_i}{\partial r} x_i \quad (7.104)$$

$$\frac{\partial x}{\partial s} = \sum_i \frac{\partial \mathcal{N}_i}{\partial s} x_i \quad (7.105)$$

$$\frac{\partial y}{\partial r} = \sum_i \frac{\partial \mathcal{N}_i}{\partial r} y_i \quad (7.106)$$

$$\frac{\partial y}{\partial s} = \sum_i \frac{\partial \mathcal{N}_i}{\partial s} y_i \quad (7.107)$$

We also have

$$\frac{\partial f}{\partial r} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial r} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial r} \quad (7.108)$$

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial s} \quad (7.109)$$

or in matrix form:

$$\begin{pmatrix} \frac{\partial f}{\partial r} \\ \frac{\partial f}{\partial s} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial y}{\partial r} \\ \frac{\partial x}{\partial s} & \frac{\partial y}{\partial s} \end{pmatrix}}_{\mathbf{J}} \cdot \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix}$$

where \mathbf{J} is called the Jacobian of the transformation. By inverting the Jacobian matrix, the desired derivatives with respect to x and y can be obtained:

$$\begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} = \mathbf{J}^{-1} \cdot \begin{pmatrix} \frac{\partial f}{\partial r} \\ \frac{\partial f}{\partial s} \end{pmatrix}$$

The inverse of the Jacobian matrix can be simply obtained in 2D (Cramer's rule for 2×2 matrices¹⁷):

$$\mathbf{J}^{-1} = \frac{1}{|\mathbf{J}|} \begin{pmatrix} \frac{\partial y}{\partial s} & -\frac{\partial y}{\partial r} \\ -\frac{\partial x}{\partial s} & \frac{\partial x}{\partial r} \end{pmatrix}$$

The presence of the determinant in the denominator implies that it cannot be zero anywhere, or in other words: the mapping is not valid if $|\mathbf{J}|$ is zero anywhere over the element.

Remark. Problems also arise when the Jacobian matrix is nearly singular due to round-off errors. To avoid problems linked to badly shaped elements, it is recommended that the inside angles of an element are larger than 15° and less than 165° .

From Eq. (??), we can also write:

$$dx = \frac{\partial x}{\partial r} dr + \frac{\partial x}{\partial s} ds \quad (7.110)$$

$$dy = \frac{\partial y}{\partial r} dr + \frac{\partial y}{\partial s} ds \quad (7.111)$$

or,

$$\begin{pmatrix} dx \\ dy \end{pmatrix} = \mathbf{J} \cdot \begin{pmatrix} dr \\ ds \end{pmatrix} \quad (7.112)$$

This means that integrating over the 'real' element in (x, y) space can be simply done by integrating of the reference element in the (r, s) space. This is the cornerstone of most of the implementation of the Finite Element Method, the second integral being carried out by means of the Gauss-Legendre quadrature.

$$\iint_{\Omega_e} \dots dx dy = \int_{-1}^{+1} \int_{-1}^{+1} \dots |\mathbf{J}| dr ds \quad (7.113)$$

7.6.2 Linear mapping on a triangle

$$\begin{aligned} x &= \sum_{i=1}^3 N_i(r, s) x_i \\ &= N_1(r, s) x_1 + N_2(r, s) x_2 + N_3(r, s) x_3 \\ &= (1 - r - s) x_1 + (r) x_2 + (s) x_3 \\ &= x_1 + (x_2 - x_1) r + (x_3 - x_1) s \\ &= a_x + b_x r + c_x s \\ y &= \sum_{i=1}^3 N_i(r, s) y_i \\ &= N_1(r, s) y_1 + N_2(r, s) y_2 + N_3(r, s) y_3 \\ &= (1 - r - s) y_1 + (r) y_2 + (s) y_3 \\ &= y_1 + (y_2 - y_1) r + (y_3 - y_1) s \\ &= a_y + b_y r + c_y s \end{aligned}$$

¹⁷https://en.wikipedia.org/wiki/Cramers_rule

Let us compute the space derivatives of these quantities:

$$\begin{aligned}\frac{\partial x}{\partial r} &= x_2 - x_1 = b_x \\ \frac{\partial x}{\partial s} &= x_3 - x_1 = c_x \\ \frac{\partial y}{\partial r} &= y_2 - y_1 = b_y \\ \frac{\partial y}{\partial s} &= y_3 - y_1 = c_y\end{aligned}$$

The jacobian matrix is then given by

$$\mathbf{J} = \begin{pmatrix} x_2 - x_1 & y_2 - y_1 \\ x_3 - x_1 & y_3 - y_1 \end{pmatrix} = \begin{pmatrix} b_x & b_y \\ c_x & c_y \end{pmatrix}$$

and its inverse

$$\mathbf{J}^{-1} = \frac{1}{b_x c_y - c_x b_y} \begin{pmatrix} c_y & -b_y \\ -c_x & b_x \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} c_y & -b_y \\ -c_x & b_x \end{pmatrix}$$

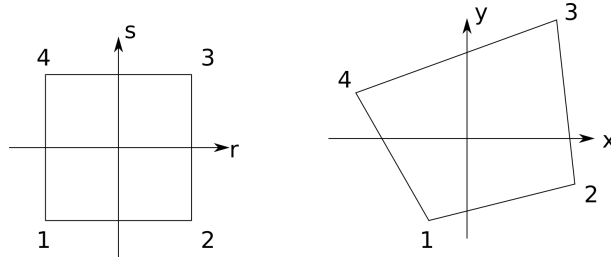
where $A = (b_x c_y - c_x b_y)/2$ is actually the area of the triangle!

The Cartesian basis function derivatives are then

$$\begin{aligned}\begin{pmatrix} \frac{\partial N_1}{\partial x} \\ \frac{\partial N_1}{\partial y} \end{pmatrix} &= \mathbf{J}^{-1} \cdot \begin{pmatrix} \frac{\partial N_1}{\partial r} \\ \frac{\partial N_1}{\partial s} \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} c_y & -b_y \\ -c_x & b_x \end{pmatrix} \cdot \begin{pmatrix} -1 \\ -1 \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} b_y - c_y \\ c_x - b_x \end{pmatrix} \\ \begin{pmatrix} \frac{\partial N_2}{\partial x} \\ \frac{\partial N_2}{\partial y} \end{pmatrix} &= \mathbf{J}^{-1} \cdot \begin{pmatrix} \frac{\partial N_2}{\partial r} \\ \frac{\partial N_2}{\partial s} \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} c_y & -b_y \\ -c_x & b_x \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} c_y \\ -c_x \end{pmatrix} \\ \begin{pmatrix} \frac{\partial N_3}{\partial x} \\ \frac{\partial N_3}{\partial y} \end{pmatrix} &= \mathbf{J}^{-1} \cdot \begin{pmatrix} \frac{\partial N_3}{\partial r} \\ \frac{\partial N_3}{\partial s} \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} c_y & -b_y \\ -c_x & b_x \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \frac{1}{2A} \begin{pmatrix} -b_y \\ b_x \end{pmatrix}\end{aligned}$$

7.6.3 Bilinear mapping (Q_1) on a quadrilateral

The reference element is in the (r, s) space. It is a square of size 2×2 centered around the origin, i.e. $(r, s) \in [-1, 1] \times [-1, 1]$. We wish to map it to the quadrilateral in the (x, y) space (and vice versa):



The coordinates of the vertices are (x_1, y_1) , (x_2, y_2) , (x_3, y_3) and (x_4, y_4) . We then simply have the following relationship, i.e. any point of the reference element can be mapped to the physical quadrilateral as follows:

$$x = \mathcal{N}_1(r, s)x_1 + \mathcal{N}_2(r, s)x_2 + \mathcal{N}_3(r, s)x_3 + \mathcal{N}_4(r, s)x_4 \quad (7.114)$$

$$y = \mathcal{N}_1(r, s)y_1 + \mathcal{N}_2(r, s)y_2 + \mathcal{N}_3(r, s)y_3 + \mathcal{N}_4(r, s)y_4 \quad (7.115)$$

where the Q_1 basis functions $\mathcal{N}_i(r, s)$ are defined in Section 5.2.

In the following example the program randomly generates 10000 points inside the reference element and computes their mapping into the (x, y) space.

```

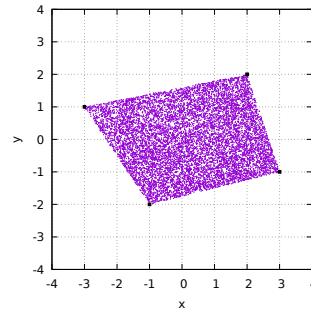
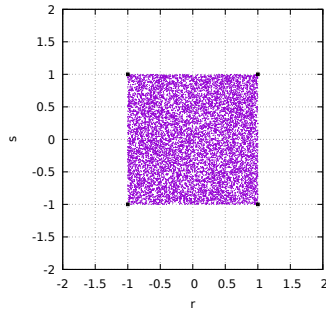
x1=-1 ; y1=-2
x2=3  ; y2=-1
x3=2  ; y3=2
x4=-3 ; y4=1

npts=10000
r=np.zeros(npts,dtype=np.float64)
s=np.zeros(npts,dtype=np.float64)
x=np.zeros(npts,dtype=np.float64)
y=np.zeros(npts,dtype=np.float64)

for i in range(0,npts):
    # compute random r,s coordinates
    r[i]=random.uniform(-1.,+1)
    s[i]=random.uniform(-1.,+1)
    # compute basis function values at r,s
    N1=0.25*(1-r[i])*(1-s[i])
    N2=0.25*(1+r[i])*(1-s[i])
    N3=0.25*(1+r[i])*(1+s[i])
    N4=0.25*(1-r[i])*(1+s[i])
    # compute x,y coordinates
    x[i]=N1*x1+N2*x2+N3*x3+N4*x4
    y[i]=N1*y1+N2*y2+N3*y3+N4*y4

np.savetxt('rs.ascii',np.array([r,s]).T)
np.savetxt('xy.ascii',np.array([x,y]).T)

```



There is also an inverse map, which is not so easily computed (see Section 9.11). However, if the quadrilateral in the (x, y) space is a rectangle of size (h_x, h_y) , the inverse mapping is trivial:

$$r = \frac{x - x_1}{x_2 - x_1} \quad (7.116)$$

$$s = \frac{y - y_1}{y_4 - y_1} \quad (7.117)$$

Also in the case of rectangular elements of size (h_x, h_y) the basis functions can easily be written as functions of (x, y) :

$$\mathcal{N}_1(x, y) = \left(\frac{x_3 - x}{h_x} \right) \left(\frac{y_3 - y}{h_y} \right)$$

$$\mathcal{N}_2(x, y) = \left(\frac{x - x_1}{h_x} \right) \left(\frac{y_3 - y}{h_y} \right)$$

$$\mathcal{N}_3(x, y) = \left(\frac{x - x_1}{h_x} \right) \left(\frac{y - y_1}{h_y} \right)$$

$$\mathcal{N}_4(x, y) = \left(\frac{x_3 - x}{h_x} \right) \left(\frac{y - y_1}{h_y} \right)$$

From Eq. (7.115) and using the expressions for the Q_1 basis functions, we can write

$$\begin{aligned} x &= \frac{1}{4}(x_1 + x_2 + x_3 + x_4) + \frac{1}{4}(-x_1 + x_2 + x_3 - x_4)r + \frac{1}{4}(-x_1 - x_2 + x_3 + x_4)s + \frac{1}{4}(x_1 - x_2 + x_3 - x_4)rs \\ y &= \frac{1}{4}(y_1 + y_2 + y_3 + y_4) + \frac{1}{4}(-y_1 + y_2 + y_3 - y_4)r + \frac{1}{4}(-y_1 - y_2 + y_3 + y_4)s + \frac{1}{4}(y_1 - y_2 + y_3 - y_4)rs \end{aligned}$$

Let us compute the space derivatives of these quantities:

$$\begin{aligned} \frac{\partial x}{\partial r} &= \frac{1}{4}(-x_1 + x_2 + x_3 - x_4) + \frac{1}{4}(x_1 - x_2 + x_3 - x_4)s = A_1 + A_2s \\ \frac{\partial x}{\partial s} &= \frac{1}{4}(-x_1 - x_2 + x_3 + x_4) + \frac{1}{4}(x_1 - x_2 + x_3 - x_4)r = A_3 + A_4r \\ \frac{\partial y}{\partial r} &= \frac{1}{4}(-y_1 + y_2 + y_3 - y_4) + \frac{1}{4}(y_1 - y_2 + y_3 - y_4)s = B_1 + B_2s \\ \frac{\partial y}{\partial s} &= \frac{1}{4}(-y_1 - y_2 + y_3 + y_4) + \frac{1}{4}(y_1 - y_2 + y_3 - y_4)r = B_3 + B_4r \end{aligned}$$

The jacobian matrix is then given by

$$\mathbf{J} = \begin{pmatrix} A_1 + A_2s & B_1 + B_2s \\ A_3 + A_4r & B_3 + B_4r \end{pmatrix}$$

and its inverse

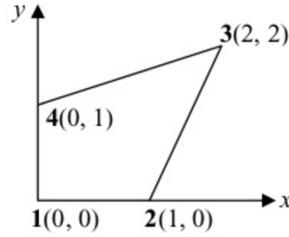
$$\mathbf{J}^{-1} = \frac{1}{C} \begin{pmatrix} B_3 + B_4r & -B_1 - B_2s \\ -A_3 - A_4r & A_1 + A_2s \end{pmatrix}$$

with C being the determinant given by

$$C = (A_1 + A_2s)(B_3 + B_4r) - (A_3 + A_4r)(B_1 + B_2s)$$

A concrete example

Let us look at this by means of a simple example and let us consider the following element:



Then a Q_1 mapping yields:

$$\begin{aligned} x(r, s) &= \sum_i \mathcal{N}_i(r, s)x_i = \mathcal{N}_2 + 2\mathcal{N}_3 = \frac{1}{4}(3 + 3r + s + rt) \\ y(r, s) &= \sum_i \mathcal{N}_i(r, s)y_i = 2\mathcal{N}_3 + \mathcal{N}_4 = \frac{1}{4}(3 + r + 3s + rt) \end{aligned} \quad (7.118)$$

The Jacobian matrix is then

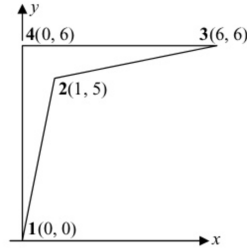
$$\mathbf{J} = \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial y}{\partial r} \\ \frac{\partial x}{\partial s} & \frac{\partial y}{\partial s} \end{pmatrix} = \frac{1}{4} \begin{pmatrix} 3 + s & 1 + s \\ 1 + r & 3 + r \end{pmatrix}$$

and its determinant is

$$|\mathbf{J}| = \frac{1}{4}[(3+s)(3+r) - (1+s)(1+r)] = \frac{1}{2} + \frac{1}{8}r + \frac{1}{8}s \quad (7.119)$$

It is clear that $|\mathbf{J}| > 0$ for $-1 \leq r \leq +1$ and $-1 \leq s \leq +1$.

Let us now consider another example, the following element:



It follows that

$$x(r, s) = \sum_i \mathcal{N}_i(r, s) x_i = \frac{1}{4}(1+r)(7+5s) \quad (7.120)$$

$$y(r, s) = \sum_i \mathcal{N}_i(r, s) y_i = \frac{1}{4}(17+5r+7s-5rs) \quad (7.121)$$

and the determinant:

$$|\mathbf{J}| = \frac{3}{2} - \frac{15r}{4} + \frac{15s}{4}$$

is zero for $r - s = 2/5$. This mapping is invalid!

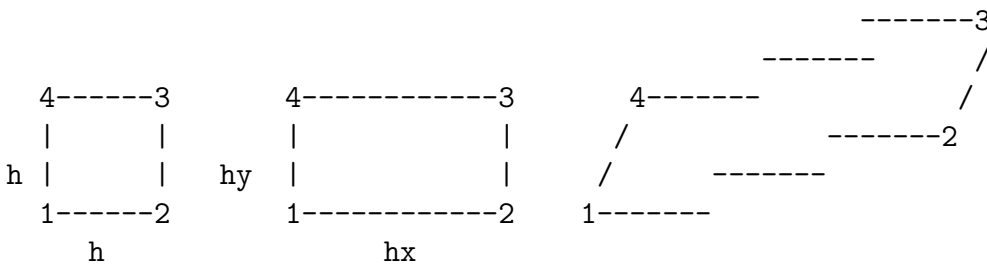
the special case of squares, quadrilaterals and parallelograms

In this case the parameters A2, A4, B2, B4 are all equal to zero. which yields:

$$\mathbf{J}^{-1} = \frac{1}{C} \begin{pmatrix} B_3 & -B_1 \\ -A_3 & A_1 \end{pmatrix} \quad \text{with} \quad C = A_1 B_3 - A_3 B_1$$

Indeed, let us draw such quadrilaterals:

redo with tikz



For the square of size h we have

$$\begin{aligned} 4A_2 &= x_1 - x_2 + x_3 - x_4 = x_1 - (x_1 + h) + (x_4 + h) - x_4 = 0 \\ 4A_4 &= x_1 - x_2 + x_3 - x_4 = x_1 - (x_1 + h) + (x_4 + h) - x_4 = 0 \\ 4B_2 &= y_1 - y_2 + y_3 - y_4 = y_1 - y_2 + (y_2 + h_y) - (y_1 + h_y) = 0 \\ 4B_4 &= y_1 - y_2 + y_3 - y_4 = y_1 - y_2 + (y_2 + h_y) - (y_1 + h_y) = 0 \end{aligned}$$

For the rectangle of size h_x, h_y we have

$$4A_2 = x_1 - x_2 + x_3 - x_4 = x_1 - (x_1 + h_x) + (x_4 + h_x) - x_4 = 0$$

$$4A_4 = x_1 - x_2 + x_3 - x_4 = x_1 - (x_1 + h_x) + (x_4 + h_x) - x_4 = 0$$

$$4B_2 = y_1 - y_2 + y_3 - y_4 = y_1 - y_2 + (y_2 + h_y) - (y_1 + h_y) = 0$$

$$4B_4 = y_1 - y_2 + y_3 - y_4 = y_1 - y_2 + (y_2 + h_y) - (y_1 + h_y) = 0$$

and the same for the parallelogram.

In the case of a rectangle we also have

$$A_1 = \frac{1}{4}(-x_1 + x_2 + x_3 - x_4) = \frac{1}{4}(-x_1 + (x_1 + h_x) + (x_4 + h_x) - x_4) = \frac{h_x}{2} \quad (7.122)$$

$$A_3 = \frac{1}{4}(-x_1 - x_2 + x_3 + x_4) = \frac{1}{4}(-x_1 - (x_1 + h_x) + (x_4 + h_x) + x_4) = 0 \quad (7.123)$$

$$B_1 = \frac{1}{4}(-y_1 + y_2 + y_3 - y_4) = \frac{1}{4}(-y_1 + y_2 + (y_2 + h_y) - (y_1 + h_y)) = 0 \quad (7.124)$$

$$B_3 = \frac{1}{4}(-y_1 - y_2 + y_3 + y_4) = \frac{1}{4}(-y_1 - y_2 + (y_2 + h_y) + (y_1 + h_y)) = \frac{h_y}{2} \quad (7.125)$$

so that the jacobian matrix is

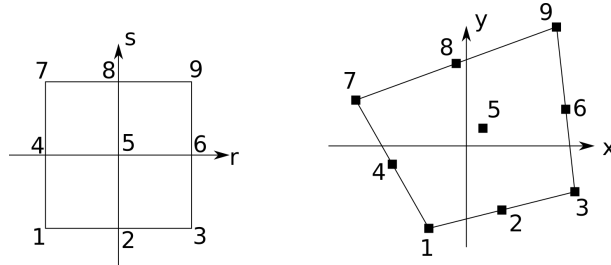
$$\mathbf{J} = \begin{pmatrix} A_1 & B_1 \\ A_3 & B_3 \end{pmatrix} = \begin{pmatrix} \frac{h_x}{2} & 0 \\ 0 & \frac{h_y}{2} \end{pmatrix}$$

The determinant is then $C = \frac{h_x h_y}{4}$ and the inverse:

$$\mathbf{J}^{-1} = \frac{1}{C} \begin{pmatrix} B_3 & -B_1 \\ -A_3 & A_1 \end{pmatrix} = \frac{4}{h_x h_y} \begin{pmatrix} \frac{h_y}{2} & 0 \\ 0 & \frac{h_x}{2} \end{pmatrix} = \begin{pmatrix} \frac{2}{h_x} & 0 \\ 0 & \frac{2}{h_y} \end{pmatrix}$$

Remark. Hua [599] (1990) has published analytical inverse transformation for quadrilateral isoparametric elements, i.e. how to compute \mathbf{J}^{-1} as a function of space coordinates and not just at the quadrature points.

7.6.4 Biquadratic mapping of a straight-edge face Q_2 element



The reference element now contains 9 nodes: 1,3,7,9 are the corners, nodes 2,4,6,8 are the mid-face points and node 5 is in the middle¹⁸. The mapping from the (r, s) space to the (x, y) space is then as follows:

$$\begin{aligned} \begin{pmatrix} x(r, s) \\ y(r, s) \end{pmatrix} &= \mathcal{N}_1(r, s) \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} + \mathcal{N}_2(r, s) \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} + \mathcal{N}_3(r, s) \begin{pmatrix} x_3 \\ y_3 \end{pmatrix} + \mathcal{N}_4(r, s) \begin{pmatrix} x_4 \\ y_4 \end{pmatrix} \\ &+ \mathcal{N}_5(r, s) \begin{pmatrix} x_5 \\ y_5 \end{pmatrix} + \mathcal{N}_6(r, s) \begin{pmatrix} x_6 \\ y_6 \end{pmatrix} + \mathcal{N}_7(r, s) \begin{pmatrix} x_7 \\ y_7 \end{pmatrix} + \mathcal{N}_8(r, s) \begin{pmatrix} x_8 \\ y_8 \end{pmatrix} \\ &+ \mathcal{N}_9(r, s) \begin{pmatrix} x_9 \\ y_9 \end{pmatrix} \end{aligned}$$

¹⁸Note that this numbering is quite arbitrary

where the Q_2 basis functions have been obtained in Section 5.3.2:

$$\begin{aligned}
\mathcal{N}_1(r, t) &= 0.5r(r-1)0.5t(t-1) \\
\mathcal{N}_2(r, t) &= (1-r^2)0.5t(t-1) \\
\mathcal{N}_3(r, t) &= 0.5r(r+1)0.5t(t-1) \\
\mathcal{N}_4(r, t) &= 0.5r(r-1)(1-t^2) \\
\mathcal{N}_5(r, t) &= (1-r^2)(1-t^2) \\
\mathcal{N}_6(r, t) &= 0.5r(r+1)(1-t^2) \\
\mathcal{N}_7(r, t) &= 0.5r(r-1)0.5t(t+1) \\
\mathcal{N}_8(r, t) &= (1-r^2)0.5t(t+1) \\
\mathcal{N}_9(r, t) &= 0.5r(r+1)0.5t(t+1)
\end{aligned}$$

```

x1=-1                ; y1=-2
x3=3                 ; y3=-1
x9=2                 ; y9=2
x7=-3                ; y7=1
x2=0.5*(x1+x3)       ; y2=0.5*(y1+y3)
x4=0.5*(x1+x7)       ; y4=0.5*(y1+y7)
x6=0.5*(x3+x9)       ; y6=0.5*(y3+y9)
x8=0.5*(x7+x9)       ; y8=0.5*(y7+y9)
x5=0.25*(x1+x3+x7+x9) ; y5=0.25*(y1+y3+y7+y9)

npts=10000
r=np.zeros(npts, dtype=np.float64)
s=np.zeros(npts, dtype=np.float64)
xQ1=np.zeros(npts, dtype=np.float64)
yQ1=np.zeros(npts, dtype=np.float64)
xQ2=np.zeros(npts, dtype=np.float64)
yQ2=np.zeros(npts, dtype=np.float64)

for i in range(0, npts):
    # compute random r, s coordinates
    r[i]=random.uniform(-1., +1)
    s[i]=random.uniform(-1., +1)
    # compute Q2 basis function values at r, s
    N1= 0.5*r[i]*(r[i]-1.) * 0.5*s[i]*(s[i]-1.)
    N2= (1.-r[i]**2) * 0.5*s[i]*(s[i]-1.)
    N3= 0.5*r[i]*(r[i]+1.) * 0.5*s[i]*(s[i]-1.)
    N4= 0.5*r[i]*(r[i]-1.) * (1.-s[i]**2)
    N5= (1.-r[i]**2) * (1.-s[i]**2)
    N6= 0.5*r[i]*(r[i]+1.) * (1.-s[i]**2)
    N7= 0.5*r[i]*(r[i]-1.) * 0.5*s[i]*(s[i]+1.)
    N8= (1.-r[i]**2) * 0.5*s[i]*(s[i]+1.)
    N9= 0.5*r[i]*(r[i]+1.) * 0.5*s[i]*(s[i]+1.)
    # compute x, y coordinates
    xQ2[i]=N1*x1+N2*x2+N3*x3+N4*x4+N5*x5+N6*x6+N7*x7+N8*x8+N9*x9
    yQ2[i]=N1*y1+N2*y2+N3*y3+N4*y4+N5*y5+N6*y6+N7*y7+N8*y8+N9*y9
    # compute Q1 basis function values at r, s
    N1=0.25*(1-r[i])*(1-s[i])
    N2=0.25*(1+r[i])*(1-s[i])
    N3=0.25*(1+r[i])*(1+s[i])
    N4=0.25*(1-r[i])*(1+s[i])
    # compute x, y coordinates
    xQ1[i]=N1*x1+N2*x3+N3*x9+N4*x7
    yQ1[i]=N1*y1+N2*y3+N3*y9+N4*y7

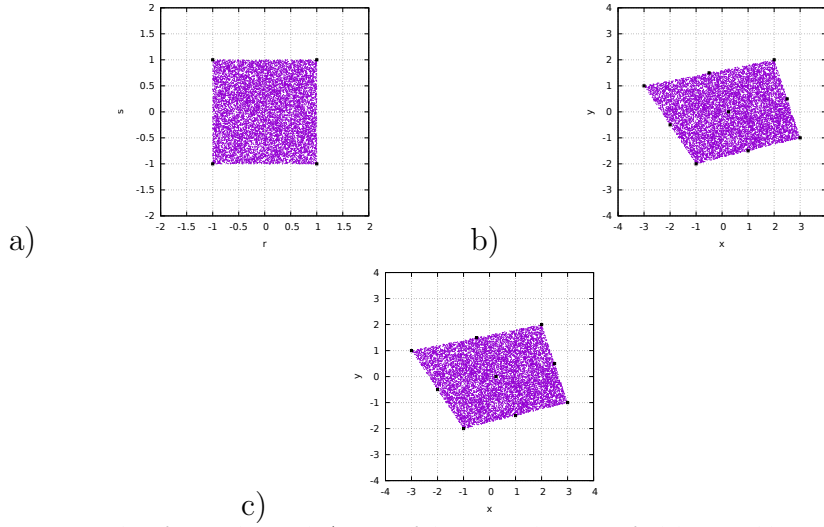
np.savetxt('rs.ascii', np.array([r, s]).T)

```



```
np.savetxt('xyQ1.ascii',np.array([xQ1,yQ1]).T)
np.savetxt('xyQ2.ascii',np.array([xQ2,yQ2]).T)
```

The code is available in /images/mappings/biquadratic Note that the coordinates of point 5 are defined being those of the barycenter of the quadrilateral. More on this choice later.

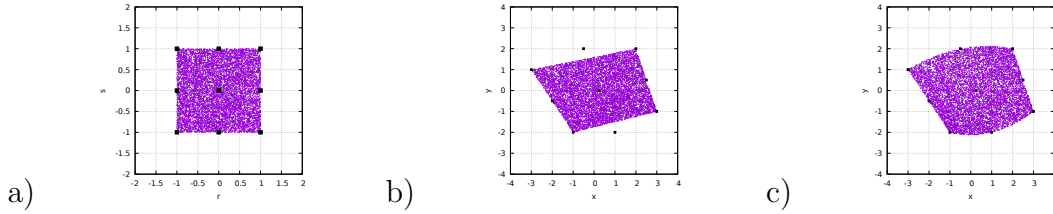


a) 10,000 random points in the reference element; b,c) image of these points by means of a bilinear and biquadratic mapping respectively.

When the sides of the element are straight we see that a Q_1 mapping is sufficient.

7.6.5 Biquadratic mapping of a not-so straight-line face Q_2 element

We now carry out the same exercise as before but nodes 2 and 8 are no more in the middle of nodes 1-3 and 7-9 respectively. The code is available in /images/mappings/biquadratic2.



a) 10,000 random points in the reference element; b,c) image of these points by means of a bilinear and biquadratic mapping respectively.

In this case we see that the Q_2 mapping manages to better capture the 'real' shape of the element. Since nodes 2 and 8 have moved, we could now ask ourselves where we should place node 5? In this example we set it as follows but it is somewhat arbitrary.

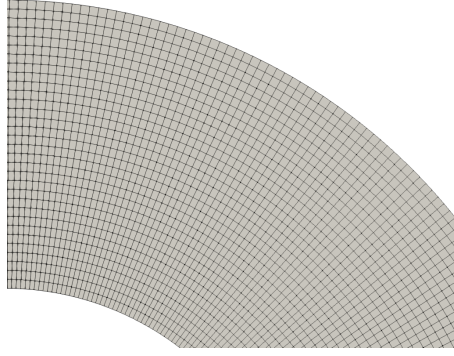
$$x5 = (x1 + x2 + x3 + x4 + x6 + x7 + x8 + x9) / 8.$$

$$y5 = (y1 + y2 + y3 + y4 + y6 + y7 + y8 + y9) / 8.$$

We will come back to this later.

7.6.6 Bilinear, biquadratic and bicubic mapping in an annulus

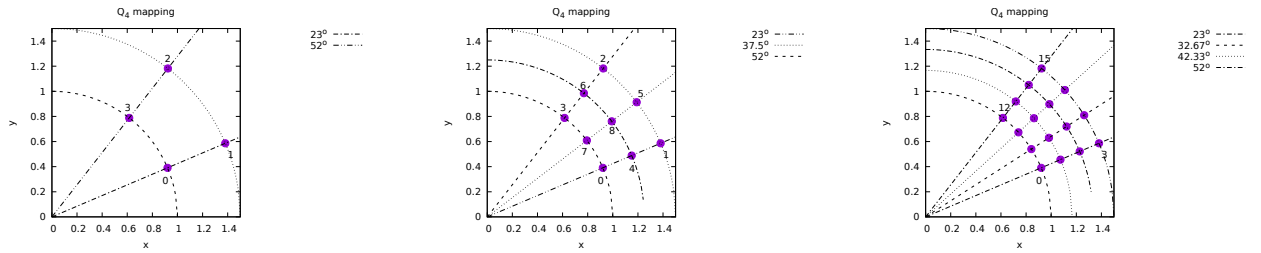
In the light of what precedes, we can now ask ourselves how this translates to a real geodynamic case. Let us then consider the case of an annular domain, a cross section of a hollow sphere. When using quadrilateral elements, the mesh will look similar to this:



We here focus on Q_1 , Q_2 and Q_3 mappings. We single out an element, and arbitrarily define it as follows in polar coordinates:

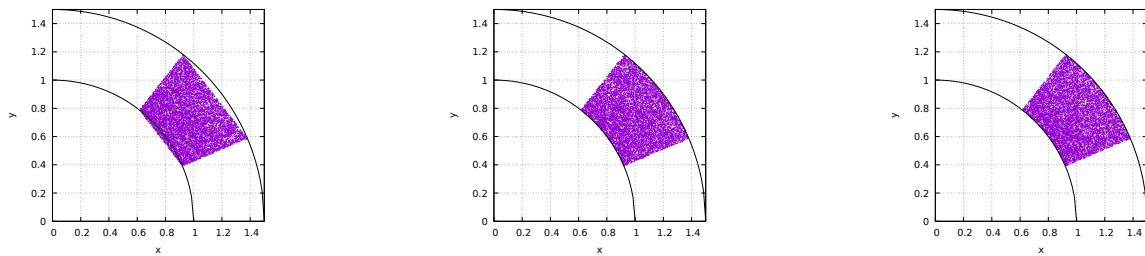
```
theta1=23./180.*np.pi
theta2=52./180.*np.pi
R1=1.
R2=1.5
```

The Q_1 mapping requires four points, the Q_2 nine points and the Q_3 sixteen points. The code used in the following is available at `./images/mappings/curved/`. These are placed equidistantly in the r, θ coordinate system, as shown hereunder:



Left to right: position of the nodes for the Q_1 , Q_2 and Q_3 mappings. Q_4 is not shown.

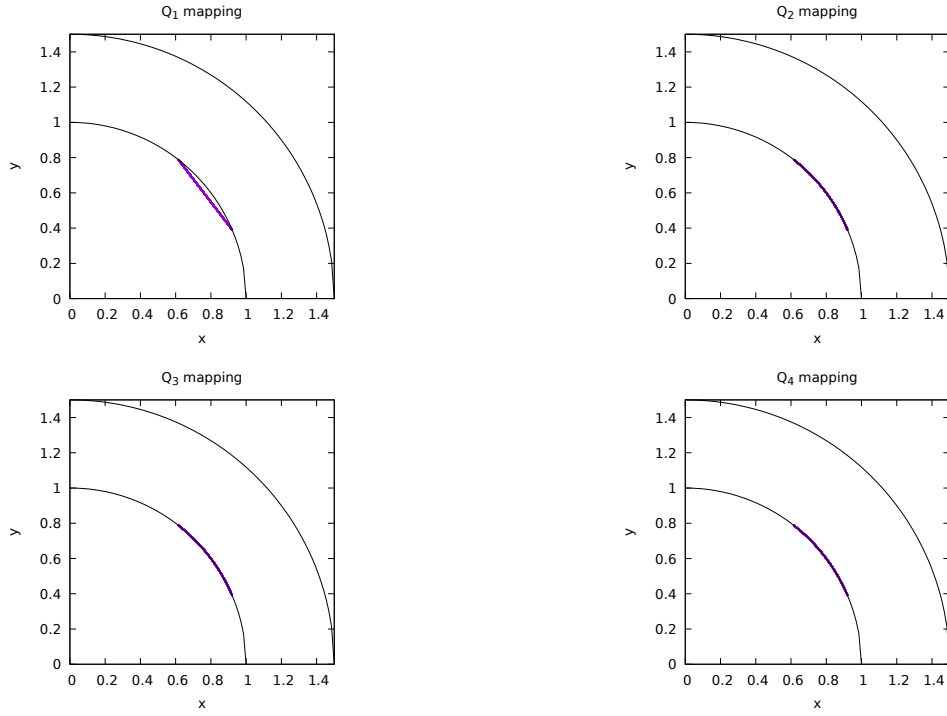
As before, we randomly shoot 10,000 points inside the reference element and map these out in the x, y space. Resulting swarms of points are shown in the following figures:



Left to right: position of the mapped points for the Q_1 , Q_2 and Q_3 mappings. Q_4 is not shown.

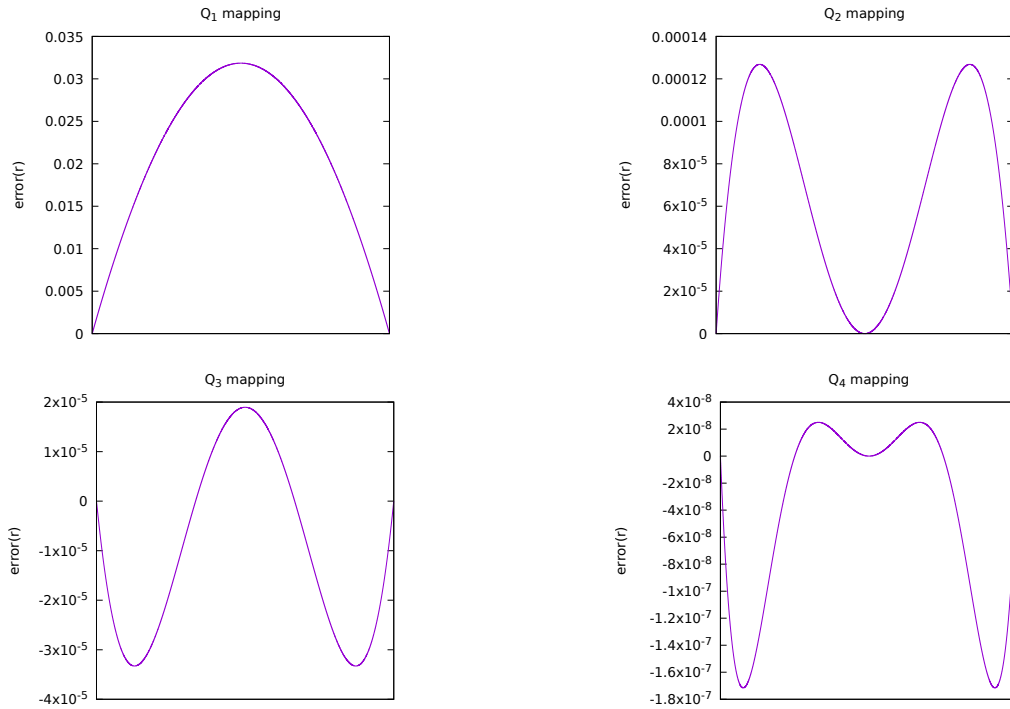
The image of a square with a Q_1 mapping is obviously a quadrilateral so that it looks like quite a few points land outside of the domain $R_1 \leq r \leq R_2$. Note that points are well within $23^\circ \leq \theta \leq 52^\circ$, which can simply be explained by the fact that the faces of the element joining R_1 to R_2 are straight lines.

However, it looks like the biquadratic and bicubic mappings are doing a much better job at mapping the region of space $R_1 \leq r \leq R_2$. In order to characterise this better, we now place 10,000 points on the bottom face of the reference element (i.e. $s = -1$) and once again compute their coordinates in the the x, y space:



Position of the mapped points for the Q_1 , Q_2 , Q_3 and Q_4 mappings.

For each point i we now compute the distance r_i to the origin, which, if the mapping was perfect, would be exactly equal to $R_1 = 1$. On the following plots are shown the error $r_i - 1$ for all points, from $r = -1$ to $r = +1$.



Radius error of the mapped points for the Q_1 , Q_2 , Q_3 and Q_4 mappings.

We see that the amplitude of the error decreases with the order of the mapping used, which is why for instance ASPECT uses a Q_4 mapping by default¹⁹. Actually, in this particular case, the equation which describes the circle is not a polynomial so that no high-order mapping will ever be able to *exactly* represent the curved boundary of the element!

¹⁹I find it also quite striking that the Q_4 mapping outperforms the Q_3 one by two orders of magnitude...

Another interesting point to keep in mind is that the location of the quadrature points in the x, y space is also determined by the mapping used, which can have consequences on the accuracy of the integration and it will be reflected (for instance) on the error convergence rate.

As already mentioned previously, the coordinates of the nodes of the element in the x, y are uniquely determined when they are on the convex hull of the element (for instance nodes 0-7 for Q_2) but we need to choose the position of the last nodes which are inside the element. Unfortunately, this choice is not neutral.

Finally, we can explore the importance of the mapping in combination with numerical quadrature. For each mapping we compute the area of the element by means of a 3x3, 4x4 or 5x5 quadrature.

```
*****Q1*****
nqperdim= 3 0.3030060126539606 rel. error -0.04215361698430029
nqperdim= 4 0.3030060126539606 rel. error -0.04215361698430012 ~ 4%
nqperdim= 5 0.3030060126539606 rel. error -0.04215361698430012
*****Q2*****
nqperdim= 3 0.3162980025394154 rel. error -0.00013569026611326453
nqperdim= 4 0.3162980025394155 rel. error -0.00013569026611308905 ~ 0.01%
nqperdim= 5 0.3162980025394154 rel. error -0.00013569026611326453
*****Q3*****
nqperdim= 3 0.3163472223929359 rel. error 1.9900899402587318e-05
nqperdim= 4 0.316347222392936 rel. error 1.9900899402938278e-05 ~ 0.002%
nqperdim= 5 0.316347222392936 rel. error 1.9900899402938278e-05
*****Q4*****
nqperdim= 3 0.3163409410866220 rel. error 4.477021014282521e-08
nqperdim= 4 0.3163409541901677 rel. error 8.619243716974044e-08 ~ 0.000008%
nqperdim= 5 0.316340954190168 rel. error 8.619243804713484e-08
```

Here again the Q_4 mapping makes quite the difference.

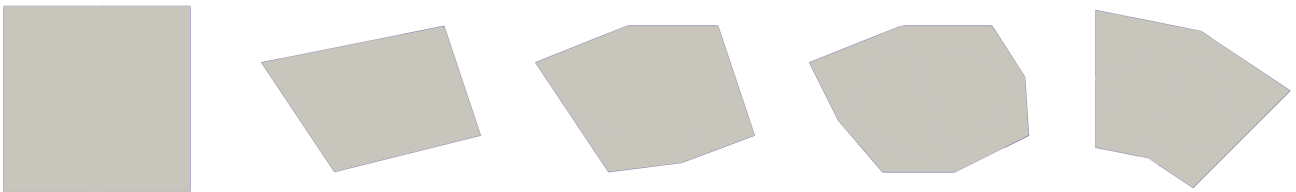
7.6.7 Biquadratic mapping - the middle node conundrum

Python code at `images/mappings/biquadratic3`.

As mentioned before, unless the element is a straight-edge quadrilateral, determining the (best) position of the middle node is not trivial. Or is it?

```
4--7--3
|      |
8  9  6  (reference element)
|      |
1--5--2
```

We will here consider 5 different elements:



From left to right: element 0,1,2,3,4.

We can think of multiple ways to come up with the 'center' of the element, i.e. the location of point I.

- center=0:

$$x_9 = (x_1 + x_2 + x_3 + x_4)/4 \quad y_9 = (y_1 + y_2 + y_3 + y_4)/4$$

- center=1:

$$x_9 = (x_1 + x_2 + x_3 + x_4 + x_5 + x_6 + x_7 + x_8)/8 \quad y_9 = (y_1 + y_2 + y_3 + y_4 + y_5 + y_6 + y_7 + y_8)/8$$

- center=2:

$$x_9 = (x_1 + x_2 + x_3 + x_4 + 3x_5 + 3x_6 + 3x_7 + 3x_8)/16. \quad y_9 = (y_1 + y_2 + y_3 + y_4 + 3y_5 + 3y_6 + 3y_7 + 3y_8)/16.$$

- center=3: (only element=4)

$$x_9 = \frac{1}{2}(R_1 + R_2) \cos(3\pi/8) \quad y_9 = \frac{1}{2}(R_1 + R_2) \sin(3\pi/8)$$

- center=4: I is the center of mass. The element is defined by $R_1 < r < R_2$ and $\theta_1 < \theta < \theta_2$.

We need to compute²⁰

$$\begin{aligned}
\vec{R} &= \frac{1}{M} \int \vec{r} \rho(\vec{r}) dV \\
&= \frac{1}{M} \rho_0 \int \vec{r} dV \\
&= \frac{1}{M} \frac{M}{V} \int \vec{r} dV \\
&= \frac{1}{V} \int \vec{r} dV \\
&= \frac{1}{V} \int \begin{pmatrix} x \\ y \end{pmatrix} dV \\
&= \frac{1}{V} \int \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix} dV \\
&= \frac{1}{V} \int_{R_1}^{R_2} \int_{\theta_1}^{\theta_2} \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix} r dr d\theta \\
&= \frac{1}{\frac{1}{2}(R_2^2 - R_1^2)(\theta_2 - \theta_1)} \frac{1}{3} (R_2^3 - R_1^3) \begin{pmatrix} \sin \theta_2 - \sin \theta_1 \\ -\cos \theta_2 + \cos \theta_1 \end{pmatrix} \\
&\simeq \begin{pmatrix} 0.5801028000103104 \\ 1.4004920473554983 \end{pmatrix}
\end{aligned} \tag{7.126}$$

which corresponds to $r = 1.5158816686291174$ and $\theta = 67.5^\circ = 3\pi/8$.

- center=5: variable position

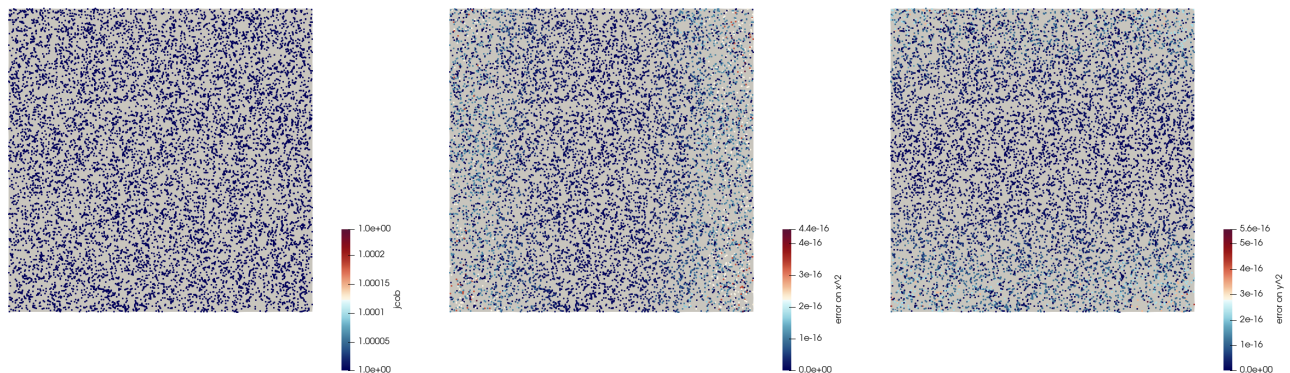
isoparametric mapping.

At each point (r, s) we compute the error $|\sum_i N_i(r, s)x_i^2 - (\sum_i N_i(r, s)x_i)^2|$.

position of edges (setting $r=\pm 1$, $s=\pm 1$) independent of position of middle node since shape functions are zero there

area indep of position middle node ?

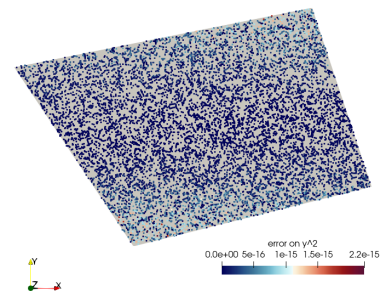
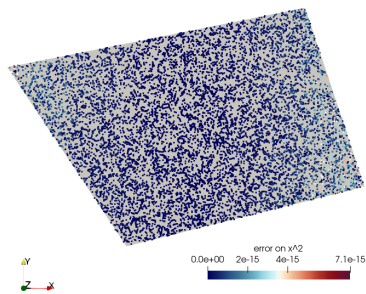
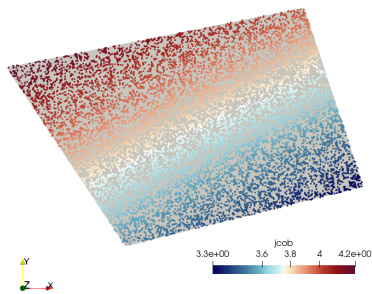
Element 0 In this case all only center=0,1,2,4 are applicable but they all lead to the same point I with $x_I = 0, y_I = 0$. This means that the position of quadrature points is also independent of the center parameter.



10,000 points at random.

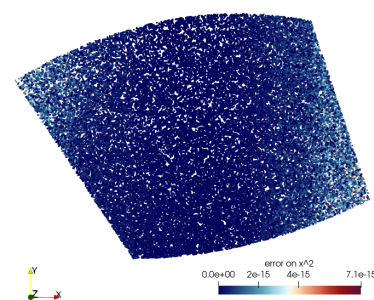
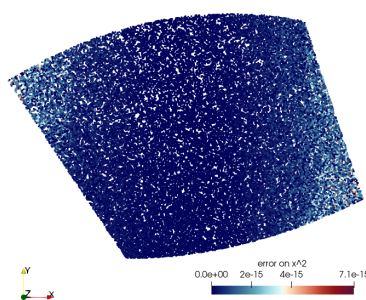
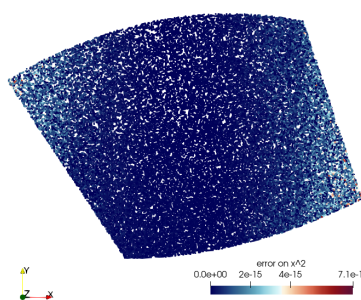
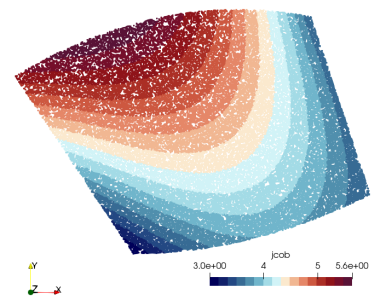
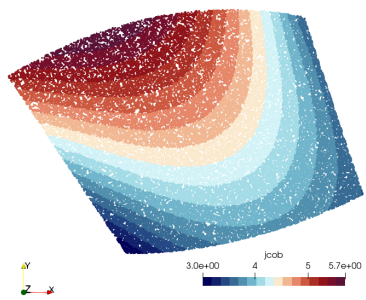
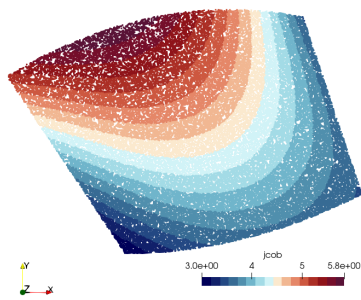
²⁰https://en.wikipedia.org/wiki/Center_of_mass

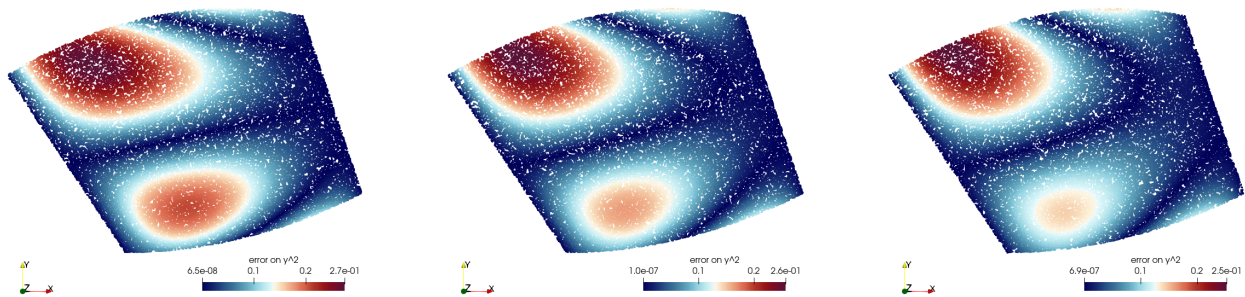
Element 1 In this case all only center=0,1,2,4 are applicable but they all lead to the same point I with $x_I = 0, y_I = 0$. This means that the position of quadrature points is also independent of the center parameter.



10,000 points at random.

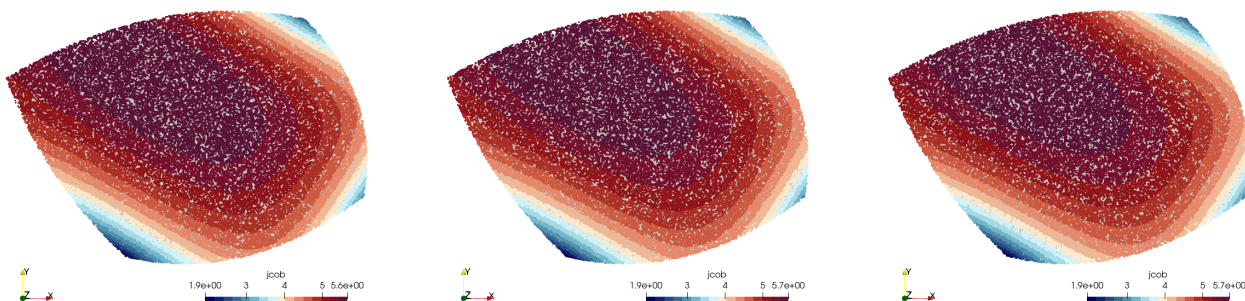
Element 2 .





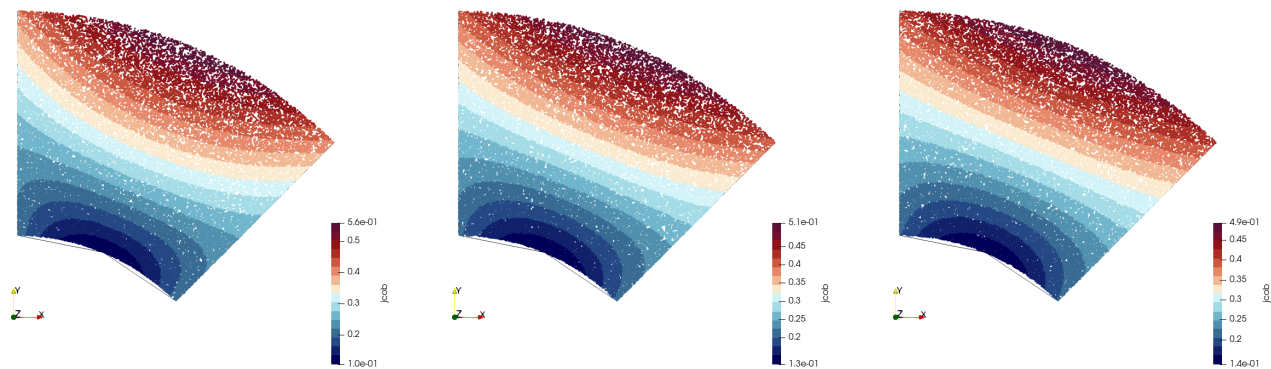
50,000 points at random. From left to right: center=0,1,2.

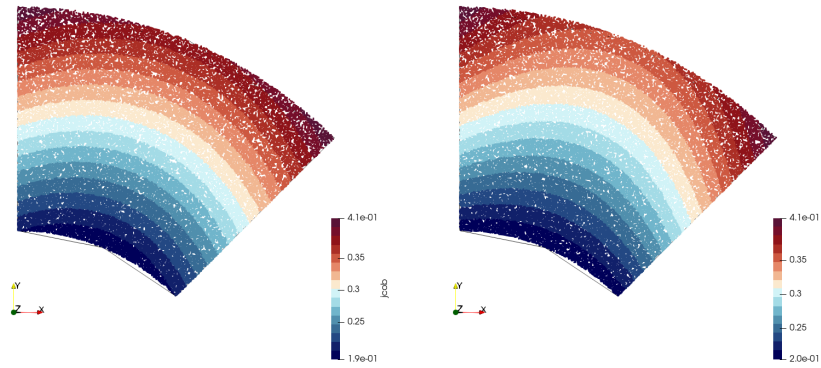
Element 3 .



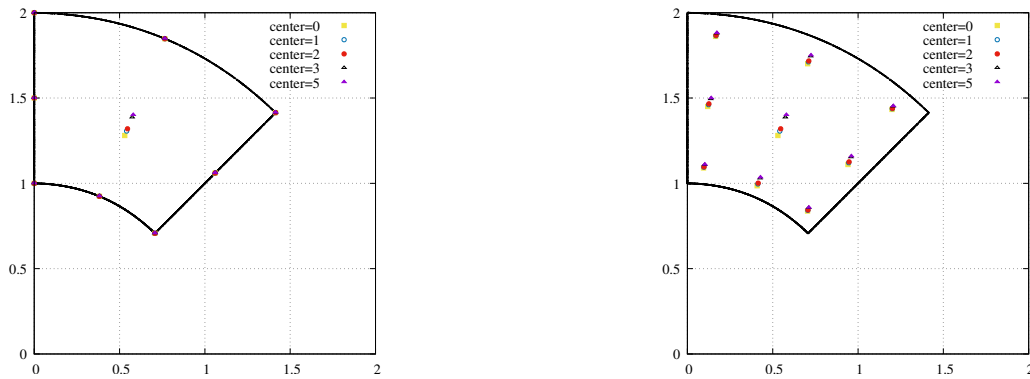
50,000 points at random. From left to right: center=0,1,2.

Element 4





50,000 points at random. From left to right: center=0,1,2,3,4.



Left: position of the nodes. Right position of quadrature points with nqperdim=3.

Area does not depend on position of middle node?!

📖 Relevant Literature

- K.Y. Yuan, Y.S. Huang, H.T. Yang, and T.H.H. Pian. “The inverse mapping and distortion measures for 8-node hexahedral isoparametric elements”. In: *Computational Mechanics* 14 (1994), pp. 189–199

7.6.8 The Double Jacobian approach

What follows is 90% borrowed from Morgan, Taramón, and Hasenclever [906] (2020) with slight changes in the notations.

The basic idea behind this approach is to compute the local to Cartesian mapping as a two-stage process, hence the name “Double Jacobian”.

1. The first stage maps from local to polar/spherical coordinates and back. This mapping is typically to a straight-edged polar or spherical element for which the Jacobian partial derivatives are constant within the element. The mapping and its inverse are given by straightforward analytical matrix expressions.
2. The second stage maps from polar/spherical to Cartesian coordinates (and back) and is also a simple analytical mapping.

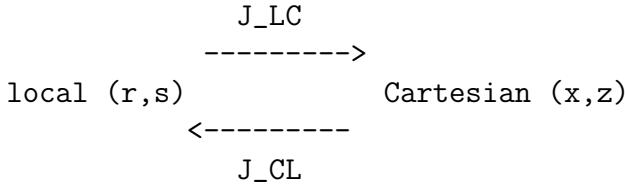
The net Jacobian from local finite element coordinates to a cylindrical or spherical element in Cartesian geometry (or any other analytically mapped geometry) is simply the matrix product of two easy-to-compute inverse Jacobian matrices. Because the net Jacobian has an analytical form, it can be more rapidly computed than a general isoparametric or superparametric finite element mapping.

Remark: in what follows I will denote the polar coordinates by ρ, θ so as to avoid confusion with the reduced/local coordinates r, s .

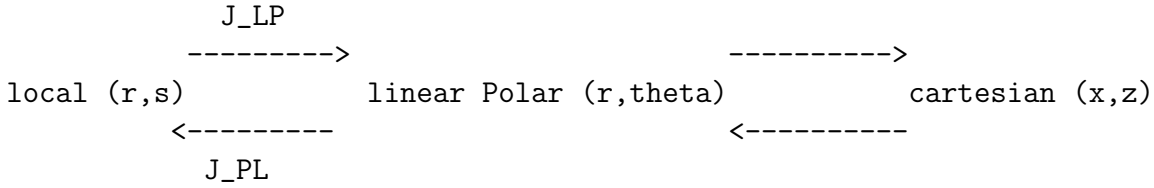
As we have seen so far the standard approach is as follows: The relationship between the derivatives of the basis functions with respect to the Cartesian coordinates and the same derivatives with respect to the local coordinates is given by:

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial r}{\partial x} & \frac{\partial s}{\partial x} \\ \frac{\partial r}{\partial z} & \frac{\partial s}{\partial z} \end{pmatrix}}_{J_{LC}} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix}$$

where \mathcal{N} are the shape functions, i is the local node numbering of the element, (x, z) are the Cartesian coordinates, and (r, s) are the local coordinates within the reference element.



DJ: First Jacobian



The Double Jacobian approach uses the standard finite element approach to first map from local to linear polar coordinates (i.e. ρ, θ). The first Jacobian for a cylindrical (polar) mapping is analogous to the standard Jacobian, where x and z are now changed to θ and ρ , respectively, θ being the angle measured from the positive z -axis in clockwise direction (i.e. the colatitude) and ρ being the radius. The global derivatives may be expressed in matrix form as

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial \theta} \\ \frac{\partial \mathcal{N}_i}{\partial \rho} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial r}{\partial \theta} & \frac{\partial s}{\partial \theta} \\ \frac{\partial r}{\partial \rho} & \frac{\partial s}{\partial \rho} \end{pmatrix}}_{J_{LP}} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} \quad (7.127)$$

where \mathcal{N}_i are the shape functions and J_{LP} is the Jacobian of the transformation from local to polar coordinates.

The derivatives of these shape functions with respect to local coordinates r, s (i.e. the rhs vector) can be computed explicitly since the basis functions are chosen/built for a given element type and formulated as a function of the local coordinates. However, the terms of the Jacobian J_{LP} cannot be directly computed since explicit expressions for $r(\theta, \rho)$ and $s(\theta, \rho)$ do not exist. A wonderful “trick” in finite element programming (discovered by Bruce Irons in the mid-60s [379]) is to make use of the

inverse coordinate transformation

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \theta}{\partial r} & \frac{\partial \rho}{\partial r} \\ \frac{\partial \theta}{\partial s} & \frac{\partial \rho}{\partial s} \end{pmatrix}}_{\mathbf{J}_{PL}} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial \theta} \\ \frac{\partial \mathcal{N}_i}{\partial \rho} \end{pmatrix} \quad (7.128)$$

where \mathbf{J}_{PL} is the Jacobian of the transformation from polar to local coordinates.

From Equations (7.127) and (7.128), we have $\mathbf{J}_{LP} = \mathbf{J}_{PL}^{-1}$. Using Cramer's rule, the inverse of the Jacobian from polar to local coordinates is given by

$$\mathbf{J}_{LP} = \mathbf{J}_{PL}^{-1} = \frac{1}{|\mathbf{J}_{PL}|} \begin{pmatrix} \frac{\partial \rho}{\partial s} & -\frac{\partial \rho}{\partial r} \\ -\frac{\partial \theta}{\partial s} & \frac{\partial \theta}{\partial r} \end{pmatrix}$$

The polar coordinates for each element are related with local coordinates through the shape functions:

$$\theta(r, s) = \sum_{i=1}^m \mathcal{N}_i(r, s) \theta_i \quad (7.129)$$

$$\rho(r, s) = \sum_{i=1}^m \mathcal{N}_i(r, s) \rho_i \quad (7.130)$$

where m is the number of nodes in the element. Then

$$\frac{\partial \theta}{\partial r}(r, s) = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, s) \theta_i \quad (7.131)$$

$$\frac{\partial \theta}{\partial s}(r, s) = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial s}(r, s) \theta_i \quad (7.132)$$

$$\frac{\partial \rho}{\partial r}(r, s) = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, s) \rho_i \quad (7.133)$$

$$\frac{\partial \rho}{\partial s}(r, s) = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial s}(r, s) \rho_i \quad (7.134)$$

DJ: Second jacobian

The second Jacobian in the DJ method is the analytical mapping from polar coordinates to Cartesian coordinates. The derivatives expressed in matrix form are given by

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \theta}{\partial x} & \frac{\partial \rho}{\partial x} \\ \frac{\partial \theta}{\partial z} & \frac{\partial \rho}{\partial z} \end{pmatrix}}_{\mathbf{J}_{PC}} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial \theta} \\ \frac{\partial \mathcal{N}_i}{\partial \rho} \end{pmatrix} \quad (7.135)$$

where \mathbf{J}_{PC} is the Jacobian from polar to Cartesian coordinates. The analytical expressions for $\theta(x, z)$ and $\rho(x, z)$ are known; however, it is again easier to use the inverse transformation

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial \theta} \\ \frac{\partial \mathcal{N}_i}{\partial \rho} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial x}{\partial \theta} & \frac{\partial z}{\partial \theta} \\ \frac{\partial x}{\partial \rho} & \frac{\partial z}{\partial \rho} \end{pmatrix}}_{\mathbf{J}_{CP}} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} \quad (7.136)$$

where \mathbf{J}_{CP} is the Jacobian from Cartesian to polar coordinates. From Equations (7.135) and (7.136), $\mathbf{J}_{PC} = \mathbf{J}_{CP}^{-1}$. The inverse of the Jacobian from Cartesian to polar coordinates is given by

$$\mathbf{J}_{PC} = \mathbf{J}_{CP}^{-1} = \frac{1}{|\mathbf{J}_{CP}|} \begin{pmatrix} \frac{\partial z}{\partial \rho} & -\frac{\partial z}{\partial \theta} \\ -\frac{\partial x}{\partial \rho} & \frac{\partial x}{\partial \theta} \end{pmatrix}$$

Cartesian coordinates are related with polar coordinates by:

$$x(\theta, \rho) = \rho \sin \theta \quad (7.137)$$

$$z(\theta, \rho) = \rho \cos \theta \quad (7.138)$$

so that

$$\frac{\partial x}{\partial \theta}(\theta, \rho) = \rho \cos \theta \quad (7.139)$$

$$\frac{\partial x}{\partial \rho}(\theta, \rho) = \sin \theta \quad (7.140)$$

$$\frac{\partial z}{\partial \theta}(\theta, \rho) = -\rho \sin \theta \quad (7.141)$$

$$\frac{\partial z}{\partial \rho}(\theta, \rho) = \cos \theta \quad (7.142)$$

The inverse of the Jacobian from Cartesian coordinates to polar coordinates can then be written as a function of polar coordinates:

$$\boxed{\mathbf{J}_{CP}^{-1} = \frac{1}{\rho} \begin{pmatrix} \cos \theta & \rho \sin \theta \\ -\sin \theta & \rho \cos \theta \end{pmatrix}} \quad (7.143)$$

where θ and ρ are evaluated at each integration point (ip).

Remark: this Jacobian is independent of the choice of basis functions.

DJ: Combining both

Making use of the matrix product of the two inverse Jacobians, global Cartesian derivatives can be expressed as a matrix product of the local derivatives in the local to polar and polar to Cartesian coordinate mappings. Substituting Equation (7.127) into Equation (7.135) yields

$$\begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} = \mathbf{J}_{PC} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial \theta} \\ \frac{\partial \mathcal{N}_i}{\partial \rho} \end{pmatrix} = \mathbf{J}_{PC} \cdot \mathbf{J}_{LP} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} = \mathbf{J}_{CP}^{-1} \cdot \mathbf{J}_{PL}^{-1} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix}$$

Note that the 2D DJ approach will ensure a perfect mapping to the circular-arc edges of the elements of a cylindrical annulus mesh. Any point on the edge of a boundary element of the mesh is mapped to its true position along a circular arc.

The major drawback from this is the fact that (as we will see) the Jacobian is no more a polynomial so special care must be taken with regards to the integration.

Triangles

For linear triangles (P_1) the basis functions are

$$\mathcal{N}_1(r, s) = 1 - r - s \quad (7.144)$$

$$\mathcal{N}_2(r, s) = r \quad (7.145)$$

$$\mathcal{N}_3(r, s) = s \quad (7.146)$$

and their derivatives:

$$\frac{\partial \theta}{\partial r}(r, s) = \theta_2 - \theta_1 = \theta_{21} \quad (7.147)$$

$$\frac{\partial \theta}{\partial s}(r, s) = \theta_3 - \theta_1 = \theta_{31} \quad (7.148)$$

$$\frac{\partial \rho}{\partial r}(r, s) = \rho_2 - \rho_1 = \rho_{21} \quad (7.149)$$

$$\frac{\partial \rho}{\partial s}(r, s) = \rho_3 - \rho_1 = \rho_{31} \quad (7.150)$$

In this case we have

$$\mathbf{J}_{LP} = \mathbf{J}_{PL}^{-1} = \frac{1}{\theta_{21}\rho_{31} - \rho_{21}\theta_{31}} \begin{pmatrix} \rho_{31} & -\rho_{21} \\ -\theta_{31} & \theta_{21} \end{pmatrix} \quad (7.151)$$

Substituting Equations (7.151) and (7.143) into the equation above, the analytical expression for this mapping is

$$\begin{aligned} \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} &= \frac{1}{\rho} \begin{pmatrix} \cos \theta & \rho \sin \theta \\ -\sin \theta & \rho \cos \theta \end{pmatrix} \cdot \frac{1}{\theta_{21}\rho_{31} - \rho_{21}\theta_{31}} \begin{pmatrix} \rho_{31} & -\rho_{21} \\ -\theta_{31} & \theta_{21} \end{pmatrix} \\ &= \frac{1}{\theta_{21}\rho_{31} - \rho_{21}\theta_{31}} \begin{pmatrix} \frac{\rho_{31}}{\rho} \cos \theta - \theta_{31} \sin \theta & -\frac{\rho_{21}}{\rho} \cos \theta + \theta_{21} \sin \theta \\ -\frac{\rho_{31}}{\rho} \sin \theta - \theta_{31} \cos \theta & \frac{\rho_{21}}{\rho} \sin \theta + \theta_{21} \cos \theta \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} \end{aligned} \quad (7.152)$$

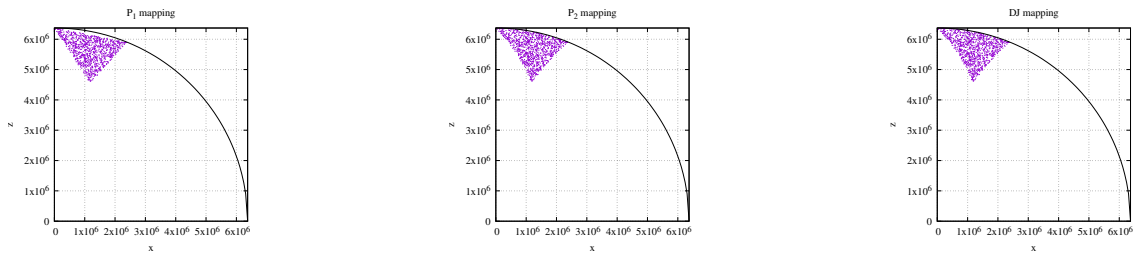
which is Eq. (19) of Morgan, Taramón, and Hasenclever [906]. This expression could be problematic if $\rho = 0$ but since it will be evaluated at the quadrature points this case is extremely unlikely.

Based on Fig 2 of Morgan, Taramón, and Hasenclever [906], we consider an annulus of outer diameter 6371km. Visually from 2A we infer $\rho_1 = 4700\text{km}$, $\rho_2 = \rho_3 = 6371\text{km}$, and $\theta_1 = 14.5$, $\theta_2 = 22$ and $\theta_3 = 0$.

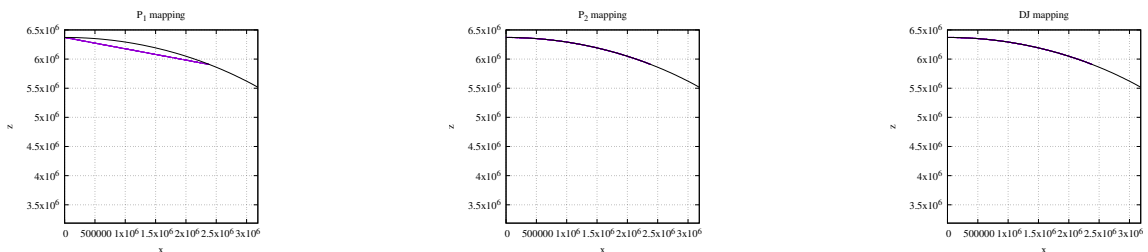
If quadratic basis functions are used, we'll define $r_4 = \frac{1}{2}(r_1 + r_2)$, $r_5 = \frac{1}{2}(r_2 + r_3)$, $r_6 = \frac{1}{2}(r_3 + r_1)$, and likewise for the $\theta_{4,5,6}$ values.

We'll then proceed to generate 1000 random points in the reference triangle and plot their image in the Cartesian plane, either using the DJ method, a linear mapping P_1 or a quadratic mapping P_2 .

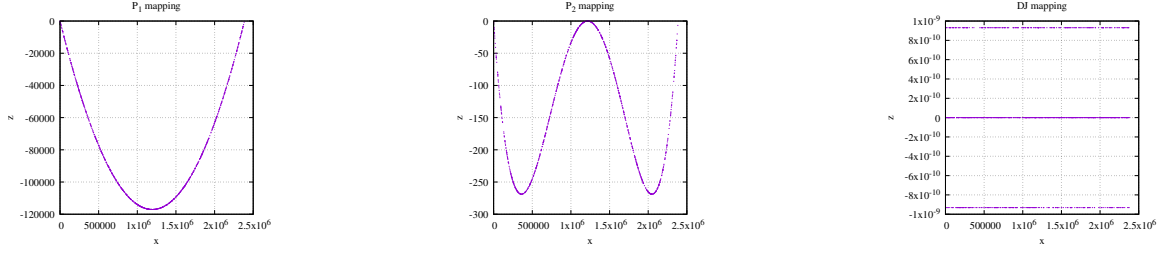
Unsurprisingly the P_1 mapping yields a triangle which 2-3 edge does not conform to the edge of the domain. The P_2 mapping does a much better job and so does the DJ mapping.



Let us now generate 1000 points on the 2-3 edge of the reference triangle and study their image with the three mappings:



We cannot see any different between P2 and DJ. Let us now plot the error, i.e. the distance of the image point to the true surface $\rho = 6371\text{km}$ as a function of x :



We find that the maximum error for P_1 is about 120km, the error for P_2 is about 260m, and the error for DJ is effectively about 10^{-9}m .

Quadrilaterals

For linear quadrilaterals (Q_1) the basis functions are:

$$\mathcal{N}_1(r, s) = \frac{1}{4}(1-r)(1-s) \quad (7.153)$$

$$\mathcal{N}_2(r, s) = \frac{1}{4}(1+r)(1-s) \quad (7.154)$$

$$\mathcal{N}_3(r, s) = \frac{1}{4}(1+r)(1+s) \quad (7.155)$$

$$\mathcal{N}_4(r, s) = \frac{1}{4}(1-r)(1+s) \quad (7.156)$$

$$\frac{\partial \theta}{\partial r}(r, s) = -\frac{1}{4}\theta_1 + \frac{1}{4}\theta_2 + \frac{1}{4}\theta_3 - \frac{1}{4}\theta_4 = \frac{1}{4}(-\theta_1 + \theta_2 + \theta_3 - \theta_4) = \tilde{\theta}_{1234} \quad (7.157)$$

$$\frac{\partial \theta}{\partial s}(r, s) = -\frac{1}{4}\theta_1 - \frac{1}{4}\theta_2 + \frac{1}{4}\theta_3 + \frac{1}{4}\theta_4 = \frac{1}{4}(-\theta_1 - \theta_2 + \theta_3 + \theta_4) = \bar{\theta}_{1234} \quad (7.158)$$

$$\frac{\partial \rho}{\partial r}(r, s) = -\frac{1}{4}\rho_1 + \frac{1}{4}\rho_2 + \frac{1}{4}\rho_3 - \frac{1}{4}\rho_4 = \frac{1}{4}(-\rho_1 + \rho_2 + \rho_3 - \rho_4) = \tilde{\rho}_{1234} \quad (7.159)$$

$$\frac{\partial \rho}{\partial s}(r, s) = -\frac{1}{4}\rho_1 - \frac{1}{4}\rho_2 + \frac{1}{4}\rho_3 + \frac{1}{4}\rho_4 = \frac{1}{4}(-\rho_1 - \rho_2 + \rho_3 + \rho_4) = \bar{\rho}_{1234} \quad (7.160)$$

In this case we have

$$\mathbf{J}_{LP} = \mathbf{J}_{PL}^{-1} = \frac{1}{\tilde{\theta}_{1234}\bar{\rho}_{1234} - \bar{\theta}_{1234}\tilde{\rho}_{1234}} \begin{pmatrix} \bar{\rho}_{1234} & -\tilde{\rho}_{1234} \\ -\bar{\theta}_{1234} & \tilde{\theta}_{1234} \end{pmatrix} \quad (7.161)$$

SPECIAL CASE: In the element we have $\theta_1 = \theta_4$, $\theta_2 = \theta_3$, $\theta_2 - \theta_1 = \theta_3 - \theta_4 = \tilde{\theta}$. Likewise, $\rho_1 = \rho_2$, $\rho_3 = \rho_4$ and $\rho_4 - \rho_1 = \rho_3 - \rho_2 = \tilde{\rho}$. In this case:

$$\frac{\partial \theta}{\partial r}(r, s) = \frac{1}{4}(-\theta_1 + \theta_2 + \theta_3 - \theta_4) = \frac{1}{2}\tilde{\theta} \quad (7.162)$$

$$\frac{\partial \theta}{\partial s}(r, s) = \frac{1}{4}(-\theta_1 - \theta_2 + \theta_3 + \theta_4) = 0 \quad (7.163)$$

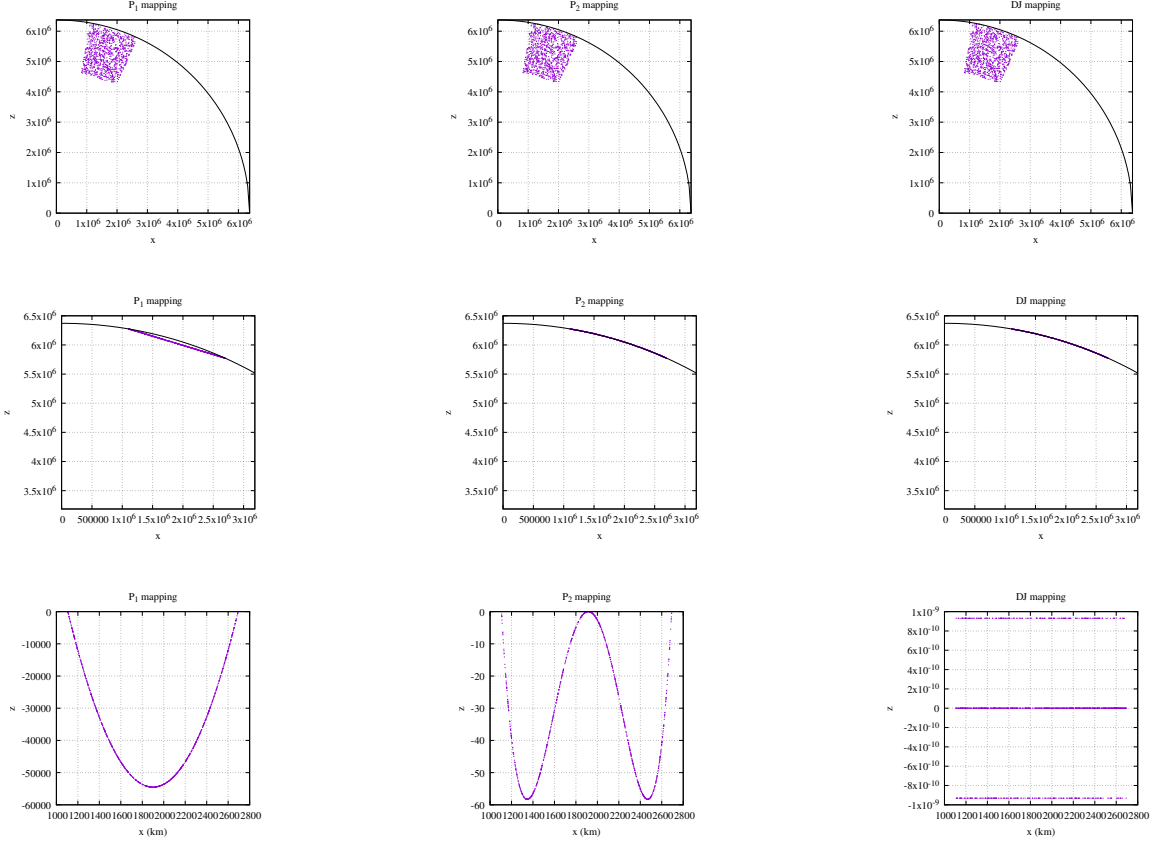
$$\frac{\partial}{\partial r}\rho(r, s) = \frac{1}{4}(-\rho_1 + \rho_2 + \rho_3 - \rho_4) = 0 \quad (7.164)$$

$$\frac{\partial}{\partial s}\rho(r, s) = \frac{1}{4}(-\rho_1 - \rho_2 + \rho_3 + \rho_4) = \frac{1}{2}\tilde{\rho} \quad (7.165)$$

and

$$\mathbf{J}_{LP} = \mathbf{J}_{PL}^{-1} = \frac{1}{\frac{1}{2}\bar{\rho}\frac{1}{2}\tilde{\theta}} \begin{pmatrix} \frac{1}{2}\bar{\rho} & 0 \\ 0 & \frac{1}{2}\tilde{\theta} \end{pmatrix} = 2 \begin{pmatrix} \tilde{\theta}^{-1} & 0 \\ 0 & \bar{\rho}^{-1} \end{pmatrix} \quad (7.166)$$

$$\begin{aligned} \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial x} \\ \frac{\partial \mathcal{N}_i}{\partial z} \end{pmatrix} &= \frac{1}{\rho} \begin{pmatrix} \cos \theta & \rho \sin \theta \\ -\sin \theta & \rho \cos \theta \end{pmatrix} \cdot 2 \begin{pmatrix} \tilde{\theta}^{-1} & 0 \\ 0 & \bar{\rho}^{-1} \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} \\ &= \frac{2}{\rho} \begin{pmatrix} \tilde{\theta}^{-1} \cos \theta & \rho \bar{\rho}^{-1} \sin \theta \\ -\tilde{\theta}^{-1} \sin \theta & \rho \bar{\rho}^{-1} \cos \theta \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial \mathcal{N}_i}{\partial r} \\ \frac{\partial \mathcal{N}_i}{\partial s} \end{pmatrix} \end{aligned} \quad (7.167)$$



7.7 Solving the elastic equations

This will be moved to Section 16.7

NOW BEING REWORKED IN OVERLEAF

In what follows \vec{v} now stands for the displacement vector, i.e. with units of length, not velocity. As before, the displacement inside an element is given by

$$\vec{v}^h(\vec{r}) = \sum_{i=1}^{m_v} N_i(\vec{r}) \vec{v}_i \quad (7.168)$$

where N_i are the polynomial basis functions for the displacement. Pressure does not appear in the equations so this is not a case of mixed FE as for the viscous Stokes flow.

Other notations are sometimes used for Eqs.(7.168):

$$u^h(\vec{r}) = \vec{N} \cdot \vec{u} \quad v^h(\vec{r}) = \vec{N} \cdot \vec{v} \quad w^h(\vec{r}) = \vec{N} \cdot \vec{w} \quad (7.169)$$

where $\vec{v} = (u, v, w)$ and \vec{N} is the vector containing all basis functions evaluated at location \vec{r} :

$$\vec{N}^v = (N_1(\vec{r}), N_2(\vec{r}), N_3(\vec{r}), \dots, N_{m_v}(\vec{r})) \quad (7.170)$$

$$\vec{N}^p = (N_1^p(\vec{r}), N_2^p(\vec{r}), N_3^p(\vec{r}), \dots, N_{m_p}^p(\vec{r})) \quad (7.171)$$

and with

$$\vec{u} = (u_1, u_2, u_3, \dots, u_{m_v}) \quad (7.172)$$

$$\vec{v} = (v_1, v_2, v_3, \dots, v_{m_v}) \quad (7.173)$$

$$\vec{w} = (w_1, w_2, w_3, \dots, w_{m_v}) \quad (7.174)$$

$$(7.175)$$

In three dimensions We start from

$$\boldsymbol{\sigma} = \lambda(\vec{\nabla} \cdot \vec{v})\mathbf{1} + 2\mu\boldsymbol{\varepsilon}$$

where μ is the shear modulus and λ the Lamé parameter.

$$\begin{aligned} \sigma_{xx} &= (\lambda + 2\mu)\varepsilon_{xx} + \lambda\varepsilon_{yy} + \lambda\varepsilon_{zz} \\ \sigma_{yy} &= \lambda\varepsilon_{xx} + (\lambda + 2\mu)\varepsilon_{yy} + \lambda\varepsilon_{zz} \\ \sigma_{zz} &= \lambda\varepsilon_{xx} + \lambda\varepsilon_{yy} + (\lambda + 2\mu)\varepsilon_{zz} \\ \sigma_{xy} &= 2\mu\varepsilon_{xy} \\ \sigma_{xz} &= 2\mu\varepsilon_{xz} \\ \sigma_{yz} &= 2\mu\varepsilon_{yz} \end{aligned} \quad (7.176)$$

or,

$$\vec{\sigma} = \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{zz} \\ 2\varepsilon_{xy} \\ 2\varepsilon_{xz} \\ 2\varepsilon_{yz} \end{pmatrix} = \vec{\varepsilon}$$

The rest of the procedure is pretty straightforward since it follows the same ideas as for the mixed viscous case, except that we here build the \mathbb{K} matrix only as follows:

$$\mathbb{K} = \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{D} \cdot \mathbf{B} dV$$

In two dimensions The above relationships simplify to

$$\sigma_{xx} = (\lambda + 2\mu)\varepsilon_{xx} + \lambda\varepsilon_{yy} \quad (7.177)$$

$$\sigma_{yy} = \lambda\varepsilon_{xx} + (\lambda + 2\mu)\varepsilon_{yy} \quad (7.178)$$

$$\sigma_{xy} = 2\mu\varepsilon_{xy} \quad (7.179)$$

so

$$\vec{\sigma} = \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ 2\varepsilon_{xy} \end{pmatrix} = \vec{\varepsilon}$$

The axisymmetric case

We start from

$$\boldsymbol{\sigma} = \lambda \vec{\nabla} \cdot \vec{u} \mathbf{1} + 2\mu \boldsymbol{\varepsilon}(\vec{u}) \quad (7.180)$$

In cylindrical coordinates the velocity gradient is given by

$$\vec{\nabla} \vec{u} = \begin{pmatrix} \frac{\partial u_r}{\partial r} & \frac{1}{r} \frac{\partial u_r}{\partial \theta} - \frac{u_\theta}{r} & \frac{\partial u_r}{\partial z} \\ \frac{\partial u_\theta}{\partial r} & \frac{1}{r} \frac{\partial u_\theta}{\partial \theta} + \frac{u_r}{r} & \frac{\partial u_\theta}{\partial z} \\ \frac{\partial u_z}{\partial r} & \frac{1}{r} \frac{\partial u_z}{\partial \theta} & \frac{\partial u_z}{\partial z} \end{pmatrix}$$

In the case of axisymmetry, and in this case symmetry about the z axis, there is invariance with respect to the rotation around the axis so stresses and other quantities are independent of the θ coordinate, or simply put $\partial_\theta \rightarrow 0$. The velocity gradient simplifies to:

$$\vec{\nabla} \vec{u} = \begin{pmatrix} \frac{\partial u_r}{\partial r} & -\frac{u_\theta}{r} & \frac{\partial u_r}{\partial z} \\ \frac{\partial u_\theta}{\partial r} & \frac{u_r}{r} & \frac{\partial u_\theta}{\partial z} \\ \frac{\partial u_z}{\partial r} & 0 & \frac{\partial u_z}{\partial z} \end{pmatrix}$$

Also, it follows logically that $u_\theta = 0$ so that ultimately:

$$\vec{\nabla} \vec{u} = \begin{pmatrix} \frac{\partial u_r}{\partial r} & 0 & \frac{\partial u_r}{\partial z} \\ 0 & \frac{u_r}{r} & 0 \\ \frac{\partial u_z}{\partial r} & 0 & \frac{\partial u_z}{\partial z} \end{pmatrix}$$

and the strain tensor is then given by

$$\boldsymbol{\varepsilon}(\vec{u}) = \frac{1}{2} \left(\vec{\nabla} \vec{u} + \vec{\nabla} \vec{u}^T \right) = \begin{pmatrix} \frac{\partial u_r}{\partial r} & 0 & \frac{1}{2} \left(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z} \right) \\ 0 & \frac{u_r}{r} & 0 \\ \frac{1}{2} \left(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z} \right) & 0 & \frac{\partial u_z}{\partial z} \end{pmatrix} \quad (7.181)$$

The term $\vec{\nabla} \cdot \vec{u}$ is simply the trace of $\boldsymbol{\varepsilon}(\vec{u})$ so

$$\vec{\nabla} \cdot \vec{u} = \frac{\partial u_r}{\partial r} + \frac{u_r}{r} + \frac{\partial u_z}{\partial z}$$

Finally the full stress tensor is then

$$\begin{aligned}\boldsymbol{\sigma} &= \begin{pmatrix} \lambda(\frac{\partial u_r}{\partial r} + \frac{u_r}{r} + \frac{\partial u_z}{\partial z}) + 2\mu\frac{\partial u_r}{\partial r} & 0 & \mu(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z}) \\ 0 & \lambda(\frac{\partial u_r}{\partial r} + \frac{u_r}{r} + \frac{\partial u_z}{\partial z}) + 2\mu\frac{u_r}{r} & 0 \\ \mu(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z}) & 0 & \lambda(\frac{\partial u_r}{\partial r} + \frac{u_r}{r} + \frac{\partial u_z}{\partial z}) + 2\mu\frac{\partial u_z}{\partial z} \end{pmatrix} \\ &= \begin{pmatrix} (\lambda + 2\mu)\frac{\partial u_r}{\partial r} + \lambda(\frac{u_r}{r} + \frac{\partial u_z}{\partial z}) & 0 & \mu(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z}) \\ 0 & (\lambda + 2\mu)\frac{u_r}{r} + \lambda(\frac{\partial u_r}{\partial r} + \frac{\partial u_z}{\partial z}) & 0 \\ \mu(\frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z}) & 0 & (\lambda + 2\mu)\frac{\partial u_z}{\partial z} + \lambda(\frac{\partial u_r}{\partial r} + \frac{u_r}{r}) \end{pmatrix}\end{aligned}$$

As we did in the 2D case, we rewrite the six independent stress terms in to a vector $\vec{\sigma}$ and we use Eq. (7.180) to arrive at:

$$\vec{\sigma} = \begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{r\theta} \\ \sigma_{rz} \\ \sigma_{\theta z} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{r\theta} \\ 2\varepsilon_{rz} \\ 2\varepsilon_{\theta z} \end{pmatrix} = \vec{\varepsilon}(\vec{u})$$

or $\vec{\sigma} = \mathbf{D} \cdot \vec{\varepsilon}(\vec{u})$. Notice the similarity of matrix \mathbf{D} with the one of Section (XXX) in the 3D penalty formulation case. The components of the $\vec{\varepsilon}$ vector are

$$\vec{\varepsilon}(\vec{u}) = \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{r\theta} \\ 2\varepsilon_{rz} \\ 2\varepsilon_{\theta z} \end{pmatrix} = \begin{pmatrix} \frac{\partial u_r}{\partial r} \\ \frac{u_r}{r} \\ \frac{\partial u_z}{\partial z} \\ 0 \\ \frac{\partial u_z}{\partial r} + \frac{\partial u_r}{\partial z} \\ 0 \end{pmatrix}$$

We see that there are two zeroes and consequently we'll find that $\sigma_{r\theta}$ and $\sigma_{\theta z}$ are also identically zero, so we discard these and end up with only four stress components :

$$\vec{\sigma} = \begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{rz} \end{pmatrix}$$

Note that in the literature the above relationship is often written

$$\begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix} = \frac{E}{(1+\nu)(1-2\nu)} \begin{pmatrix} 1-\nu & \lambda & \nu & 0 \\ \nu & 1-\nu & \nu & 0 \\ \nu & \nu & 1-\nu & 0 \\ 0 & 0 & 0 & (1-2\nu)/2 \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{rz} \end{pmatrix}$$

which is equivalent since $E = 2\mu(1+\nu)$ and $\lambda = \frac{\nu E}{(1+\nu)(1-2\nu)}$ (see for instance Section 5.2.4 in [1430]).

Only displacements in the r and z directions remain (note that $\varepsilon_{\theta\theta}$ is in fact equal to u_r/r). In what follows I rename $u = u_r$ and $u_z = w$ to simplify notations. Then, inside an element we have

$$\begin{aligned}u^h(r, z) &= \sum_{i=1}^m N_i(r, z) u_i \\ w^h(r, z) &= \sum_{i=1}^m N_i(r, z) w_i\end{aligned}\tag{7.182}$$

where N_i are the basis functions attached to the m nodes of the element. We compute the elements of the $\boldsymbol{\varepsilon}$ tensor of Eq. (7.181) as follows:

$$\varepsilon_{rr} = \frac{\partial u^h}{\partial r} = \sum_{i=1}^m \frac{\partial N_i}{\partial r}(r, z) u_i \quad (7.183)$$

$$\varepsilon_{\theta\theta} = \frac{u_r^h}{r} = \frac{1}{r} \sum_{i=1}^m N_i(r, z) u_i \quad (7.184)$$

$$\varepsilon_{zz} = \frac{\partial w^h}{\partial z} = \sum_{i=1}^m \frac{\partial N_i}{\partial z}(r, z) w_i \quad (7.185)$$

$$\varepsilon_{rz} = \frac{1}{2} \frac{\partial u^h}{\partial z} + \frac{1}{2} \frac{\partial w^h}{\partial r} = \sum_{i=1}^m \frac{\partial N_i}{\partial z}(r, z) u_i + \sum_{i=1}^m \frac{\partial N_i}{\partial r}(r, z) w_i \quad (7.186)$$

Let us take $m = 3$, i.e. linear triangles, for simplicity. Then the strain vector $\vec{\varepsilon}^h$ is given by

$$\vec{\varepsilon}^h = \begin{pmatrix} \frac{\partial u^h}{\partial r} \\ \frac{u_r^h}{r} \\ \frac{\partial w^h}{\partial z} \\ \frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial r} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial N_1}{\partial r} & 0 & \frac{\partial N_2}{\partial r} & 0 & \frac{\partial N_3}{\partial r} & 0 \\ \frac{N_1}{r} & 0 & \frac{N_2}{r} & 0 & \frac{N_3}{r} & 0 \\ 0 & \frac{\partial N_1}{\partial z} & 0 & \frac{\partial N_2}{\partial z} & 0 & \frac{\partial N_3}{\partial z} \\ \frac{\partial N_1}{\partial z} & \frac{\partial N_1}{\partial r} & \frac{\partial N_2}{\partial z} & \frac{\partial N_2}{\partial r} & \frac{\partial N_3}{\partial z} & \frac{\partial N_3}{\partial r} \end{pmatrix}}_{\mathbf{B}(4 \times 6)} \cdot \underbrace{\begin{pmatrix} u1 \\ w1 \\ u2 \\ w2 \\ u3 \\ w3 \end{pmatrix}}_{\vec{U}(6 \times 1)}$$

or $\vec{\varepsilon}^h = \mathbf{B} \cdot \vec{U}$ and finally

$$\underbrace{\begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix}}_{\vec{\sigma}} = \underbrace{\begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix}}_{\mathbf{D}} \cdot \underbrace{\begin{pmatrix} \frac{\partial N_1}{\partial r} & 0 & \frac{\partial N_2}{\partial r} & 0 & \frac{\partial N_3}{\partial r} & 0 \\ \frac{N_1}{r} & 0 & \frac{N_2}{r} & 0 & \frac{N_3}{r} & 0 \\ 0 & \frac{\partial N_1}{\partial z} & 0 & \frac{\partial N_2}{\partial z} & 0 & \frac{\partial N_3}{\partial z} \\ \frac{\partial N_1}{\partial z} & \frac{\partial N_1}{\partial r} & \frac{\partial N_2}{\partial z} & \frac{\partial N_2}{\partial r} & \frac{\partial N_3}{\partial z} & \frac{\partial N_3}{\partial r} \end{pmatrix}}_{\mathbf{B}(4 \times 6)} \cdot \underbrace{\begin{pmatrix} u1 \\ w1 \\ u2 \\ w2 \\ u3 \\ w3 \end{pmatrix}}_{\vec{U}(6 \times 1)}$$

or,

$$\boxed{\vec{\sigma} = \mathbf{D} \cdot \mathbf{B} \cdot \vec{U}}$$

Note that in 2D, the matrix \mathbf{D} is 3×3 and \mathbf{B} is 3×6 .

I do not know yet how to arrive at what follows

The 6×6 stiffness matrix is then

$$\mathbb{K} = \iiint \mathbf{B}^T \cdot \mathbf{D} \cdot \mathbf{B} dV$$

with $dV = r dr d\theta dz$ in cylindrical coordinates. The integral over the θ coordinate yields a factor 2π so

$$\mathbb{K} = 2\pi \iint \mathbf{B}^T \cdot \mathbf{D} \cdot \mathbf{B} r dr dz$$

The integration can now be performed as simply as was the case in the plane stress problem.

Note that in practice the matrix \mathbf{D} is computed as follows (see for example Stone 63):

$$\mathbf{D} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix} = \lambda \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} + \mu \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

The divergence of the stress tensor is given by

$$\vec{\nabla} \cdot \boldsymbol{\sigma} = \left[\frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{rr}) + \frac{1}{r} \frac{\partial \sigma_{r\theta}}{\partial \theta} + \frac{\partial \sigma_{rz}}{\partial z} - \frac{\sigma_{\theta\theta}}{r} \right] \vec{e}_r \quad (7.187)$$

$$+ \left[\frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{r\theta}) + \frac{1}{r} \frac{\partial \sigma_{\theta\theta}}{\partial \theta} + \frac{\partial \sigma_{\theta z}}{\partial z} + \frac{\sigma_{r\theta}}{r} \right] \vec{e}_\theta \quad (7.188)$$

$$+ \left[\frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{rz}) + \frac{1}{r} \frac{\partial \sigma_{\theta z}}{\partial \theta} + \frac{\partial \sigma_{zz}}{\partial z} \right] \vec{e}_z \quad (7.189)$$

Since $\sigma_{r\theta} = \sigma_{\theta r} = 0$ and $\sigma_{z\theta} = \sigma_{\theta z} = 0$ and since $\partial_\theta \rightarrow 0$ then

$$\vec{\nabla} \cdot \boldsymbol{\sigma} = \left[\frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{rr}) + \frac{\partial \sigma_{rz}}{\partial z} - \frac{\sigma_{\theta\theta}}{r} \right] \vec{e}_r \quad (7.190)$$

$$+ \left[\frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{rz}) + \frac{\partial \sigma_{zz}}{\partial z} \right] \vec{e}_z \quad (7.191)$$

Then

$$\vec{\nabla} \cdot \boldsymbol{\sigma}|_r = \frac{1}{r} \frac{\partial}{\partial r} (r \sigma_{rr}) + \frac{\partial \sigma_{rz}}{\partial z} - \frac{\sigma_{\theta\theta}}{r} \quad (7.192)$$

$$= \frac{\partial \sigma_{rr}}{\partial r} + \frac{1}{r} (\sigma_{rr} - \sigma_{\theta\theta}) + \frac{\partial \sigma_{rz}}{\partial z} \quad (7.193)$$

$$= \frac{\partial \sigma_{rr}}{\partial r} + \frac{2\mu}{r} \left(\frac{\partial u_r}{\partial r} - \frac{u_r}{r} \right) + \frac{\partial \sigma_{rz}}{\partial z} \quad (7.194)$$

$$\vec{\nabla} \cdot \boldsymbol{\sigma}|_z = \frac{\partial \sigma_{rz}}{\partial r} + \frac{\sigma_{rz}}{r} + \frac{\partial \sigma_{zz}}{\partial z} \quad (7.195)$$

7.8 The case against the $Q_1 \times P_0$ element

What follows was written by Dave May and sent to me by email in May 2014. It captures so well the problem at hand that I have decided to reproduce it hereunder.

In the case of the incompressible Stokes equations, we would like to solve

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ 0 \end{pmatrix}$$

with an iterative method which is algorithmically scalable and optimal. Scalable here would mean that the number of iterations doesn't grow as the mesh is refined. Optimal means the solution time varies linearly with the total number of unknowns. When using a stable element, If we right precondition the above system with

$$P = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ 0 & -\mathbb{S} \end{pmatrix}$$

then convergence will occur in 2 iterations, however this requires an exact solve on \mathbb{K} and on $\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$ (\mathbb{S} is the pressure schur complement). In practice, people relax the ideal "two iteration" scenario by first replacing \mathbb{S} via $\mathbb{S}^* = \int \eta^{-1} \vec{N}^T \vec{N} dv$ (e.g. the pressure mass matrix scaled by the local inverse of viscosity).

$$P^* = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ 0 & -\mathbb{S}^* \end{pmatrix}$$

Using P^* , we obtain iteration counts which are larger than 2, but likely less than 10 - *however*, the number of iterations is independent of the mesh size. Replacing the exact \mathbb{K} solve in P^* again increases the iterations required to solve Stokes, but it's still independent of the number of elements. When you have this behaviour, we say the preconditioner (P^*) is spectrally equivalent to the operator (which here is Stokes)

The problem with $Q_1 \times P_0$ is that there are no approximations for \mathbb{S} which can be generated that ensure a spectrally equivalent P^* . Thus, as you refine the mesh using $Q_1 \times P_0$ elements, the iteration count ALWAYS grows. I worked on this problem during my thesis, making some improvements to the situation - however the problem still remains, it cannot be completely fixed and stems entirely from using unstable elements.

Citcom solvers works like this:

1. Solve $\mathbb{S} \cdot \mathcal{P} = \vec{f}'$ for pressure
2. Solve $\mathbb{K} \cdot \mathcal{V} = \vec{f} - \mathbb{G} \cdot \mathcal{P}$ for velocity

To obtain a scalable method, we need the number of iterations performed in (1) and (2) to be independent of the mesh. This means we need a spectrally equivalent preconditioner for \mathbb{S} and \mathbb{K} . Thus, we have the same issue as when you iterate on the full stokes system.

When we don't have a scalable method, it means increasing the resolution requires more cpu time in a manner which cannot be predicted. The increase in iteration counts as the mesh is refined can be dramatic.

If we can bound the number of iterations, AND ensure that the cost per iteration is linearly related to the number of unknowns, then we have a good method which can run on any mesh resolution with a predictable cpu time. Obtaining scalable and optimal preconditioners for \mathbb{K} is somewhat easier. Multi-grid will provide us with this.

The reason citcom doesn't run with 400^3 elements is exactly due to this issue. I've added petsc support in citcom (when i was young and naive) - but the root cause of the non-scalable solve is directly caused by the element choice. Note that many of the high resolution citcom jobs are single time step calculations— there is a reason for that.

For many lithosphere dynamics problems, we need a reasonable resolution (at least 200^3 and realistically 400^3 to 800^3). Given the increase in cost which occurs when using Q1P0, this is not achievable, as the citcom code has demonstrated. Note that citcom is 20 years old now and for its time, it was great, but we know much more now and we know how to improve on it. As a result of this realization, I dumped all my old Q1P0 codes (and Q1Q1 codes, but for other reasons) in the trash and started from scratch. The only way to make something like 800^3 tractable is via iterative, scalable and optimal methods and that mandates stable elements. I can actually run at something like 1000^3 (nodal points) these days because of such design choices.

7.9 Isoviscous Stokes for incompressible flow

isoviscous_stokes.tex

We start from the momentum equation:

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}^d(\vec{v})) + \rho \vec{g} = \vec{0} \quad (7.196)$$

When the viscosity is constant in space, it can be taken out of the divergence operator:

$$-\vec{\nabla} p + 2\eta \vec{\nabla} \cdot \dot{\epsilon}^d(\vec{v}) + \rho \vec{g} = \vec{0} \quad (7.197)$$

Let us for simplicity look at a 2D Cartesian formulation of this equation and for incompressible flow:

$$\begin{aligned} 2\vec{\nabla} \cdot \dot{\epsilon}^d(\vec{v}) &= \vec{\nabla} \cdot (\vec{\nabla} \vec{v} + \vec{\nabla} \vec{v}^T) \\ &= (\partial_x \partial_y) \cdot \begin{pmatrix} \partial_x u & \partial_x v \\ \partial_y u & \partial_y v \end{pmatrix} + (\partial_x \partial_y) \cdot \begin{pmatrix} \partial_x u & \partial_y u \\ \partial_x v & \partial_y v \end{pmatrix} \\ &= (\partial_x^2 u + \partial_y^2 u, \partial_x^2 v + \partial_y^2 v) + (\partial_x \partial_x u + \partial_y \partial_x v, \partial_x \partial_y u + \partial_y \partial_y v) \\ &= (\partial_x^2 u + \partial_y^2 u, \partial_x^2 v + \partial_y^2 v) + (\underbrace{\partial_x (\partial_x u + \partial_y v)}_{=0}, \underbrace{\partial_y (\partial_x u + \partial_y v)}_{=0}) \\ &= (\partial_x^2 u + \partial_y^2 u, \partial_x^2 v + \partial_y^2 v) \end{aligned} \quad (7.198)$$

and then finally the Stokes equation is:

$$-\vec{\nabla} p + \eta \Delta \vec{v} + \rho \vec{g} = \vec{0} \quad (7.199)$$

The mass conservation equation remains unchanged and so does the pressure gradient term. We shall then focus on the weak form of the previously obtained term. We multiply it by a velocity test function \mathcal{N}_i^γ and integrate over an element²¹:

$$\begin{aligned} &\int_{\Omega_e} \mathcal{N}_i^\gamma \Delta \vec{v}^h dV \\ &= \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_i^\gamma \Delta u^h \\ \mathcal{N}_i^\gamma \Delta v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_i^\gamma \vec{\nabla} \cdot \vec{\nabla} u^h \\ \mathcal{N}_i^\gamma \vec{\nabla} \cdot \vec{\nabla} v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \vec{\nabla} \mathcal{N}_i^\gamma \cdot \vec{\nabla} u^h \\ \vec{\nabla} \mathcal{N}_i^\gamma \cdot \vec{\nabla} v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \partial_x \mathcal{N}_i^\gamma \partial_x u^h + \partial_y \mathcal{N}_i^\gamma \partial_y u^h \\ \partial_x \mathcal{N}_i^\gamma \partial_x v^h + \partial_y \mathcal{N}_i^\gamma \partial_y v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \partial_x \mathcal{N}_i^\gamma & \partial_y \mathcal{N}_i^\gamma & 0 & 0 \\ 0 & 0 & \partial_x \mathcal{N}_i^\gamma & \partial_y \mathcal{N}_i^\gamma \end{pmatrix} \cdot \begin{pmatrix} \partial_x u^h \\ \partial_y u^h \\ \partial_x v^h \\ \partial_y v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \partial_x \mathcal{N}_i^\gamma & \partial_y \mathcal{N}_i^\gamma & 0 & 0 \\ 0 & 0 & \partial_x \mathcal{N}_i^\gamma & \partial_y \mathcal{N}_i^\gamma \end{pmatrix} \cdot \begin{pmatrix} \partial_x \mathcal{N}_1^\gamma & 0 & \partial_x \mathcal{N}_2^\gamma & 0 & \dots & \partial_x \mathcal{N}_{m_v}^\gamma & 0 \\ \partial_y \mathcal{N}_1^\gamma & 0 & \partial_y \mathcal{N}_2^\gamma & 0 & \dots & \partial_y \mathcal{N}_{m_v}^\gamma & 0 \\ 0 & \partial_x \mathcal{N}_1^\gamma & 0 & \partial_x \mathcal{N}_2^\gamma & \dots & 0 & \partial_x \mathcal{N}_{m_v}^\gamma \\ 0 & \partial_y \mathcal{N}_1^\gamma & 0 & \partial_y \mathcal{N}_2^\gamma & \dots & 0 & \partial_y \mathcal{N}_{m_v}^\gamma \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix} dV \end{aligned}$$

²¹As per usual we discard the surface term when integrating by parts

Writing this equation for $i = 1, \dots, m_v$, we obtain:

$$\int \begin{pmatrix} \partial_x \mathcal{N}_1^\nu & \partial_y \mathcal{N}_1^\nu & 0 & 0 \\ 0 & 0 & \partial_x \mathcal{N}_1^\nu & \partial_y \mathcal{N}_1^\nu \\ \partial_x \mathcal{N}_2^\nu & \partial_y \mathcal{N}_2^\nu & 0 & 0 \\ 0 & 0 & \partial_x \mathcal{N}_2^\nu & \partial_y \mathcal{N}_2^\nu \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \partial_x \mathcal{N}_{m_v}^\nu & \partial_y \mathcal{N}_{m_v}^\nu & 0 & 0 \\ 0 & 0 & \partial_x \mathcal{N}_{m_v}^\nu & \partial_y \mathcal{N}_{m_v}^\nu \end{pmatrix} \cdot \begin{pmatrix} \partial_x \mathcal{N}_1^\nu & 0 & \partial_x \mathcal{N}_2^\nu & 0 & \dots & \partial_x \mathcal{N}_{m_v}^\nu & 0 \\ \partial_y \mathcal{N}_1^\nu & 0 & \partial_y \mathcal{N}_2^\nu & 0 & \dots & \partial_y \mathcal{N}_{m_v}^\nu & 0 \\ 0 & \partial_x \mathcal{N}_1^\nu & 0 & \partial_x \mathcal{N}_2^\nu & \dots & 0 & \partial_x \mathcal{N}_{m_v}^\nu \\ 0 & \partial_y \mathcal{N}_1^\nu & 0 & \partial_y \mathcal{N}_2^\nu & \dots & 0 & \partial_y \mathcal{N}_{m_v}^\nu \end{pmatrix} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix}}_{\vec{v}}$$

or,

$$\mathbf{K}_\eta = \eta \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{B} dV$$

where \mathbf{B} is a $(ndim * ndim) \times (m_v * ndofV)$ matrix (see also Eq. 6.24 of Donea and Huerta [341]).

In three dimensions, the matrix \mathbf{B} is given by

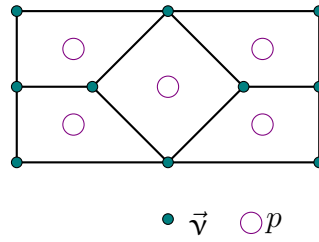
$$\begin{pmatrix} \partial_x \mathcal{N}_1^\nu & 0 & \partial_x \mathcal{N}_2^\nu & 0 & \dots & \partial_x \mathcal{N}_{m_v}^\nu & 0 \\ \partial_y \mathcal{N}_1^\nu & 0 & \partial_y \mathcal{N}_2^\nu & 0 & \dots & \partial_y \mathcal{N}_{m_v}^\nu & 0 \\ \partial_z \mathcal{N}_1^\nu & 0 & \partial_z \mathcal{N}_2^\nu & 0 & \dots & \partial_z \mathcal{N}_{m_v}^\nu & 0 \\ 0 & \partial_x \mathcal{N}_1^\nu & 0 & \partial_x \mathcal{N}_2^\nu & \dots & 0 & \partial_x \mathcal{N}_{m_v}^\nu \\ 0 & \partial_y \mathcal{N}_1^\nu & 0 & \partial_y \mathcal{N}_2^\nu & \dots & 0 & \partial_y \mathcal{N}_{m_v}^\nu \\ 0 & \partial_z \mathcal{N}_1^\nu & 0 & \partial_z \mathcal{N}_2^\nu & \dots & 0 & \partial_z \mathcal{N}_{m_v}^\nu \end{pmatrix}$$

7.10 $Q_1 \times P_0$ macro-elements

The Stenberg macro-element

This macro-element is introduced in Stenberg (1984) [1206].

(tikz_stenberg.tex)



Gresho & Sani [488] state: "For fans of $Q_1 \times Q_0$ who want guaranteed optimal convergence of both u and p (with however larger error constants caused by the distorted shapes?), one way to assure this is to discretise via the macro elements above, each composed of five $Q_1 \times Q_0$ quadrilaterals. Such checkerboard-killer meshes have been employed in practice by (at least) Bathé [220]. Both the macro-element and the proof are due to Stenberg [1206]."

Chapelle & Bathe [220]: "the numerical inf-sup test is passed for this mesh and in fact, this behavior was proven analytically (see Brezzi & Fortin [148], see also Le Tallec & Ruas [774])."

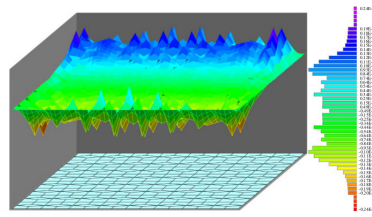



Figure 14. Error of component 1 of u_h on the Stenberg mesh [15].

Taken from Qin & Zhang (2007) [1025].

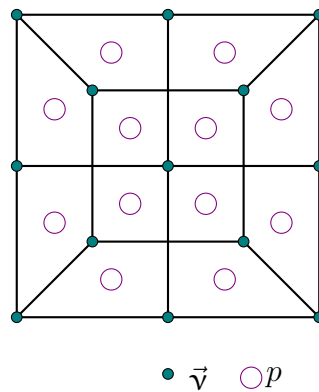
Implemented in [STONE](#) 78.

 Relevant Literature: Fig 3.12 of Elman *et al.* book [371]. Mentioned in Qin and Zhang [1026] (2007).

The Le Tallec macro-element

This macro-element is introduced in Le Tallec (1981) [773].

(tikz_letallec.tex)

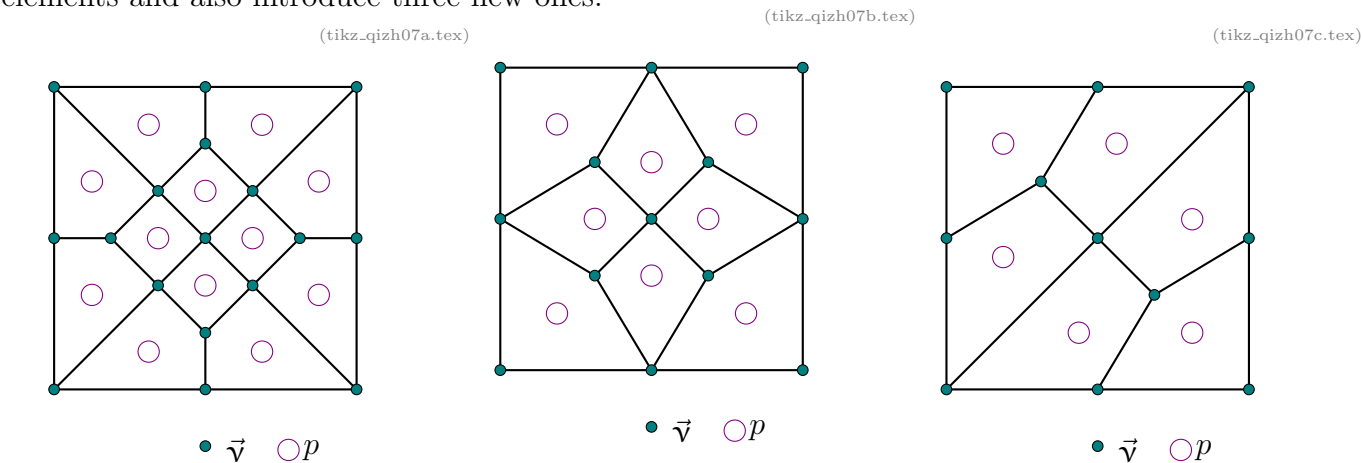


This macro-element has been proven stable in [773, 774], i.e. it satisfies the stability condition (see Section 7.3). It is also mentioned in Qin & Zhang (2007) [1025].

Implemented in [STONE](#) 78.

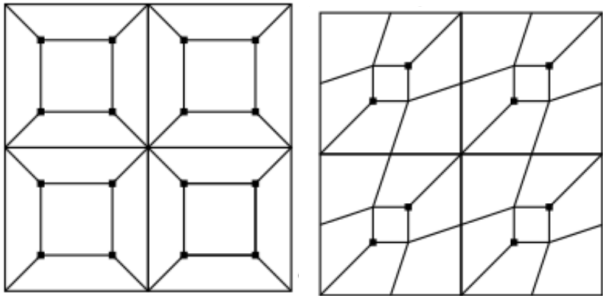
The Qin & Zhang macro-elements

In their paper Qin & Zhang (2007) [1025] the authors mention the Stenberg and Le Tallec macro-elements and also introduce three new ones:



They also indicate that although stable, these macro-elements are inferior to the above two (Stenberg & Le Tallec).

New macro-elements ?



I came up with these, no idea whether these are stable/usable or better than the others.

7.11 Solving the Stokes system

Let us start again from the (full) Stokes system:

$$\underbrace{\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix}}_{\mathcal{A}} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (7.200)$$

We need to solve this system in order to obtain the solution, i.e. the $\vec{\mathcal{V}}$ and $\vec{\mathcal{P}}$ vectors. But how? Unfortunately, this question is not simple to answer and the appropriate method depends on many parameters, but mainly on how big the matrix blocks are and what the condition number of the matrix \mathbb{K} is.

First let us start with an obvious question: couldn't we just compute the inverse of the matrix \mathcal{A} ? Under the assumption that the inverse of \mathbb{K} and \mathbb{S} exists, we can and we find²²

$$\mathcal{A}^{-1} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix}^{-1} = \begin{pmatrix} \mathbb{K}^{-1} + \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \mathbb{S}^{-1} \cdot \mathbb{G}^T \cdot \mathbb{K}^{-1} & -\mathbb{K}^{-1} \cdot \mathbb{G} \cdot \mathbb{S}^{-1} \\ -\mathbb{S}^{-1} \cdot \mathbb{G}^T \cdot \mathbb{K}^{-1} & \mathbb{S}^{-1} \end{pmatrix}$$

However, such an expression is of limited interest in the numerical solution of saddle point problems since it showcases 5 times the inverse of \mathbb{K} and more importantly the inverse of the Schur complement matrix \mathbb{S} which is likely to be a full matrix so that we never want to compute it explicitly.

As concisely explained in Clevenger & Heister (2021) [261], there are three common approaches used in the literature for solving the above equation on large scales:

- a pressure corrected, Schur complement CG scheme, using multigrid as an approximation to the velocity block;
- a block-preconditioned Krylov method, also using multigrid on the velocity block. For this method, there are two main types:
 - GMRES[845, 1088] (or any Krylov method not requiring symmetry) with block-triangular preconditioner (This is what ASPECT does):

$$\mathbf{P} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ 0 & -\mathbb{S} \end{pmatrix}$$

- MINRES[469] with block-diagonal preconditioner

$$\mathbf{P} = \begin{pmatrix} \mathbb{K} & 0 \\ 0 & -\mathbb{S} \end{pmatrix}$$

- an all-at-once multigrid performed on the entire Stokes system, using Uzawa-type smoothers.



Relevant Literature: Preconditioners for Incompressible Navier-Stokes Solvers [1148]

Saddle point preconditioners have been extensively discussed and studied [73], [941]

Diagonal preconditioners in [1145], [48].

Pragmatic solvers for 3D Stokes problems with heterogeneous coefficients [1106]

²²The matrix \mathbb{C} is here omitted but it bears no consequences on the conclusion.

7.11.1 When using the penalty formulation

In this case we are only solving for velocity since pressure has been eliminated and is later recovered in a post-processing step:

$$(\mathbb{K}_\eta + \mathbb{K}_\lambda) \cdot \vec{\mathcal{V}} = \vec{f}$$

We also know that the penalty factor λ is many orders of magnitude higher than the viscosity and in combination with the use of the $Q_1 \times P_0$ element the resulting matrix condition number is very high so that the use of iterative solvers is precluded. Indeed codes such as SOPALE [426], DOUAR [136], FANTOM [1258] or SULEC [1028] relying on the penalty formulation all use direct solvers. The most popular are BLKFCT²³, MUMPS²⁴[13, 15, 14, 16, 12], PaSTiX [563], WSMP²⁵ [510, 511], UMFPACK and CHOLMOD²⁶, SuperLU²⁷, PARDISO²⁸ [321, 1318, 725], or those inside PETSc²⁹.

Braun *et al.* (2008) [136] list the following features of direct solvers:

- Robust
- Black-box operation
- Difficult to parallelize
- Memory consumption
- Limited scalability

The main advantage of direct solvers is used in this case: They can solve ill-conditioned matrices. However, memory requirements for the storage of number of nonzeros in the Cholesky matrix grow very fast as the number of equations/grid size increases, especially in 3D, to the point that even modern computers with tens of Gb of RAM cannot deal with a $\sim 100^3$ element mesh. This explains why direct solvers are often used for 2D problems and rarely in 3D with noticeable exceptions [1261, 1377, 137, 807, 10, 9, 11, 1352, 933].

7.11.2 Uzawa algorithms and the Schur complement approach

Let us write the above system as two equations:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (7.201)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} - \mathbb{C} \cdot \vec{\mathcal{P}} = \vec{h} \quad (7.202)$$

The first line can be re-written $\vec{\mathcal{V}} = \mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}})$ and can be inserted in the second:

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} = \mathbb{G}^T \cdot [\mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}})] - \mathbb{C} \cdot \vec{\mathcal{P}} = \vec{h} \quad (7.203)$$

or,

$$(\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot \vec{\mathcal{P}} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h} \quad (7.204)$$

²³<http://dm.unife.it/blkfclt/>

²⁴<http://mumps.enseeiht.fr/>

²⁵<http://www.research.ibm.com/projects/wsmp>

²⁶<http://faculty.cse.tamu.edu/davis/suitesparse.html>

²⁷<https://portal.nersc.gov/project/sparse/superlu/>

²⁸<https://www.pardiso-project.org/>

²⁹<https://www.mcs.anl.gov/petsc/>

The matrix $\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}$ is called the Schur complement. It is Symmetric (since \mathbb{K} is symmetric) and Positive-Definite³⁰ (SPD) if $\text{Ker}(\mathbb{G}) = 0$. Having solved this equation (i.e. we have obtained $\vec{\mathcal{P}}$), the velocity can be recovered by solving $\mathbb{K} \cdot \vec{\mathcal{V}} = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}$.

Remark. The Schur complement matrix naturally occurs when the Stokes matrix is decomposed using a LDU block-factorisation. Indeed, we have

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} = \begin{pmatrix} \mathbf{I} & 0 \\ \mathbb{G}^T \cdot \mathbb{K}^{-1} & \mathbf{I} \end{pmatrix} \cdot \begin{pmatrix} \mathbb{K} & 0 \\ 0 & -\mathbb{S} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I} & \mathbb{K}^{-1} \cdot \mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix}$$

For now, let us assume that we have built the \mathbb{S} matrix³¹ and the right hand side $\vec{f} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h}$. We must solve $\mathbb{S} \cdot \vec{\mathcal{P}} = \vec{f}$. It is easy to see that \mathbb{S} is actually a full matrix (i.e. not sparse) and aside from the costs of building it explicitly using a direct solver would require a lot of memory so that we must then turn to iterative methods.

One can resort to so-called Richardson iterations, defined as follows (e.g., see Varga [1312], p141): in solving the matrix equation $\mathbf{A} \cdot \vec{X} = \vec{b}$, the Richardson iterative method is defined by:

$$\vec{X}_{k+1} = \vec{X}_k + \alpha_k(-\mathbf{A} \cdot \vec{X}_k + \vec{b}) \quad m \geq 0 \quad (7.205)$$

where the α_k 's are real scalars. It is easy to see that when the method converges then $\vec{X}_{k+1} \simeq \vec{X}_k$ and then for $\alpha_k \neq 0$ then $\mathbf{A} \cdot \vec{X} = \vec{b}$ is satisfied. In our case, it writes:

$$\begin{aligned} \vec{\mathcal{P}}_{k+1} &= \vec{\mathcal{P}}_k + \alpha_k(-\mathbb{S} \cdot \vec{\mathcal{P}}_k + \vec{f}) \\ &= \vec{\mathcal{P}}_k + \alpha_k \left[-(\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot \vec{\mathcal{P}}_k + (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h}) \right] \\ &= \vec{\mathcal{P}}_k + \alpha_k \left[\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (-\mathbb{G} \cdot \vec{\mathcal{P}}_k + \vec{f}) - \mathbb{C} \cdot \vec{\mathcal{P}}_k - \vec{h} \right] \\ &= \vec{\mathcal{P}}_k + \alpha_k \left[\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (\mathbb{K} \cdot \vec{\mathcal{V}}_k) - \mathbb{C} \cdot \vec{\mathcal{P}}_k - \vec{h} \right] \\ &= \vec{\mathcal{P}}_k + \alpha_k \left(\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \mathbb{C} \cdot \vec{\mathcal{P}}_k - \vec{h} \right) \end{aligned} \quad (7.206)$$

The above iterations are then carried out and for each new pressure field the associated velocity field is computed. The method of using Richardson iterations applied to the Schur complement is commonly called the Uzawa algorithm (see Braess [128, p221]³²).

Uzawa algorithm (1): assume $\vec{\mathcal{P}}_0$ known

$$\text{solve} \quad \mathbb{K} \cdot \vec{\mathcal{V}}_k = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_k \quad (7.207)$$

$$\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k(\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \mathbb{C} \cdot \vec{\mathcal{P}}_k - \vec{h}) \quad k = 0, 1, 2, \dots \quad (7.208)$$

This method is rather simple to implement, although what makes an appropriate set of α_k values is not straightforward, which is why the conjugate gradient is often preferred, as detailed in the next section.

It is known that such iterations will converge for $0 < \alpha < \rho(\mathbb{S}) = \lambda_{max}(\mathbb{S})$ where $\rho(\mathbb{S})$ is the spectral radius of the matrix \mathbb{S} which is essentially the largest, in absolute value, eigenvalue of \mathbb{S} (neither of which can be computed easily). It can also be proven that the rate of convergence depends on the condition number of the matrix.

³⁰ M positive definite $\iff x^T M x > 0 \forall x \in \mathbb{R}^n \setminus \mathbf{0}$

³¹We will revisit this topic later on, but be aware that we never build \mathbb{S} in practice.

³²I have slightly altered the indices of the velocities wrt the book

Richardson iterations are part of the family of stationary iterative methods³³, since it can be rewritten

$$\vec{X}_{k+1} = (\mathbf{I} - \alpha_k \mathbf{A}) \cdot \vec{X}_k + \alpha_k \vec{b} \quad (7.209)$$

which is the definition of a stationary method. The four main stationary methods are the Jacobi method, Gauss-Seidel method, successive overrelaxation method (SOR), and symmetric successive overrelaxation method (SSOR)

Since the α parameter is the key to a successful Uzawa algorithm, this issue has of course been looked into. What follows is presented in p221 of Braess [128]. For the analysis of the Uzawa algorithm, we define the residue

$$\vec{\mathcal{R}}_k = \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}_k + \mathbb{C} \cdot \vec{\mathcal{P}}_k$$

In addition, suppose the solution of the saddle point problem is denoted by $(\vec{\mathcal{V}}^*, \vec{\mathcal{P}}^*)$ so that we have

$$\vec{f} = \mathbb{K} \cdot \vec{\mathcal{V}}^* + \mathbb{G} \cdot \vec{\mathcal{P}}^* \quad \text{and} \quad \vec{h} = \mathbb{G}^T \cdot \vec{\mathcal{V}}^* - \mathbb{C} \cdot \vec{\mathcal{P}}^*$$

Now substituting the iteration formula for \mathcal{V}_k , and inserting \vec{f} and \vec{h} from above, we get

$$\begin{aligned} \vec{\mathcal{R}}_k &= \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}_k + \mathbb{C} \cdot \vec{\mathcal{P}}_k \\ &= \vec{h} - \mathbb{G}^T \cdot \mathbb{K}^{-1}(\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_k) + \mathbb{C} \cdot \vec{\mathcal{P}}_k \end{aligned} \quad (7.210)$$

$$= (\mathbb{G}^T \cdot \vec{\mathcal{V}}^* - \mathbb{C} \cdot \vec{\mathcal{P}}^*) - \mathbb{G}^T \cdot \mathbb{K}^{-1}(\mathbb{K} \cdot \vec{\mathcal{V}}^* + \mathbb{G} \cdot \vec{\mathcal{P}}^* - \mathbb{G} \cdot \vec{\mathcal{P}}_k) + \mathbb{C} \cdot \vec{\mathcal{P}}_k \quad (7.211)$$

$$= (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot (\vec{\mathcal{P}}_k - \vec{\mathcal{P}}^*) \quad (7.212)$$

From Eq. (7.208) it follows that:

$$\vec{\mathcal{P}}_{k+1} - \vec{\mathcal{P}}_k = \alpha (\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \mathbb{C} \cdot \vec{\mathcal{P}}_k - \vec{h}) \quad (7.213)$$

$$= -\alpha \vec{\mathcal{R}}_k \quad (7.214)$$

$$= -\alpha (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot (\vec{\mathcal{P}}_k - \vec{\mathcal{P}}^*) \quad (7.215)$$

$$= \alpha (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot (\vec{\mathcal{P}}^* - \vec{\mathcal{P}}_k) \quad (7.216)$$

Thus the Uzawa algorithm is equivalent to applying the gradient method to the reduced equation using a fixed step size. In particular, the iteration converges for $\alpha < 2\|\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}\|^{-1}$ and one can show that the good step size α_k is given by

$$\alpha_k = \frac{\vec{\mathcal{R}}_k \cdot \vec{\mathcal{R}}_k}{(\mathbb{G} \cdot \vec{\mathcal{R}}_k) \cdot (\mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{\mathcal{R}}_k)} \quad (7.217)$$

include matrix C!

However, if we were to use this rule formally, we would need an additional multiplication by \mathbb{K}^{-1} in every step of the iteration. This can be avoided by storing an auxiliary vector. Note that this algorithm is presented in Zienkiewicz *et al.* (1985) [1434] in the context of viscoplastic flow.

As mentioned above, there is a way to rework the original Uzawa algorithm to include Eq. (7.217). It yields a modified Uzawa algorithm (see p222 of Braess [128]³⁴):

³³<https://mathworld.wolfram.com/StationaryIterativeMethod.html>

³⁴I have slightly altered the indices of the velocities wrt the book

Uzawa algorithm (2): assume $\vec{\mathcal{P}}_0$ known. Solve $\mathbb{K} \cdot \vec{\mathcal{V}}_0 = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0$. For $k = 0, 1, 2, \dots$, compute

$$\vec{\mathcal{R}}_k = \vec{q}_k = \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}_k + \mathbb{C} \cdot \vec{\mathcal{P}}_k \quad (7.218)$$


$$\vec{p}_k = \mathbb{G} \cdot \vec{q}_k \quad (7.219)$$

$$\vec{H}_k = \mathbb{K}^{-1} \cdot \vec{p}_k \quad (7.220)$$

$$\alpha_k = \frac{\vec{q}_k \cdot \vec{q}_k}{\vec{p}_k \cdot \vec{H}_k} \quad (7.221)$$

$$\vec{\mathcal{P}}_k = \vec{\mathcal{P}}_{k-1} - \alpha_k \vec{q}_k \quad (7.222)$$

$$\vec{\mathcal{V}}_k = \vec{\mathcal{V}}_{k-1} + \alpha_k \vec{H}_k \quad (7.223)$$

 Relevant Literature: Cahouet & Chabard (1988) [201], Cao (2003) [207].

7.11.3 Conjugate gradient and the Schur complement approach

Since the Schur matrix \mathbb{S} is Symmetric Positive Definite, the Conjugate Gradient (CG) method³⁵ [567] is very appropriate to solve this system.

Indeed, looking at the definition of Wikipedia: *"In mathematics, the conjugate gradient method is an algorithm for the numerical solution of particular systems of linear equations, namely those whose matrix is symmetric and positive-definite. The conjugate gradient method is often implemented as an iterative algorithm, applicable to sparse systems that are too large to be handled by a direct implementation or other direct methods such as the Cholesky decomposition. Large sparse systems often arise when numerically solving partial differential equations or optimization problems."*

A simple Google search tells us that the Conjugate Gradient algorithm is as follows:

Algorithm as obtained from Wikipedia.

```

r0 := b - Ax0
if r0 is sufficiently small, then return x0 as the result
p0 := r0
k := 0
repeat
     $\alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$ 
    xk+1 := xk +  $\alpha_k \mathbf{p}_k$ 
    rk+1 := rk -  $\alpha_k \mathbf{A} \mathbf{p}_k$ 
    if rk+1 is sufficiently small, then exit loop
     $\beta_k := \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}$ 
    pk+1 := rk+1 +  $\beta_k \mathbf{p}_k$ 
    k := k + 1
end repeat
return xk+1 as the result

```

The same algorithm with our notations:

```

 $\vec{r}_0 = \vec{f} - \mathbb{S} \cdot \vec{\mathcal{P}}_0$ 
 $\vec{p}_0 = \vec{r}_0$ 
k = 0
repeat
     $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\vec{p}_k^T \cdot \mathbb{S} \cdot \vec{p}_k)$ 
     $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$ 
     $\vec{r}_{k+1} = \vec{r}_k - \alpha_k \mathbb{S} \cdot \vec{p}_k$ 
     $\beta_k = (\vec{r}_{k+1}^T \cdot \vec{r}_{k+1}) / (\vec{r}_k^T \cdot \vec{r}_k)$ 
     $\vec{p}_{k+1} = \vec{r}_{k+1} + \beta_k \vec{p}_k$ 
    k = k + 1
end repeat
return  $\vec{\mathcal{P}}_{k+1}$  as the result

```

This algorithm is of course explained in detail in many textbooks such as Saad [1092], in Zhong, Yuen, Moresi & Knepley (2012) [1415], and in Section 9.33.

Let us look at this algorithm more closely. The parts which may prove to be somewhat tricky are those involving the matrix the Schur complement matrix since we wish never to build it explicitly. We start the iterations with a guess pressure $\vec{\mathcal{P}}_0$ (and an initial guess velocity which could be obtained

³⁵https://en.wikipedia.org/wiki/Conjugate_gradient_method

by solving $\mathbb{K} \cdot \vec{\mathcal{V}}_0 = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0$.

$$\vec{r}_0 = \vec{f} - \mathbb{S} \cdot \vec{\mathcal{P}}_0 \quad (7.224)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h} - (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot \vec{\mathcal{P}}_0 \quad (7.225)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0) - \vec{h} \quad (7.226)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{K} \cdot \vec{\mathcal{V}}_0 - \mathbb{C} \cdot \vec{\mathcal{P}}_0 - \vec{h} \quad (7.227)$$

$$= \mathbb{G}^T \cdot \vec{\mathcal{V}}_0 - \mathbb{C} \cdot \vec{\mathcal{P}}_0 - \vec{h} \quad (7.228)$$

We see that we were able to compute $\mathbb{S} \cdot \vec{\mathcal{P}}_0$ without ever forming the Schur complement matrix explicitly. We now turn to the α_k coefficient:

$$\alpha_k = \frac{\vec{r}_k^T \cdot \vec{r}_k}{\vec{p}_k \cdot \mathbb{S} \cdot \vec{p}_k} = \frac{\vec{r}_k^T \cdot \vec{r}_k}{\vec{p}_k \cdot (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \cdot \vec{p}_k} = \frac{\vec{r}_k^T \cdot \vec{r}_k}{(\mathbb{G} \cdot \vec{p}_k)^T \cdot \mathbb{K}^{-1} \cdot (\mathbb{G} \cdot \vec{p}_k) + \vec{p}_k \cdot \mathbb{C} \cdot \vec{p}_k}$$

We then define $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$, so that α_k can be computed as follows:

1. compute $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$
2. solve $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
3. compute

$$\alpha_k = \frac{\vec{r}_k^T \cdot \vec{r}_k}{\tilde{\vec{p}}_k^T \cdot \vec{d}_k + \vec{p}_k^T \cdot \mathbb{C} \cdot \vec{p}_k}$$

Then we need to look at the term $\mathbb{S} \cdot \vec{p}_k$:

$$\mathbb{S} \cdot \vec{p}_k = (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}) \vec{p}_k = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \tilde{\vec{p}}_k + \mathbb{C} \cdot \vec{p}_k = \mathbb{G}^T \cdot \vec{d}_k + \mathbb{C} \cdot \vec{p}_k$$

We can then rewrite the CG algorithm as follows:

- choose $\vec{\mathcal{P}}_0$
- compute $\vec{\mathcal{V}}_0$ solution of $\mathbb{K} \cdot \vec{\mathcal{V}}_0 = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0$
- $\vec{r}_0 = \mathbb{G}^T \cdot \vec{\mathcal{V}}_0 - \mathbb{C} \cdot \vec{\mathcal{P}}_0 - \vec{h}$
- if \vec{r}_0 is sufficiently small, then return $(\vec{\mathcal{V}}_0, \vec{\mathcal{P}}_0)$ as the result
- $\vec{p}_0 = \vec{r}_0$
- $k = 0$
- repeat
 - compute $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$
 - solve $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
 - compute $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k + \vec{p}_k^T \cdot \mathbb{C} \cdot \vec{p}_k)$
 - $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$
 - $\vec{r}_{k+1} = \vec{r}_k - \alpha_k (\mathbb{G}^T \cdot \vec{d}_k + \mathbb{C} \cdot \vec{p}_k)$
 - if \vec{r}_{k+1} is sufficiently small, then exit loop
 - $\beta_k = (\vec{r}_{k+1}^T \cdot \vec{r}_{k+1}) / (\vec{r}_k^T \cdot \vec{r}_k)$
 - $\vec{p}_{k+1} = \vec{r}_{k+1} + \beta_k \vec{p}_k$

– $k = k + 1$

- return $\vec{\mathcal{P}}_{k+1}$ as result

We see that we have managed to solve the Schur complement equation with the Conjugate Gradient method without ever building the matrix \mathbb{S} . Having obtained the pressure solution $\vec{\mathcal{P}}_{k+1}$, we can easily recover the corresponding velocity with $\mathbb{K} \cdot \vec{\mathcal{V}}_{k+1} = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_{k+1}$. However, this is rather unfortunate because it requires yet another solve with the \mathbb{K} matrix. As it turns out, we can slightly alter the above algorithm to have it update the velocity as well so that this last solve is unnecessary.

We have

$$\vec{\mathcal{V}}_{k+1} = \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot \vec{\mathcal{P}}_{k+1}) \quad (7.229)$$

$$= \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot (\vec{\mathcal{P}}_k + \alpha_k \vec{p}_k)) \quad (7.230)$$

$$= \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot \vec{\mathcal{P}}_k) - \alpha_k \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{p}_k \quad (7.231)$$

$$= \vec{\mathcal{V}}_k - \alpha_k \mathbb{K}^{-1} \cdot \tilde{\vec{p}}_k \quad (7.232)$$

$$= \vec{\mathcal{V}}_k - \alpha_k \vec{d}_k \quad (7.233)$$

and we can insert this minor extra calculation inside the algorithm and get the velocity solution nearly for free. The final CG algorithm is then

solver_cg: assume $\vec{\mathcal{P}}_0$ known

- compute $\vec{\mathcal{V}}_0 = \mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0)$
- $\vec{r}_0 = \mathbb{G}^T \cdot \vec{\mathcal{V}}_0 - \mathbb{C} \cdot \vec{\mathcal{P}}_0 - \vec{h}$
- if \vec{r}_0 is sufficiently small, then return $(\vec{\mathcal{V}}_0, \vec{\mathcal{P}}_0)$ as the result
- $\vec{p}_0 = \vec{r}_0$
- $k = 0$
- repeat
 - compute $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$
 - solve $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
 - compute $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k + \vec{p}_k^T \cdot \mathbb{C} \cdot \vec{p}_k)$
 - $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$
 - $\vec{\mathcal{V}}_{k+1} = \vec{\mathcal{V}}_k - \alpha_k \vec{d}_k$
 - $\vec{r}_{k+1} = \vec{r}_k - \alpha_k (\mathbb{G}^T \cdot \vec{d}_k + \mathbb{C} \cdot \vec{p}_k)$
 - if \vec{r}_{k+1} is sufficiently small ($\|\vec{r}_{k+1}\|_2 / \|\vec{r}_0\|_2 < tol$), then exit loop
 - $\beta_k = (\vec{r}_{k+1}^T \cdot \vec{r}_{k+1}) / (\vec{r}_k^T \cdot \vec{r}_k)$
 - $\vec{p}_{k+1} = \vec{r}_{k+1} + \beta_k \vec{p}_k$
 - $k = k + 1$
- return $\vec{\mathcal{P}}_{k+1}$ as result

Remark. The matrix \mathbb{C} is rarely present unless for example when stabilised elements are used such as the stabilised $Q_1 \times P_0$ or the stabilised $Q_1 \times Q_1$ elements.

This iterative algorithm will converge to the solution with a rate which depends on the condition number of the \mathbb{S} matrix, which is not easy to compute since \mathbb{S} is never built. However, it has been established that large viscosity contrasts in the domain will have a negative impact on the convergence.

Remark. *This algorithm requires one solve with matrix \mathbb{K} per iteration but says nothing about the method employed to do so (direct or iterative solver) nor the corresponding preconditioner.*

One thing we know improves the convergence of any iterative solver is the use of a preconditioner matrix and therefore now focus on the Preconditioned Conjugate Gradient (PCG) method. Once again we turn to Wikipedia³⁶:

Algorithm as obtained from Wikipedia.

```

r0 := b − Ax0
z0 := M−1r0
p0 := z0
k := 0
repeat
   $\alpha_k := \frac{\mathbf{r}_k^\top \mathbf{z}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k}$ 
  xk+1 := xk +  $\alpha_k \mathbf{p}_k$ 
  rk+1 := rk −  $\alpha_k \mathbf{A} \mathbf{p}_k$ 
  if rk+1 is sufficiently small then exit loop end if
  zk+1 := M−1rk+1
   $\beta_k := \frac{\mathbf{z}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{z}_k^\top \mathbf{r}_k}$ 
  pk+1 := zk+1 +  $\beta_k \mathbf{p}_k$ 
  k := k + 1
end repeat
The result is xk+1

```

The same algorithm with our notations:

```

 $\vec{r}_0 = \vec{f} - \mathbb{S} \cdot \vec{\mathcal{P}}_0$ 
 $\vec{z}_0 = \vec{M}^{-1} \cdot \vec{r}_0$ 
 $\vec{p}_0 = \vec{z}_0$ 
k = 0
repeat
   $\alpha_k = (\vec{r}_k^\top \cdot \vec{z}_k) / (\vec{p}_k^\top \cdot \mathbb{S} \cdot \vec{p}_k)$ 
   $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$ 
   $\vec{r}_{k+1} = \vec{r}_k - \alpha_k \mathbb{S} \cdot \vec{p}_k$ 
   $\vec{z}_{k+1} = \vec{M}^{-1} \cdot \vec{r}_{k+1}$ 
   $\beta_k = (\vec{z}_{k+1}^\top \cdot \vec{r}_{k+1}) / (\vec{z}_k^\top \cdot \vec{r}_k)$ 
   $\vec{p}_{k+1} = \vec{z}_{k+1} + \beta_k \vec{p}_k$ 
  k = k + 1
end repeat
return  $\vec{\mathcal{P}}_{k+1}$  as the result

```

Unsurprisingly we find the same algorithm in Saad [1092]:

```

ALGORITHM 9.1 Preconditioned Conjugate Gradient
1. Compute  $r_0 := b - Ax_0$ ,  $z_0 = M^{-1}r_0$ , and  $p_0 := z_0$ 
2. For  $j = 0, 1, \dots$ , until convergence Do:
3.    $\alpha_j := (r_j, z_j) / (Ap_j, p_j)$ 
4.    $x_{j+1} := x_j + \alpha_j p_j$ 
5.    $r_{j+1} := r_j - \alpha_j Ap_j$ 
6.    $z_{j+1} := M^{-1}r_{j+1}$ 
7.    $\beta_j := (r_{j+1}, z_{j+1}) / (r_j, z_j)$ 
8.    $p_{j+1} := z_{j+1} + \beta_j p_j$ 
9. EndDo

```

Note that in the algorithm above the preconditioner matrix \mathbf{M} has to be symmetric positive-definite and fixed, i.e., cannot change from iteration to iteration. We see that this algorithm introduces an additional vector \vec{z} and a solve with the matrix \mathbf{M} at each iteration, which means that \mathbf{M} must be such that solving $\mathbf{M} \cdot \vec{x} = \vec{f}$ where \vec{f} is a given rhs vector must be cheap. Ultimately, the PCG algorithm applied to the Schur complement equation takes the form:

solver_pcg: assume $\vec{\mathcal{P}}_0$ known

- compute $\mathcal{V}_0 = \mathbb{K}^{-1}(f - \mathbb{G}\mathcal{P}_0)$
- $\vec{r}_0 = \mathbb{G}^T \mathcal{V}_0 - \mathbb{C} \cdot \vec{\mathcal{P}}_0 - \vec{h}$
- if \vec{r}_0 is sufficiently small, then return $(\vec{\mathcal{V}}_0, \vec{\mathcal{P}}_0)$ as the result

³⁶https://en.wikipedia.org/wiki/Conjugate_gradient_method

- $\vec{z}_0 = M^{-1} \cdot \vec{r}_0$
- $\vec{p}_0 = \vec{z}_0$
- $k = 0$
- repeat
 - compute $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$
 - solve $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
 - compute $\alpha_k = (\vec{r}_k^T \cdot \vec{z}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k + \vec{p}_k^T \cdot \mathbb{C} \cdot \vec{p}_k)$
 - $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$
 - $\vec{\mathcal{V}}_{k+1} = \vec{\mathcal{V}}_k - \alpha_k \vec{d}_k$
 - $\vec{r}_{k+1} = \vec{r}_k - \alpha_k (\mathbb{G}^T \cdot \vec{d}_k + \mathbb{C} \cdot \vec{p}_k)$
 - if \vec{r}_{k+1} is sufficiently small (i.e. $\|\vec{r}_{k+1}\|_2 / \|\vec{r}_0\|_2 < tol$), then exit loop
 - $\vec{z}_{k+1} = M^{-1} \cdot \vec{r}_{k+1}$
 - $\beta_k = (\vec{z}_{k+1}^T \cdot \vec{r}_{k+1}) / (\vec{z}_k^T \cdot \vec{r}_k)$
 - $\vec{p}_{k+1} = \vec{z}_{k+1} + \beta_k \vec{p}_k$
 - $k = k + 1$
- return $\vec{\mathcal{P}}_{k+1}$ as result

Following Zhong *et al.* [1415] one can define the following matrix as preconditioner:

$$\mathbf{M} = \text{diag} [\mathbb{G}^T (\text{diag}[\mathbb{K}])^{-1} \mathbb{G}]$$

which is the preconditioner used for the Citcom codes (see appendix ??). It can be constructed while the FEM matrix is being built/assembled and it is trivial to invert. The entries in $\text{diag}[\mathbb{K}]$ are the average viscosity in the elements associated with a given degree of freedom.

Another very cheap way of building \mathbf{M} for $Q_1 \times P_0$ elements is to realise that the matrix \mathbb{S} has dimensions element surface/volume divided by viscosity. We can then postulate

$$M_{e,e} = \frac{|\Omega|_e}{\eta_e}$$

where e is an element and η_e is the (average viscosity) inside the element. For higher order elements, we need to use the pressure mass matrix.

These two preconditioners and two other variants are implemented in [STONE](#) 16 for $Q_1 \times P_0$ elements.

7.11.4 Generalized Conjugate Residual approach (Geenen *et al.* (2009))

This approach is presented in Geenen *et al.* (2009) [443]. The saddle point problem arising from the constrained Stokes equation is solved with a Krylov method, GCR [1329], right preconditioned (postconditioned) with a block triangular preconditioner (BTR) [132].

The preconditioner \mathbf{P} is given by

$$\mathbf{P} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ 0 & -\tilde{\mathbb{S}} \end{pmatrix}$$

The GCR algorithm [365] in this case is taken from Vuik *et al.* (2000) [1331] and makes use of the block triangular preconditioner as follows:

$$\vec{r}_0 = \vec{b} - \mathbf{A} \cdot \vec{x}^0$$

for $k=0,1,2,\dots$

$$\begin{aligned} & - \vec{s}^{k+1} = \mathbf{P}^{-1} \cdot \vec{r}^k \\ & - \vec{v}^{k+1} = \mathbf{A} \cdot \vec{s}^{k+1} \\ & - \text{for } i=0,1,\dots,k \\ & \quad * \vec{v}^{k+1} = \vec{v}^{k+1} - (\vec{v}^{k+1}, \vec{v}^i) \vec{v}^i \\ & \quad * \vec{s}^{k+1} = \vec{s}^{k+1} - (\vec{v}^{k+1}, \vec{v}^i) \vec{s}^i \\ & - \text{end for} \\ & - \vec{v}^{k+1} = \vec{v}^{k+1} / \|\vec{v}^{k+1}\|_2 \\ & - \vec{s}^{k+1} = \vec{s}^{k+1} / \|\vec{v}^{k+1}\|_2 \\ & - \vec{x}^{k+1} = \vec{x}^k + (\vec{v}^{k+1}, \vec{r}^k) \vec{s}^{k+1} \\ & - \vec{r}^{k+1} = \vec{r}^k - (\vec{v}^{k+1}, \vec{r}^k) \vec{v}^{k+1} \end{aligned}$$

end for

As explained in Geenen *et al.*, instead of constructing \mathbf{P}^{-1} explicitly and applying it to \vec{r} , we instead solve the system $\mathbf{P} \cdot \vec{s} = \vec{r}$. We first decompose \vec{r} and \vec{s} as follows:

$$\vec{r}^k = \begin{pmatrix} \vec{r}_v^k \\ \vec{r}_p^k \end{pmatrix} \quad \vec{s}^{k+1} = \begin{pmatrix} \vec{s}_v^{k+1} \\ \vec{s}_p^{k+1} \end{pmatrix}$$

so that we have to solve

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ 0 & -\tilde{\mathbb{S}} \end{pmatrix} \cdot \begin{pmatrix} \vec{s}_v^{k+1} \\ \vec{s}_p^{k+1} \end{pmatrix} = \begin{pmatrix} \vec{r}_v^k \\ \vec{r}_p^k \end{pmatrix}$$

This is actually rather trivial because of the upper triangular nature of the preconditioner \mathbf{P} . It immediately follows:

$$\tilde{\mathbb{S}} \cdot \vec{s}_p^{k+1} = -\vec{r}_p^k \tag{7.234}$$

$$\mathbb{K} \cdot \vec{s}_v^{k+1} = \vec{r}_v^k - \mathbb{G} \cdot \vec{s}_p^{k+1} \tag{7.235}$$

As before we now must specify how we solve the above two equations (and we must therefore make a choice about the approximate Schur complement $\tilde{\mathbb{S}}$).

In the paper they take \mathbf{M}_p , the pressure mass matrix scaled with the inverse of viscosity as an approximation to the Schur complement $\tilde{\mathbb{S}}$, which is spectrally equivalent. Note that sometimes this mass matrix can be lumped which makes solving with it trivial and fast.

The inner solve with \mathbb{K} is carried out with a CG solvers preconditioned with AMG. They state that “Using AMG as a preconditioner to CG for the subsystem solution guarantees h -independent convergence of the solver during the preconditioner construction phase.”

7.11.5 Using MINRES a la Burstedde *et al.* (2008)

This approach is presented in Burstedde *et al.* (2008) [190]. They state that neglecting the off-diagonal blocks motivates use of the symmetric positive definite preconditioner:

$$\mathbf{P} = \begin{pmatrix} \tilde{\mathbb{K}} & 0 \\ 0 & \tilde{\mathbb{S}} \end{pmatrix}$$

where $\tilde{\mathbb{K}}$ is a variable-viscosity discrete vector Laplacian approximation of \mathbb{K} (see explanations in [188]), which is motivated by the fact that for constant viscosity and Dirichlet boundary conditions, \mathbb{K} and $\tilde{\mathbb{K}}$ are equivalent. $\tilde{\mathbb{S}}$ is an approximation of the Schur complement given by a lumped mass matrix weighted by the inverse viscosity η^{-1} . The resulting diagonal matrix $\tilde{\mathbb{S}}$ is spectrally equivalent to \mathbb{S} [371]. They also use AMG as preconditioner for the inner solves.

Note that Burstedde *et al.* (2008) [190] relies on stabilised $Q_1 \times Q_1$ elements from Dohrmann & Bochev [336] so that their Stokes matrix does feature the associated $-\mathbb{C}$ block. Subsequent papers do so too, see Burstedde *et al.* (2009) [188], Burstedde *et al.* (2013) [189]. The same solver structure based on MINRES is used in these articles too.

7.11.6 The Augmented Lagrangian approach

see LaCoDe paper [322].

We start from the saddle point Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (7.236)$$

The AL method consists of subtracting $\lambda^{-1}\mathbb{M}_p \cdot \vec{\mathcal{P}}$ from the left and right-side of the mass conservation equation (where \mathbb{M}_p is the pressure mass matrix) and introducing the following iterative scheme:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\lambda^{-1}\mathbb{M}_p \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}}^{k+1} \\ \vec{\mathcal{P}}^{k+1} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} - \lambda^{-1}\mathbb{M}_p \cdot \vec{\mathcal{P}}^k \end{pmatrix} \quad (7.237)$$

where k is the iteration counter and λ is an artificial compressibility term which has the dimensions of dynamic viscosity. The choice of λ can be difficult as too low or too high a value yields either erroneous results and/or terribly ill-conditioned matrices. LaCoDe paper (!) use such a method and report that $\lambda = \max_{\Omega}(\eta)$ works well. Note that at convergence we have $\|\vec{\mathcal{P}}^{k+1} - \vec{\mathcal{P}}^k\| < \epsilon$ and then Eq.(7.237) converges to Eq.(7.236) and the velocity and pressure fields are solution of the unmodified system Eq.(7.236).

The introduction of this term serves one purpose: allowing us to solve the system in a segregated manner (i.e. computing successive iterates of the velocity and pressure fields until convergence is reached). The second line of Eq. (7.237) is

$$\mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} - \lambda^{-1}\mathbb{M}_p \cdot \vec{\mathcal{P}}^{k+1} = \vec{h} - \lambda^{-1}\mathbb{M}_p \cdot \vec{\mathcal{P}}^k$$

and can therefore be rewritten

$$\vec{\mathcal{P}}^{k+1} = \vec{\mathcal{P}}^k + \lambda\mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} - \vec{h})$$

We can then substitute this expression of $\vec{\mathcal{P}}^{k+1}$ in the first equation. This yields:

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}^{k+1} \quad (7.238)$$

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k + \lambda\mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} - \vec{h})) \quad (7.239)$$

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} + \lambda\mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k - \lambda\mathbb{M}_p^{-1}\vec{h}) \quad (7.240)$$

$$\underbrace{(\mathbb{K} + \lambda\mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T)}_{\tilde{\mathbb{K}}} \cdot \vec{\mathcal{V}}^{k+1} = \underbrace{\vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k - \lambda\mathbb{M}_p^{-1}\vec{h})}_{\vec{f}^{k+1}} \quad (7.241)$$

$$(7.242)$$

The iterative algorithm goes as follows:

1. if it is the first timestep, set $\vec{\mathcal{P}}^0 = 0$, otherwise set it to the pressure of the previous timestep.
2. calculate $\tilde{\mathbb{K}}$
3. calculate \vec{f}^{k+1}
4. solve $\tilde{\mathbb{K}} \cdot \vec{\mathcal{V}}^{k+1} = \vec{f}^{k+1}$
5. update pressure with $\vec{\mathcal{P}}^{k+1} = \vec{\mathcal{P}}^k + \lambda \mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} - \vec{h})$

Remark. *If discontinuous pressures are used, the pressure mass matrix can be inverted element by element which is cheaper than inverting \mathbb{M}_p as a whole.*

Remark. *This method has obvious ties with the penalty method.*

Remark. *If $\lambda \gg \max_{\Omega} \eta$ then the matrix $\tilde{\mathbb{K}}$ is ill-conditioned and an iterative solver must be used.*

7.11.7 The SIMPLE method

simple.tex

What follows is borrowed from Volker John. *Finite Element Methods for Incompressible Flow Problems*. Springer, 2016. ISBN: 978-3-319-45749-9. DOI: 10.1007/978-3-319-45750-5, page 666.

The SIMPLE method (Semi-Implicit Method for Pressure-Linked Equations) has been introduced by Patankar & Spalding (1972) [981] as an iterative method to solve the finite volume discretized incompressible Navier-Stokes equations.

The algorithm is based on the following steps:

- First the pressure is assumed to be known from the previous iteration.
- Then the velocity is solved from the momentum equations. The newly obtained velocities do not satisfy the continuity equation since the pressure is only a guess.
- In the next substeps the velocities and pressures are corrected in order to satisfy the discrete continuity equation.

SIMPLE relies on the block LU decomposition

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\mathbb{S} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I} & \mathbb{K}^{-1} \cdot \mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (7.243)$$

The approximation \mathbb{K}^{-1} as $\mathbf{D}_{\mathbb{K}}^{-1} = (\text{diag}(\mathbb{K}))^{-1}$ leads to the SIMPLE algorithm. In this case the approximation of the Schur complement matrix is given by $\tilde{\mathbb{S}} = \mathbb{G}^T \cdot \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G}$ and the decomposition looks like

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \simeq \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\tilde{\mathbb{S}} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I} & \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix}$$

Thus one iteration of SIMPLE solves the following system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \simeq \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\tilde{\mathbb{S}} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I} & \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix}$$

Before we can write out the SIMPLE algorithm, we must first take a small detour via so-called distributive iterative methods [1331, 626]. Let us consider the linear system

$$\mathbf{A} \cdot \vec{x} = \vec{b}$$

A stationary iterative method is defined as follows:

$$\vec{x}^{k+1} = \mathbf{B} \cdot \vec{x}^k + \vec{c}$$

where $\vec{c} = (\mathbf{I} - \mathbf{B}) \cdot \mathbf{A}^{-1} \cdot \vec{b}$. Left-multiplying all terms by $(\mathbf{I} - \mathbf{B})^{-1}$ first and then left-multiplying again by \mathbf{A} we arrive at:

$$\mathbf{A} \cdot (\mathbf{I} - \mathbf{B})^{-1} \cdot \vec{x}^{k+1} = \mathbf{A} \cdot (\mathbf{I} - \mathbf{B})^{-1} \cdot \mathbf{B} \cdot \vec{x}^k + \mathbf{A} \cdot (\mathbf{I} - \mathbf{B})^{-1} \cdot \vec{c}$$

We define $\mathbf{M} = \mathbf{A} \cdot (\mathbf{I} - \mathbf{B})^{-1}$ so that now

$$\mathbf{M} \cdot \vec{x}^{k+1} = \mathbf{M} \cdot \mathbf{B} \cdot \vec{x}^k + \vec{b}$$

We define $\mathbf{N} = \mathbf{M} \cdot \mathbf{B}$ and finally

$$\mathbf{M} \cdot \vec{x}^{k+1} = \mathbf{N} \cdot \vec{x}^k + \vec{b}$$

Note that $\mathbf{M} - \mathbf{N} = \mathbf{M} - \mathbf{M} \cdot \mathbf{B} = \mathbf{M} \cdot (\mathbf{I} - \mathbf{B}) = \mathbf{A} \cdot (\mathbf{I} - \mathbf{B})^{-1} \cdot (\mathbf{I} - \mathbf{B}) = \mathbf{A}$. Let us now write the original system $\mathbf{A} \cdot \vec{x} = \vec{b}$ as $(\mathbf{A} \cdot \mathbf{B}) \cdot (\mathbf{B}^{-1} \cdot \vec{x}) = \vec{b}$ or, $\underline{\mathbf{A}} \cdot \underline{\vec{x}} = \vec{b}$ with $\vec{x} = \mathbf{B} \cdot \underline{\vec{x}}$ and $\underline{\mathbf{A}} = \mathbf{A} \cdot \mathbf{B}$. Splitting $\underline{\mathbf{A}} = \mathbf{M} - \mathbf{N}$ again yields

$$\mathbf{M} \cdot \underline{\vec{x}}^{k+1} = \mathbf{N} \cdot \underline{\vec{x}}^k + \vec{b}$$

Using $\vec{x} = \mathbf{B} \cdot \underline{\vec{x}}$, we get

$$\mathbf{M} \cdot \mathbf{B}^{-1} \cdot \vec{x}^{k+1} = \mathbf{N} \cdot \mathbf{B}^{-1} \cdot \vec{x}^k + \vec{b}$$

We then have

$$\vec{x}^{k+1} = \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [\mathbf{N} \cdot \mathbf{B}^{-1} \cdot \vec{x}^k + \vec{b}] \quad (7.244)$$

$$= \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [(\mathbf{M} - \underline{\mathbf{A}}) \cdot \mathbf{B}^{-1} \cdot \vec{x}^k + \vec{b}] \quad (7.245)$$

$$= \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [(\mathbf{M} - \mathbf{A} \cdot \mathbf{B}) \cdot \mathbf{B}^{-1} \cdot \vec{x}^k + \vec{b}] \quad (7.246)$$

$$= \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [\mathbf{M} \cdot \mathbf{B}^{-1} \cdot \vec{x}^k - \mathbf{A} \cdot \mathbf{B} \cdot \mathbf{B}^{-1} \cdot \vec{x}^k + \vec{b}] \quad (7.247)$$

$$= \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [\mathbf{M} \cdot \mathbf{B}^{-1} \cdot \vec{x}^k - \mathbf{A} \cdot \vec{x}^k + \vec{b}] \quad (7.248)$$

$$= \vec{x}^k + \mathbf{B} \cdot \mathbf{M}^{-1} \cdot [\vec{b} - \mathbf{A} \cdot \vec{x}^k] \quad (7.249)$$

Finally, we have the following recursion:

$$\boxed{\vec{x}^{k+1} = \vec{x}^k + \mathbf{B} \cdot \mathbf{M}^{-1} \cdot (\vec{b} - \mathbf{A} \cdot \vec{x}^k)} \quad (7.250)$$

Coming back to the SIMPLE algorithm, we start from

$$\mathbf{A} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix}$$

The matrix \mathbf{B} is then chosen to be

$$\mathbf{B} = \begin{pmatrix} \mathbf{I} & -\mathbb{K}^{-1}\mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix}$$

We then have

$$\mathbf{A} \cdot \mathbf{B} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I} & -\mathbb{K}^{-1}\mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\mathbb{S} \end{pmatrix}$$

where $\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$. Let us recall that we define $\mathbf{D}_{\mathbb{K}} = \text{diag}(\mathbb{K})$ and $\hat{\mathbb{S}} = \mathbb{G}^T \cdot \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G}$. We further define

$$\mathbf{M} = \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\hat{\mathbb{S}} \end{pmatrix}$$

and \mathbf{N} follows from the splitting $\mathbf{A} \cdot \mathbf{B} = \mathbf{M} - \mathbf{N}$. (Note that we do not need to form nor use \mathbf{N}).

The standard SIMPLE algorithm also replaces \mathbb{K}^{-1} by $\mathbf{D}_{\mathbb{K}}^{-1}$ in \mathbf{B} so that \mathbf{B} is approximated by:

$$\mathbf{B} = \begin{pmatrix} \mathbf{I} & -\mathbf{D}_{\mathbb{K}}^{-1}\mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix}$$

in the iterations. We can define

$$\vec{r}^k = \vec{b} - \mathbf{A} \cdot \vec{x}^k = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} - \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{v}^k \\ \vec{p}^k \end{pmatrix} = \begin{pmatrix} \vec{r}_{\mathcal{V}}^k \\ \vec{r}_{\mathcal{P}}^k \end{pmatrix}$$

The iteration loop of Eq. (7.250) then takes the form

$$\begin{pmatrix} \vec{v}^{k+1} \\ \vec{p}^{k+1} \end{pmatrix} = \begin{pmatrix} V^k \\ P^k \end{pmatrix} + \mathbf{B} \mathbf{M}^{-1} \begin{pmatrix} r_V^k \\ r_P^k \end{pmatrix} = \begin{pmatrix} V^k \\ P^k \end{pmatrix} + \begin{pmatrix} \delta V^k \\ \delta P^k \end{pmatrix} \quad \text{with} \quad \begin{pmatrix} \delta V^k \\ \delta P^k \end{pmatrix} = \mathbf{B} \mathbf{M}^{-1} \begin{pmatrix} \vec{r}_{\mathcal{V}}^k \\ \vec{r}_{\mathcal{P}}^k \end{pmatrix}$$

This last equation can be rewritten³⁷:

$$\mathbf{M} \cdot \left[\mathbf{B}^{-1} \cdot \begin{pmatrix} \delta \vec{\mathcal{V}}^k \\ \delta \vec{\mathcal{P}}^k \end{pmatrix} \right] = \begin{pmatrix} \vec{r}_{\mathcal{V}}^* \\ \vec{r}_{\mathcal{P}}^* \end{pmatrix}$$

We then have to solve

$$\mathbf{M} \cdot \begin{pmatrix} \delta^* \vec{\mathcal{V}}^k \\ \delta^* \vec{\mathcal{P}}^k \end{pmatrix} = \begin{pmatrix} \mathbb{K} & 0 \\ \mathbb{G}^T & -\hat{\mathbb{S}} \end{pmatrix} \cdot \begin{pmatrix} \delta^* \vec{\mathcal{V}}^k \\ \delta^* \vec{\mathcal{P}}^k \end{pmatrix} = \begin{pmatrix} \vec{r}_{\mathcal{V}}^* \\ \vec{r}_{\mathcal{P}}^* \end{pmatrix} \quad (7.251)$$

and then compute

$$\begin{pmatrix} \delta \vec{\mathcal{V}}^k \\ \delta \vec{\mathcal{P}}^k \end{pmatrix} = \mathbf{B} \begin{pmatrix} \delta^* \vec{\mathcal{V}}^k \\ \delta^* \vec{\mathcal{P}}^k \end{pmatrix} \quad (7.252)$$

Fortunately Eq. (7.251) translates into:

$$\mathbb{K} \cdot \delta^* \vec{\mathcal{V}}^k = \vec{r}_{\mathcal{V}}^* \quad (7.253)$$

$$\hat{\mathbb{S}} \cdot \delta^* \vec{\mathcal{P}}^k = -\vec{r}_{\mathcal{P}}^* + \mathbb{G}^T \cdot \delta^* \vec{\mathcal{V}}^k \quad (7.254)$$

and Eq. (7.252) translates into:

$$\begin{pmatrix} \delta \vec{\mathcal{V}}^k \\ \delta \vec{\mathcal{P}}^k \end{pmatrix} = \begin{pmatrix} \mathbf{I} & -\mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G} \\ 0 & \mathbf{I} \end{pmatrix} \cdot \begin{pmatrix} \delta^* \vec{\mathcal{V}}^k \\ \delta^* \vec{\mathcal{P}}^k \end{pmatrix}$$

or,

$$\delta \vec{\mathcal{V}}_k = \delta^* \vec{\mathcal{V}}^k - \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G} \cdot \delta^* \vec{\mathcal{P}}_k \quad (7.255)$$

$$\delta \vec{\mathcal{P}}_k = \delta^* \vec{\mathcal{P}}^k \quad (7.256)$$

The final algorithm will then look as follows:

1. compute the residuals

$$\begin{aligned} \vec{r}_{\mathcal{V}} &= \vec{f} - \mathbb{K} \cdot \vec{\mathcal{V}}^{(k)} - \mathbb{G} \cdot \vec{\mathcal{P}}^{(k)} \\ \vec{r}_{\mathcal{P}} &= \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}^{(k)} \end{aligned} \quad (7.257)$$

2. Solve $\mathbb{K} \cdot \delta^* \vec{\mathcal{V}}^k = \vec{r}_{\mathcal{V}}^*$
3. Solve $\hat{\mathbb{S}} \cdot \delta^* \vec{\mathcal{P}}^k = \vec{r}_{\mathcal{P}}^* - \mathbb{G}^T \cdot \delta^* \vec{\mathcal{V}}^k$
4. Compute $\delta \vec{\mathcal{V}}^k = \delta^* \vec{\mathcal{V}}^k - \mathbf{D}_{\mathbb{K}}^{-1} \cdot \mathbb{G} \cdot \delta^* \vec{\mathcal{P}}_k$
5. Update $\delta \vec{\mathcal{P}}^k = \delta^* \vec{\mathcal{P}}^k$
6. Update

$$\begin{aligned} \vec{\mathcal{V}}^{(k+1)} &= \vec{\mathcal{V}}^{(k)} + \omega_{\mathcal{V}} \delta \vec{\mathcal{V}}^k \\ \vec{\mathcal{P}}^{(k+1)} &= \vec{\mathcal{P}}^{(k)} + \omega_{\mathcal{P}} \delta \vec{\mathcal{P}}^k \end{aligned} \quad (7.258)$$

where the parameters $\omega_{\mathcal{V}}$ and $\omega_{\mathcal{P}}$ are between 0 and 1.

Note that SIMPLE can be used as left and as right preconditioner, see page 669 of John [650].

³⁷Remember that $(\mathbf{A} \cdot \mathbf{B})^{-1} = \mathbf{B}^{-1} \cdot \mathbf{A}^{-1}$

Also, John states that: “SIMPLE is easily to implement, which makes it attractive. It relies on the already assembled matrix blocks. Only the approximation $\hat{\mathbb{S}}$ of the Schur complement matrix has to be computed. This matrix couples pressure degrees of freedom that are usually not coupled in finite element approximations of the diffusion operator, but it is still a sparse matrix. The efficiency of SIMPLE depends on how good \mathbb{K}^{-1} is approximated by its diagonal.”




Relevant Literature:

- C. Echevarria Serur. “Fast iterative methods for solving the incompressible Navier-Stokes equations”. PhD thesis. TU Delft, 2013,
- D. Braess and R. Sarazin. “An Efficient Smoother for the Stokes Problem”. In: *Applied Numerical Math.* 23 (1997), pp. 3–20,
- M. ur Rehman, C. Vuik, and G. Segal. “SIMPLE-type preconditioners for the Oseen problem”. In: *International Journal for Numerical Methods in Fluids* 61 (2009), pp. 432–452 for SIMPLE(R) algorithm,
- Alik Ismail-Zadeh and Paul Tackley. *Computational Methods for Geodynamics*. Cambridge University Press, 2010

7.11.8 The GMRES approach - NOT FINISHED

The Generalized Minimal Residual method [1093] is an extension of MINRES (which is only applicable to symmetric systems) to unsymmetric systems. Like MINRES, it generates a sequence of orthogonal vectors and combines these through a least-squares solve and update. However, in the absence of symmetry this can no longer be done with short recurrences. As a consequence, all previously computed vectors in the orthogonal sequence have to be retained and for this reason "restarted" versions of the method are used.

It must be said that the (preconditioned) GMRES method is actually much more difficult to implement than the (preconditioned) Conjugate Gradient method. However, since it can deal with unsymmetric matrices, it means that it can be applied directly to the Stokes system matrix (as opposed to the CG method which is used on the Schur complement equation).

 Relevant Literature: [364, p208] [1092, 1091] [48] [34]


finish GMRES algo description. not sure what to do, hard to explain, not easy to code.

Algorithm 3 GMRES algorithm

1. x_0 is an initial guess and $r_0 = b - Ax_0$
 2. For $j = 1, 2, 3, \dots$
 3. $s_i = P_{m,i-1}(A)r_{i-1}$,
 s_i be the approximate solution of $As = r_{i-1}$
obtained after m steps of an iterative method
 4. $v_i = As_i$
 5. For $j = 1$ to $i - 1$
 6. $\alpha = (v_i, v_j)$,
 7. $v_i = v_i - \alpha v_j, s_i = s_i - \alpha s_j$,
 8. End
 9. $v_i = v_i / \|v_i\|_2, s_i = s_i / \|v_i\|_2$
 10. $x_i = x_{i-1} + (r_{i-1}, v_i)s_i$;
 11. $r_i = r_{i-1} - (r_{i-1}, v_i)v_i$;
 12. End
-

Taken from ur Rehman, vuik & Segal.

the FGMRES approach [331]

 Relevant Literature[970, 846, 431, 716, 717, 720]

7.12 Boundary conditions

7.12.1 Imposing Dirichlet boundary conditions

howtobc.tex

Let us consider a quadrilateral element with one degree of freedom per node and let us assume that we are solving the temperature equation. The local matrix and right-hand side vector are given by

$$A_{el}(4 \times 4) \quad \text{and} \quad B_{el}(4)$$

Let us assume that we want to impose $\tilde{T} = 10$ on the third node (local coordinates numbering). For instance, having built A_{el} and B_{el} , the system looks like :

$$\begin{pmatrix} 3 & 1 & 6 & 9 \\ 5 & 2 & 2 & 8 \\ 7 & 4 & 11 & 2 \\ 9 & 6 & 4 & 3 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ 1 \\ 2 \end{pmatrix}$$

which can be rewritten

$$3T_1 + T_2 + 6T_3 + 9T_4 = 4$$

$$5T_1 + 2T_2 + 2T_3 + 8T_4 = 3$$

$$7T_1 + 4T_2 + 11T_3 + 2T_4 = 1$$

$$9T_1 + 6T_2 + 4T_3 + 3T_4 = 2$$

or,

$$3T_1 + T_2 + \quad + 9T_4 = 4 - 6T_3$$

$$5T_1 + 2T_2 + \quad + 8T_4 = 3 - 2T_3$$

$$7T_1 + 4T_2 + 11T_3 + 2T_4 = 1$$

$$9T_1 + 6T_2 + \quad + 3T_4 = 2 - 4T_3$$

- Technique 1: Replace the hereabove system by

$$\begin{pmatrix} 3 & 1 & 6 & 9 \\ 5 & 2 & 2 & 8 \\ 7 & 4 & 11 + 10^{12} & 2 \\ 9 & 6 & 4 & 3 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix} = \begin{pmatrix} 4 \\ 3 \\ \tilde{T} \times (11 + 10^{12}) \\ 2 \end{pmatrix}$$

- Technique 2: One can choose not to solve for T_3 anymore, i.e. not to consider it as a degree of freedom and therefore write:

$$3T_1 + T_2 + 9T_4 = 4 - 6T_3$$

$$5T_1 + 2T_2 + 8T_4 = 3 - 2T_3$$

$$9T_1 + 6T_2 + 3T_4 = 2 - 4T_3$$

- Technique 3: Since we want to impose $T_3 = 10$, then we can write

$$3T_1 + T_2 + \quad + 9T_4 = 4 - 6T_3$$

$$5T_1 + 2T_2 + \quad + 8T_4 = 3 - 2T_3$$

$$0 + 0 + T_3 + 0 = 10$$

$$9T_1 + 6T_2 + \quad + 3T_4 = 2 - 4T_3$$

and in matrix form :

$$\begin{pmatrix} 3 & 1 & 0 & 9 \\ 5 & 2 & 0 & 8 \\ 0 & 0 & 1 & 0 \\ 9 & 6 & 0 & 3 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{pmatrix} = \begin{pmatrix} 4 - A_{13}T_3 \\ 3 - A_{23}T_3 \\ 10 \\ 2 - A_{43}T_3 \end{pmatrix}$$

The first technique is not a good idea in practice as it introduces very large values and will likely derail the solver. The second option is somewhat difficult to implement as it means that elemental matrix and rhs sizes will change from element to element and it therefore requires more book-keeping. The third technique is the one adopted throughout this document.

As shown in Wu, Xu, and Li [1371] (2008), it is better to replace the 1 on the diagonal by the former diagonal term as it reduces the condition number of the matrix. The rhs must then be modified accordingly.

 **Relevant Literature** Behr (2004) [69]

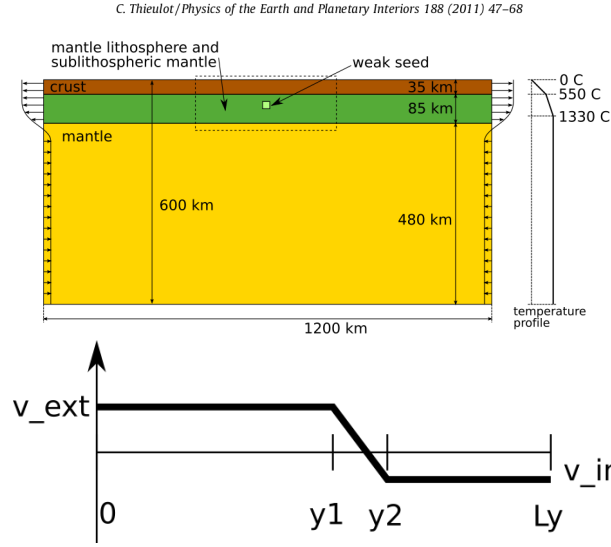
This is an excerpt of an email sent to me by Dave May in May 2014: *Never ever ever impose bc's using a penalty approach. For problems with a fixed mesh topology and time dependent Dirichlet domain (e.g. the segment of the boundary with Dirichlet bc's maybe change size/shape over time - for example with a true stick/slip type interface), it's nice to define the matrix with the dimension associated with the mesh+basis and leave all bc's in the operator. Leaving the bc's in the operator can be implemented in a manner which still retains the operators symmetry (assuming it was symmetric to begin with). This leaves the choice of what to stick on the diagonal. Simply using "1" could screw up the spectrum of the matrix and kill the iterative solver performance. A better choice would be to insert a diagonal entry closely related to the operator; e.g. something that looks like the diagonal entry of $\int 2\eta\epsilon(u) : \epsilon(v)dV$ (for the discrete stress tensor term).*

Removing Dirichlet bc's entirely for the discrete operator sounds attractive. The code will like the FE theory and you will only be solving for variables which are "unknowns" (compared with the above). However, introducing a time dependent Dirichlet domain means the matrix must be re-sized, as should its non-zero structure be re-allocated. Also, implemented multi-grid is annoying when the Dirichlet entries are removed. In fact, most of the code associated with stripping out Dirichlet bc's is annoying and ugly. However, removing the bcs ensures symmetry, it ensures the discrete operator will have a nice spectrum (c.f. the above option). Also, stripping out bcs usually increases overall storage as you have one representation of the discrete vectors given to the solver which will be of size $(N-n)$ and in your mesh you will have a representation of length N . "N" being the total number of dofs in your system, "n" being the number of Dirichlet constrained dofs in your system.

ASK
DAVE
for per-
mission

7.12.2 In-out flux boundary conditions for lithospheric models

kinematic_bc.tex



The velocity on the side is given by

$$\begin{aligned}
 u(y) &= v_{ext} & y < y_1 \\
 u(y) &= \frac{v_{in} - v_{ext}}{y_2 - y_1}(y - y_1) + v_{ext} & y_1 < y < y_2 \\
 u(y) &= v_{in} & y > y_2
 \end{aligned}$$

The requirement for volume conservation is:

$$\Phi = \int_0^{L_y} u(y) dy = 0$$

Having chosen v_{in} (the velocity of the plate), one can then compute v_{ext} as a function of y_1 and y_2 .

$$\begin{aligned}
 \Phi &= \int_0^{y_1} u(y) dy + \int_{y_1}^{y_2} u(y) dy + \int_{y_2}^{L_y} u(y) dy \\
 &= v_{ext} y_1 + \frac{1}{2}(v_{in} + v_{ext})(y_2 - y_1) + (L_y - y_2)v_{in} \\
 &= v_{ext} \left[y_1 + \frac{1}{2}(y_2 - y_1) \right] + v_{in} \left[\frac{1}{2}(y_2 - y_1) + (L_y - y_2) \right] \\
 &= v_{ext} \frac{1}{2}(y_1 + y_2) + v_{in} \left[L_y - \frac{1}{2}(y_1 + y_2) \right]
 \end{aligned}$$

and finally

$$v_{ext} = -v_{in} \frac{L_y - \frac{1}{2}(y_1 + y_2)}{\frac{1}{2}(y_1 + y_2)}$$

7.12.3 Periodic boundary conditions

This type of boundary conditions can be handy in some specific cases such as infinite domains. The idea is simple: when material leaves the domain through a boundary it comes back in through the opposite boundary (which of course presupposes a certain topology of the domain).

For instance, if one wants to model a gas at the molecular level and wishes to avoid interactions of the molecules with the walls of the container, such boundary conditions can be used, mimicking an infinite domain in all directions.

Let us consider the small mesh depicted hereunder:

missing picture

We wish to implement horizontal boundary conditions so that

$$u_5 = u_1 \quad u_{10} = u_6 \quad u_{15} = u_{11} \quad u_{20} = u_{16}$$

One could of course rewrite these conditions as constraints and extend the Stokes matrix but this approach turns out to be not practical at all.

Instead, the method is rather simple: replace in the connectivity array the dofs on the right side (nodes 5, 10, 15, 20) by the dofs on the left side. In essence, we wrap the system upon itself in the horizontal direction so that elements 4, 8 and 12 'see' and are 'made of' the nodes 1, 6, 11 and 16. In fact, this is only necessary during the assembly. Everywhere in the loops nodes 5, 10, 15 and 20 appear one must replace them by their left pendants 1, 6, 11 and 16. This automatically generates a matrix with lines and columns corresponding to the u_5 , u_{10} , u_{15} and u_{20} being exactly zero. The Stokes matrix is the same size, the blocks are the same size and the symmetric character of the matrix is respected. However, there remains a problem. There are zeros on the diagonal of the above mentioned lines and columns. One must then place there 1 or a more appropriate value.

Another way of seeing this is as follows: let us assume we have built and assembled the Stokes matrix, and we want to impose periodic b.c. so that dof j and i are the same. The algorithm is composed of four steps:

1. add col j to col i
2. add row j to row i (including rhs)
3. zero out row j , col j
4. put average diagonal value on diagonal (j, j)

Remark. Unfortunately the non-zero pattern of the matrix with periodic b.c. is not the same as the matrix without periodic b.c.

7.12.4 Free-slip boundary conditions on annulus

fsbc_annulus.tex

In the context of geodynamical modelling we often wish to prescribed free-slip boundary conditions on a given boundary of the domain. If the domain is a rectangle which sides align with the Cartesian axis, then fixing $\mathbf{v}_x = 0$ or $\mathbf{v}_y = 0$ is simple and does indeed insure free-slip boundary conditions.

However the situation is much more complicated in the case of a curved boundary, such as for instance the inner and outer boundaries of an annulus or spherical shell.

If the curved boundary is a circular, the procedure is as follows:

1. identify the node on the boundary which is to be fixed.
2. compute its coordinate angle θ (and ϕ in 3D)
3. do a rotation so as to bring it back onto the x-axis (2D) or z-axis (3D)
4. apply free slip boundary condition (now easy since parallel or perpendicular to axis)
5. rotate back

This technique is implemented in `STONE ??`, `STONE ??` and `STONE ??`.

A few remarks about rotation matrices In a given plane, the counter-clockwise rotation matrix by and angle θ is defined by

$$\mathcal{R} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

The image of vector \vec{V} by a rotation of angle θ is given by

$$\vec{V}' = \mathcal{R} \cdot \vec{V}$$

Coordinate transformations of second-rank tensors involve the very same matrix as vector transforms. A transformation of the stress tensor $\boldsymbol{\sigma}$, from the reference xy -coordinate system to $\boldsymbol{\sigma}'$ in a new $x'y'$ -system is done as follows:

$$\boldsymbol{\sigma}' = \mathcal{R} \cdot \boldsymbol{\sigma} \cdot \mathcal{R}^T$$

[from Wikipedia] A basic rotation (also called elemental rotation) is a rotation about one of the axes of a Coordinate system. The following three basic rotation matrices rotate vectors by an angle α about the x-, y-, or z-axis, in three dimensions, using the right-hand rule which codifies their alternating signs.

$$\mathcal{R}_x(\alpha) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \alpha & -\sin \alpha \\ 0 & \sin \alpha & \cos \alpha \end{pmatrix}$$

$$\mathcal{R}_y(\alpha) = \begin{pmatrix} \cos \alpha & 0 & \sin \alpha \\ 0 & 1 & 0 \\ -\sin \alpha & 0 & \cos \alpha \end{pmatrix}$$

$$\mathcal{R}_z(\alpha) = \begin{pmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

In my ELEFANT code I first rotate around the z axis by and angle $-\phi$ and then around axis y by an angle $-\theta$ in the case of a spherical shell.

$$\mathcal{R}_y(-\theta) = \begin{pmatrix} \cos(-\theta) & 0 & \sin(-\theta) \\ 0 & 1 & 0 \\ -\sin(-\theta) & 0 & \cos(-\theta) \end{pmatrix} = \begin{pmatrix} \cos \theta & 0 & -\sin \theta \\ 0 & 1 & 0 \\ \sin \theta & 0 & \cos \theta \end{pmatrix}$$

$$\mathcal{R}_z(-\phi) = \begin{pmatrix} \cos(-\phi) & -\sin(-\phi) & 0 \\ \sin(-\phi) & \cos(-\phi) & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} \cos \phi & \sin \phi & 0 \\ -\sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

These are the **Rott** and **Rotp** matrices in the routines.

Relevant Literature

- Note that in some cases applying free slip boundary conditions on a curved boundary with a triangular mesh can be problematic as explained in Dione, Tibirna, and Urquiza [334] (2013).
- M.S. Engelman, R.L. Sani, and P.M. Gresho. “The implementation of normal and/or tangential boundary conditions in finite element codes for incompressible fluid flow”. In: *Int. J. Num. Meth. Fluids* 2 (1982), pp. 225–238. DOI: 10.1002/fld.1650020302
- M. Behr. “On the Application of Slip Boundary Condition on Curved Boundaries”. In: *Int. J. Num. Meth. Fluids* 45 (2004), pp. 43–51. DOI: 10.1002/fld.663 in which it is stated:
 1. *If the slip boundary coincides with a Cartesian coordinate plane, the implementation is trivial, with the equations corresponding to the normal component of velocity simply being dropped from the equation system.*
 2. *If the slip boundary does not coincide with a Cartesian coordinate plane, the equations corresponding to the velocity components at the boundary are locally aligned with the normal- tangent-bi-tangent coordinate system, and the normal component of velocity is set to zero. This procedure is described by Engelman, Sani, and Gresho [372] (1982), who also advocate the use of consistent normals for proper mass conservation.*

7.13 Open boundary conditions

openbc.tex

So-called open boundary conditions have a special meaning in computational geodynamics. They usually refer to the boundary conditions on the sides of Cartesian models, usually looking at subduction or rifting processes.

In the literature boundary conditions on the vertical sidewalls are usually

- no-slip (no flow at the boundary),
- free slip (impermeable);
- open to some particular form of through-flow.

Free slip is the most commonly used boundary condition while prescribed in- and outflow or periodic boundary conditions are also common. (REF?)

Taken from Chertova, Geenen, Berg, and Spakman [231] (2012): “Open boundaries for which the horizontal in- and outflow are defined by a fully internally developed flow, have hardly been used [...]. Such open boundaries basically prescribe a hydrostatic pressure condition on the boundary preventing the model to collapse while horizontal in and outflow is free, in the sense that it is driven by the internal dynamics and the usual condition of incompressible flow. Among the range of boundary conditions used, open boundaries may fit best to real-mantle flow conditions surrounding subduction zones.”

Two examples of the use of such boundary conditions were found in the literature: Quinteros, Sobolev, and Popov [1030] (2010) and Chertova, Geenen, Berg, and Spakman [231] (212).

We start again from the variational form of the momentum equation, and focus on the term containing the full stress tensor $\boldsymbol{\sigma}$. Let us look at the stress tensor gradient, multiplied by the basis function \mathcal{N} , integrated over the domain:

$$\begin{aligned} \int_V \mathcal{N} \vec{\nabla} \cdot \boldsymbol{\sigma} \, dV &= \int_V \left[\vec{\nabla} \cdot (\mathcal{N} \boldsymbol{\sigma}) - \vec{\nabla} \mathcal{N} \cdot \boldsymbol{\sigma} \right] \, dV \\ &= \int_V \vec{\nabla} \cdot (\mathcal{N} \boldsymbol{\sigma}) \, dV - \int_V \vec{\nabla} \mathcal{N} \cdot \boldsymbol{\sigma} \, dV \end{aligned} \quad (7.259)$$

The right term yields the \mathbb{K} and \mathbb{G} matrices after discretisation, as seen in Section ???. Turning to the left term, we then make use of the Green-Gauss divergence theorem³⁸ which states that for a continuously differentiable vector field \vec{F} :

$$\int_V (\vec{\nabla} \cdot \vec{F}) \, dV = \int_S \vec{F} \cdot \vec{n} \, dS$$

so that (applying it now to tensors):

$$\int_V \vec{\nabla} \cdot (\mathcal{N} \boldsymbol{\sigma}) \, dV = \int_S \mathcal{N} \boldsymbol{\sigma} \cdot \vec{n} \, dS$$

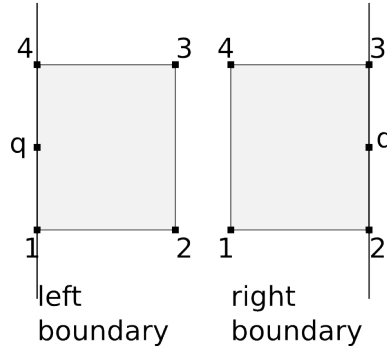
This right hand side term is responsible for the surface boundary conditions and cannot be neglected if one wishes to implement stress boundary conditions, such as the so-called open boundary conditions.

Note that in Liao and Gerya [786] (2017) the authors describe an iterative algorithm that allows them to control the actual force applied at the boundary by scaling the kinematical boundary conditions

³⁸https://en.wikipedia.org/wiki/Divergence_theorem

7.13.1 Two-dimensional case - $Q_1 \times P_0$ elements

On the following figure two elements are represented, one on the left boundary, one on the right boundary:



The prescribed traction on the left boundary is

$$\vec{t} = \boldsymbol{\sigma} \cdot \vec{n} = \begin{pmatrix} -p_{bc} & 0 \\ 0 & -p_{bc} \end{pmatrix} \cdot \begin{pmatrix} -1 \\ 0 \end{pmatrix} = \begin{pmatrix} p_{bc} \\ 0 \end{pmatrix}$$

The integral on the side of the element is then

$$\int_{\Gamma} \mathcal{N}_i \vec{t} dS$$

for $i = 1, 2, 3, 4$, which yields the following elemental rhs vector:

$$\vec{F}_{el} = \int_{\Gamma_{14}} \begin{pmatrix} \mathcal{N}_1(x, y) t_x(x, y) \\ \mathcal{N}_1(x, y) t_y(x, y) \\ \mathcal{N}_2(x, y) t_x(x, y) \\ \mathcal{N}_2(x, y) t_y(x, y) \\ \mathcal{N}_3(x, y) t_x(x, y) \\ \mathcal{N}_3(x, y) t_y(x, y) \\ \mathcal{N}_4(x, y) t_x(x, y) \\ \mathcal{N}_4(x, y) t_y(x, y) \end{pmatrix} dS$$

It is worth noting that the integral takes place on the edge Γ_{14} so that \mathcal{N}_2 and \mathcal{N}_3 are identically zero on this edge and also $t_y = 0$ so

$$\vec{F}_{el} = \begin{pmatrix} \int_{\Gamma_{14}} \mathcal{N}_1(x, y) t_x(x, y) dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{14}} \mathcal{N}_4(x, y) t_x(x, y) dS \\ 0 \end{pmatrix}$$

If the traction (applied pressure) is constant over the element, then

$$\vec{F}_{el} = t_x \begin{pmatrix} \int_{\Gamma_{14}} \mathcal{N}_1(x, y) dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{14}} \mathcal{N}_4(x, y) dS \\ 0 \end{pmatrix} = t_x \begin{pmatrix} \int_{y_1}^{y_4} \mathcal{N}_1(x, y) dy \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{y_1}^{y_4} \mathcal{N}_4(x, y) dy \\ 0 \end{pmatrix} = \frac{t_x h_y}{2} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

where h_y is the height of the element along the segment.

On the right boundary, we have $\mathcal{N}_2 = 0$ and $\mathcal{N}_3 = 0$, and since $t_y = 0$ then the corresponding additional elemental right hand side vector writes:

$$\vec{F}_{el} = -\frac{t_x h_y}{2} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

In the case where the traction is not constant over the edge, a numerical quadrature rule must be employed to integrate $\int_{\Gamma} \mathcal{N}_i t_x dS$.

7.13.2 Three-dimensional case - $Q_1 \times P_0$ elements

The right hand side is $ndof \times ndim = 8 \times 3 = 24$ long.

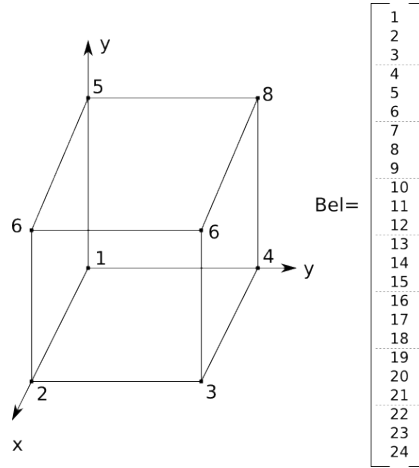


figure above has two node 6! there are two y axis in the figure and no z. redo with tikz

- The face $r = -1$ is made of nodes 1,4,5,8, so $\vec{n} = (-1, 0, 0)$. Since $t_y = 0$ and $t_z = 0$ and $\mathcal{N}_2 = \mathcal{N}_3 = \mathcal{N}_6 = \mathcal{N}_7$ on this face:

$$\vec{F}_{el} = \int_{\Gamma_{1458}} \begin{pmatrix} \mathcal{N}_1(x, y, z) t_x(x, y, z) \\ \mathcal{N}_1(x, y, z) t_y(x, y, z) \\ \mathcal{N}_1(x, y, z) t_z(x, y, z) \\ \mathcal{N}_2(x, y, z) t_x(x, y, z) \\ \mathcal{N}_2(x, y, z) t_y(x, y, z) \\ \mathcal{N}_2(x, y, z) t_z(x, y, z) \\ \mathcal{N}_3(x, y, z) t_x(x, y, z) \\ \mathcal{N}_3(x, y, z) t_y(x, y, z) \\ \mathcal{N}_3(x, y, z) t_z(x, y, z) \\ \mathcal{N}_4(x, y, z) t_x(x, y, z) \\ \mathcal{N}_4(x, y, z) t_y(x, y, z) \\ \mathcal{N}_4(x, y, z) t_z(x, y, z) \\ \mathcal{N}_5(x, y, z) t_x(x, y, z) \\ \mathcal{N}_5(x, y, z) t_y(x, y, z) \\ \mathcal{N}_5(x, y, z) t_z(x, y, z) \\ \mathcal{N}_6(x, y, z) t_x(x, y, z) \\ \mathcal{N}_6(x, y, z) t_y(x, y, z) \\ \mathcal{N}_6(x, y, z) t_z(x, y, z) \\ \mathcal{N}_7(x, y, z) t_x(x, y, z) \\ \mathcal{N}_7(x, y, z) t_y(x, y, z) \\ \mathcal{N}_7(x, y, z) t_z(x, y, z) \\ \mathcal{N}_8(x, y, z) t_x(x, y, z) \\ \mathcal{N}_8(x, y, z) t_y(x, y, z) \\ \mathcal{N}_8(x, y, z) t_z(x, y, z) \end{pmatrix} dS = \int_{\Gamma_{1458}} \begin{pmatrix} \mathcal{N}_1(x, y, z) t_x \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{N}_4(x, y, z) t_x \\ 0 \\ 0 \\ \mathcal{N}_5(x, y, z) t_x \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \mathcal{N}_8(x, y, z) t_x \\ 0 \\ 0 \end{pmatrix} dS = t_x \begin{pmatrix} \int_{\Gamma_{1458}} \mathcal{N}_1(x, y, z) dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{1458}} \mathcal{N}_4(x, y, z) dS \\ 0 \\ 0 \\ \int_{\Gamma_{1458}} \mathcal{N}_5(x, y, z) dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{1458}} \mathcal{N}_8(x, y, z) dS \\ 0 \\ 0 \end{pmatrix} = \frac{h_y h_z t_x}{4} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

- The face $r = +1$ is made of nodes 2,3,6,7, so $\vec{n} = (1, 0, 0)$. so the non-zero terms are in positions (4, 7, 16, 19).

- The face $s = -1$ is made of nodes 1,2,5,6, so $\vec{n} = (0, -1, 0)$. so the non-zero terms are in positions (2, 5, 14, 17).
- The face $s = +1$ is made of nodes 3,4,7,8, so $\vec{n} = (0, +1, 0)$. so the non-zero terms are in positions (8, 11, 20, 23).

7.13.3 Two-dimensional case - $Q_2 \times Q_1$ elements

We here too assume that we wish to prescribe a traction on the sides of a 2D domain which are aligned with the vertical axis.

constant traction

It is not fundamentally different, except that the element counts 9 nodes, so the vector is $9 \times 2 = 18$ long. The internal numbering of the nodes is as follows:

| velocity | pressure |
|-------------|----------|
| 3---6---2 | 3-----2 |
| | |
| 7 8 5 | |
| | |
| 0---4---1 | 0-----1 |

On the left boundary nodes 0,3,7 are involved while on the right boundary nodes 1,2,5 are. Assuming once again t_x constant over the edge and $t_y = 0$, we have on the left side:

$$\vec{F}_{el} = \int_{\Gamma_{073}} \begin{pmatrix} \mathcal{N}_0(x, y)t_x(x, y) \\ \mathcal{N}_0(x, y)t_y(x, y) \\ \mathcal{N}_1(x, y)t_x(x, y) \\ \mathcal{N}_1(x, y)t_y(x, y) \\ \mathcal{N}_2(x, y)t_x(x, y) \\ \mathcal{N}_2(x, y)t_y(x, y) \\ \mathcal{N}_3(x, y)t_x(x, y) \\ \mathcal{N}_3(x, y)t_y(x, y) \\ \mathcal{N}_4(x, y)t_x(x, y) \\ \mathcal{N}_4(x, y)t_y(x, y) \\ \mathcal{N}_5(x, y)t_x(x, y) \\ \mathcal{N}_5(x, y)t_y(x, y) \\ \mathcal{N}_6(x, y)t_x(x, y) \\ \mathcal{N}_6(x, y)t_y(x, y) \\ \mathcal{N}_7(x, y)t_x(x, y) \\ \mathcal{N}_7(x, y)t_y(x, y) \\ \mathcal{N}_8(x, y)t_x(x, y) \\ \mathcal{N}_8(x, y)t_y(x, y) \end{pmatrix} dS = t_x \begin{pmatrix} \int_{\Gamma_{073}} \mathcal{N}_0(x, y)dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{073}} \mathcal{N}_3(x, y)dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{073}} \mathcal{N}_7(x, y)dS \\ 0 \\ 0 \\ 0 \end{pmatrix} = t_x \frac{h_y}{2} \begin{pmatrix} \int_{-1}^{+1} \mathcal{N}_0(r = -1, s)ds \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{-1}^{+1} \mathcal{N}_3(r = -1, s)ds \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{-1}^{+1} \mathcal{N}_7(r = -1, s)ds \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

We then compute

$$\int_{-1}^{+1} \mathcal{N}_0(r = -1, s)ds = \int_{-1}^{+1} \frac{1}{2}s(s-1)ds = \frac{1}{3} \quad (7.260)$$

$$\int_{-1}^{+1} \mathcal{N}_3(r = -1, s)ds = \int_{-1}^{+1} \frac{1}{2}s(s+1)ds = \frac{1}{3} \quad (7.261)$$

$$\int_{-1}^{+1} \mathcal{N}_7(r = -1, s)ds = \int_{-1}^{+1} (1-s^2)ds = \frac{4}{3} \quad (7.262)$$

Note that the sum of the three terms is 2, as expected: on the edge we have $\mathcal{N}_0 + \mathcal{N}_3 + \mathcal{N}_7 = 1$ so that the integral of the sum over the interval $[-1,1]$ yields 2. Finally

$$\vec{F}_{el} = \frac{t_x h_y}{6} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 4 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

This is implemented in [STONE](#) 61, 64, 146, 148.

On the right boundary, we need to compute (careful with sign when implementing!)

$$\vec{F}_{el} = \int_{\Gamma_{152}} \begin{pmatrix} \mathcal{N}_0(x,y)t_x(x,y) \\ \mathcal{N}_0(x,y)t_y(x,y) \\ \mathcal{N}_1(x,y)t_x(x,y) \\ \mathcal{N}_1(x,y)t_y(x,y) \\ \mathcal{N}_2(x,y)t_x(x,y) \\ \mathcal{N}_2(x,y)t_y(x,y) \\ \mathcal{N}_3(x,y)t_x(x,y) \\ \mathcal{N}_3(x,y)t_y(x,y) \\ \mathcal{N}_4(x,y)t_x(x,y) \\ \mathcal{N}_4(x,y)t_y(x,y) \\ \mathcal{N}_5(x,y)t_x(x,y) \\ \mathcal{N}_5(x,y)t_y(x,y) \\ \mathcal{N}_6(x,y)t_x(x,y) \\ \mathcal{N}_6(x,y)t_y(x,y) \\ \mathcal{N}_7(x,y)t_x(x,y) \\ \mathcal{N}_7(x,y)t_y(x,y) \\ \mathcal{N}_8(x,y)t_x(x,y) \\ \mathcal{N}_8(x,y)t_y(x,y) \end{pmatrix} dS = t_x \begin{pmatrix} 0 \\ 0 \\ \int_{\Gamma_{125}} \mathcal{N}_1(x,y)dS \\ 0 \\ \int_{\Gamma_{125}} \mathcal{N}_2(x,y)dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{\Gamma_{125}} \mathcal{N}_5(x,y)dS \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = t_x \frac{h_y}{2} \begin{pmatrix} 0 \\ 0 \\ \int_{-1}^{+1} \mathcal{N}_1(-1,s)ds \\ 0 \\ \int_{-1}^{+1} \mathcal{N}_2(-1,s)ds \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \int_{-1}^{+1} \mathcal{N}_5(-1,s)ds \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \frac{t_x h_y}{6} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 4 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

linear traction

Let us now turn to the case where the traction we wish to apply on the boundary is not piecewise constant but linear. We set $t_x(y) = ay + b$, so that on the right side (nodes, 1,2,5), we have to compute

$$\begin{aligned}
\int_{\Gamma_{125}} \mathcal{N}_1(x, y) t_x(y) dS &= \int_{\Gamma_{125}} \mathcal{N}_1(x, y) (ay + b) dy \\
&= \frac{h_y}{2} \int_{-1}^1 \mathcal{N}_1(r = 1, s) [ay(s) + b] ds
\end{aligned}$$

We have

$$s(y) = \frac{2}{h_y}(y - y_1) - 1 \quad \text{or} \quad y(s) = \frac{h_y}{2}(s + 1) + y_1$$

and

$$\mathcal{N}_1(r, s) = \frac{1}{2}r(r + 1)\frac{1}{2}s(s - 1) \quad \Rightarrow \quad \mathcal{N}_1(r = 1, s) = \frac{1}{2}1(1 + 1)\frac{1}{2}s(s - 1) = \frac{1}{2}s(s - 1)$$

Then

$$\begin{aligned}
\int_{\Gamma_{125}} \mathcal{N}_1(x, y) t_x(y) dS &= \frac{h_y}{2} \int_{-1}^1 \mathcal{N}_1(r = 1, s) \left[a \left(\frac{h_y}{2}(s + 1) + y_1 \right) + b \right] ds \\
&= \frac{h_y}{2} \int_{-1}^1 \frac{1}{2}s(s - 1) \left[a \left(\frac{h_y}{2}(s + 1) + y_1 \right) + b \right] ds \\
&= \frac{h_y}{4} \int_{-1}^1 s(s - 1) \left[\frac{ah_y}{2}(s + 1) + (ay_1 + b) \right] ds \\
&= \frac{h_y}{4} \left[\frac{ah_y}{2} \int_{-1}^1 s(s - 1)(s + 1) ds + (ay_1 + b) \int_{-1}^1 s(s - 1) ds \right] \\
&= \frac{h_y}{4} \frac{ah_y}{2} \underbrace{\int_{-1}^1 s(s^2 - 1) ds}_{=0} + \frac{h_y}{4} (ay_1 + b) \underbrace{\int_{-1}^1 s(s - 1) ds}_{=-2/3} \\
&= \frac{h_y}{6} (ay_1 + b)
\end{aligned} \tag{7.263}$$

Let us now turn to \mathcal{N}_2 :

$$\mathcal{N}_2(r, s) = \frac{1}{2}r(r + 1)\frac{1}{2}s(s + 1) \quad \Rightarrow \quad \mathcal{N}_2(r = 1, s) = \frac{1}{2}s(s + 1)$$

Then

$$\begin{aligned}
\int_{\Gamma_{125}} \mathcal{N}_2(x, y) t_x(y) dS &= \frac{h_y}{2} \int_{-1}^1 \mathcal{N}_2(r = 1, s) \left[a \left(\frac{h_y}{2}(s + 1) + y_1 \right) + b \right] ds \\
&= \frac{h_y}{2} \int_{-1}^1 \frac{1}{2}s(s + 1) \left[a \left(\frac{h_y}{2}(s + 1) + y_1 \right) + b \right] ds \\
&= \frac{h_y}{4} \int_{-1}^1 s(s + 1) \left[\frac{ah_y}{2}(s + 1) + (ay_1 + b) \right] ds \\
&= \frac{h_y}{4} \left[\frac{ah_y}{2} \int_{-1}^1 s(s + 1)(s + 1) ds + (ay_1 + b) \int_{-1}^1 s(s + 1) ds \right] \\
&= \frac{h_y}{4} \frac{ah_y}{2} \underbrace{\int_{-1}^1 s(s + 1)^2 ds}_{=4/3} + \frac{h_y}{4} (ay_1 + b) \underbrace{\int_{-1}^1 s(s + 1) ds}_{=2/3} \\
&= \frac{h_y}{6} (ah_y + ay_1 + b)
\end{aligned} \tag{7.264}$$

And finally let us turn to \mathcal{N}_5 :

$$\mathcal{N}_5(r, s) = \frac{1}{2}r(r+1)(1-s^2) \quad \Rightarrow \quad \mathcal{N}_5(r=1, s) = (1-s^2)$$

then

$$\begin{aligned} \int_{\Gamma_{125}} \mathcal{N}_5(x, y) t_x(y) dS &= \frac{h_y}{2} \int_{-1}^1 \mathcal{N}_2(r=1, s) \left[a \left(\frac{h_y}{2}(s+1) + y_1 \right) + b \right] ds \\ &= \frac{h_y}{2} \int_{-1}^1 (1-s^2) \left[a \left(\frac{h_y}{2}(s+1) + y_1 \right) + b \right] ds \\ &= \frac{h_y}{2} \int_{-1}^1 (1-s^2) \left[\frac{ah_y}{2}(s+1) + (ay_1 + b) \right] ds \\ &= \frac{h_y}{2} \left[\frac{ah_y}{2} \int_{-1}^1 (1-s^2)(s+1) ds + (ay_1 + b) \int_{-1}^1 (1-s^2) ds \right] \\ &= \frac{h_y}{2} \frac{ah_y}{2} \underbrace{\int_{-1}^1 (1-s^2)(1+s) ds}_{=4/3} + \frac{h_y}{2} (ay_1 + b) \underbrace{\int_{-1}^1 (1-s^2) ds}_{=4/3} \\ &= \frac{h_y}{6} (2ah_y + 4ay_1 + 4b) \end{aligned} \tag{7.265}$$

Note that by setting $a = 0$ and $b = t_x$ we recover the expressions above for a piecewise constant value.

If we know p_1 and p_2 (say, for example that the lithostatic pressure has been computed on these nodes and we wish to prescribe it on the side) then

$$t_y = ay + b = \underbrace{\frac{p_2 - p_1}{y_2 - y_1}}_{=a} y + \underbrace{p_1 - \frac{p_2 - p_1}{y_2 - y_1} y_1}_{=b}$$

On the left side (nodes 0,7,3), we have to compute

$$\int_{073} \mathcal{N}_0(x, y) (ay + b) dS = \frac{h_y}{2} \int_{-1}^{+1} \mathcal{N}_0(r = -1, s) [ay(s) + b] ds \tag{7.266}$$

$$\int_{073} \mathcal{N}_3(x, y) (ay + b) dS = \frac{h_y}{2} \int_{-1}^{+1} \mathcal{N}_3(r = -1, s) [ay(s) + b] ds \tag{7.267}$$

$$\int_{073} \mathcal{N}_7(x, y) (ay + b) dS = \frac{h_y}{2} \int_{-1}^{+1} \mathcal{N}_7(r = -1, s) [ay(s) + b] ds \tag{7.268}$$

with

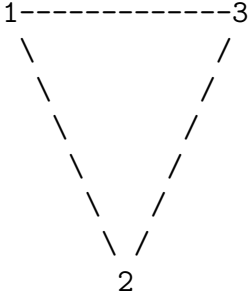
$$\begin{aligned} \mathcal{N}_0(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s-1) \rightarrow \mathcal{N}_0(-1, s) &= \frac{1}{2}s(s-1) \\ \mathcal{N}_3(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}(1-s^2) \rightarrow \mathcal{N}_3(-1, s) &= (1-s^2) \\ \mathcal{N}_7(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s+1) \rightarrow \mathcal{N}_7(-1, s) &= \frac{1}{2}s(s+1) \end{aligned}$$

If we know p_0 and p_3 then

$$t_x = ay + b = \underbrace{\frac{p_3 - p_0}{y_3 - y_0}}_{=a} y + \underbrace{p_0 - \frac{p_3 - p_0}{y_3 - y_0} y_0}_{=b}$$

7.13.4 Two-dimensional case - Linear triangle elements

Let us assume we want to apply a stress on the face 13 of the following element:



The integral on the side of the element is $\int_{\Gamma} \mathcal{N}_i \vec{t} dS$ for $i = 1, 2, 3$, which yields the following elemental rhs vector:

$$\vec{F}_{el} = \int_{\Gamma_{13}} \begin{pmatrix} \mathcal{N}_1(x, y)t_x(x, y) \\ \mathcal{N}_1(x, y)t_y(x, y) \\ \mathcal{N}_2(x, y)t_x(x, y) \\ \mathcal{N}_2(x, y)t_y(x, y) \\ \mathcal{N}_3(x, y)t_x(x, y) \\ \mathcal{N}_3(x, y)t_y(x, y) \end{pmatrix} dS = \int_{\Gamma_{13}} \begin{pmatrix} 0 \\ \mathcal{N}_1(x, y)t_y(x, y) \\ 0 \\ 0 \\ 0 \\ \mathcal{N}_3(x, y)t_y(x, y) \end{pmatrix} dS$$

since $t_x = 0$ and there function \mathcal{N}_2 will be zero on the edge.

We also arbitrarily set $y_1 = y_3 = 0$. We have seen in Section ?? that the basis functions (expressed as a function of the real coordinates x, y) for a linear triangle are given by:

$$\begin{aligned} \mathcal{N}_1(x, y) &= \frac{1}{D}[(x_2y_3 - x_3y_2) + (y_2 - y_3)x + (x_3 - x_2)y] \\ \mathcal{N}_2(x, y) &= \frac{1}{D}[(x_3y_1 - x_1y_3) + (y_3 - y_1)x + (x_1 - x_3)y] \\ \mathcal{N}_3(x, y) &= \frac{1}{D}[(x_1y_2 - x_2y_1) + (y_1 - y_2)x + (x_2 - x_1)y] \end{aligned}$$

with

$$D = \begin{vmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ 1 & x_3 & y_3 \end{vmatrix} = \begin{vmatrix} 1 & x_1 & 0 \\ 1 & x_2 & y_2 \\ 1 & x_3 & 0 \end{vmatrix} = -x_3y_2 + x_1y_2 = y_2(x_1 - x_3)$$

$$\begin{aligned}
\int_{x_1}^{x_3} \mathcal{N}_1(x, y=0) dx &= \frac{1}{D} \int_{x_1}^{x_3} [(x_2 y_3 - x_3 y_2) + (y_2 - y_3)x] \\
&= \frac{1}{D} \int_{x_1}^{x_3} [-x_3 y_2 + y_2 x] \quad \text{since } y_1 = y_3 = 0 \\
&= \frac{y_2}{D} \int_{x_1}^{x_3} (-x_3 + x) dx \\
&= \frac{y_2}{y_2(x_1 - x_3)} \left[-x_3 x + \frac{1}{2} x^2 \right]_{x_1}^{x_3} \\
&= \frac{1}{x_1 - x_3} \left[-x_3(x_3 - x_1) + \frac{1}{2}(x_3^2 - x_1^2) \right] \\
&= \frac{1}{x_1 - x_3} \left[x_3(x_1 - x_3) + \frac{1}{2}(x_3 - x_1)(x_3 + x_1) \right] \\
&= x_3 - \frac{1}{2}(x_3 + x_1) \\
&= \frac{1}{2}(x_3 - x_1) \tag{7.269}
\end{aligned}$$

$$\begin{aligned}
\int_{x_1}^{x_3} \mathcal{N}_3(x, y=0) dx &= \frac{1}{D} \int_{x_1}^{x_3} [(x_1 y_2 - x_2 y_1) + (y_1 - y_2)x] dx \\
&= \frac{1}{D} \int_{x_1}^{x_3} [x_1 y_2 - y_2 x] dx \\
&= \frac{y_2}{D} \int_{x_1}^{x_3} [x_1 - x] dx \\
&= \frac{y_2}{y_2(x_1 - x_3)} \left[x_1 x - \frac{1}{2} x^2 \right]_{x_1}^{x_3} \\
&= \frac{1}{x_1 - x_3} \left[x_1(x_3 - x_1) - \frac{1}{2}(x_3^2 - x_1^2) \right] \\
&= -x_1 + \frac{1}{2}(x_3 + x_1) \\
&= \frac{1}{2}(x_3 - x_1) \tag{7.270}
\end{aligned}$$

Finally

$$\vec{F}_{el} = \frac{ht_y}{2} \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

7.14 About nullspaces

7.14.1 Pressure normalisation, nullspace

pressure_normlisation.tex

Basic idea and naive implementation

When Dirichlet boundary conditions are imposed everywhere on the boundary, pressure is only present by its gradient in the equations. It is thus determined up to an arbitrary constant (one speaks then of a nullspace of size 1). In such a case, one commonly impose the average of the pressure over the whole domain or on a subset of the boundary to have a zero average, i.e.

$$\int_{\Omega} p \, dV = 0 \quad (7.271)$$

Let us assume for example/simplicity that we are using $Q_1 \times P_0$ elements. The pressure is constant inside each element so the integral above becomes:

$$\int_{\Omega} p \, dV = \sum_e \int_{\Omega_e} p \, dV = \sum_e p_e \int_{\Omega_e} dV = \sum_e p_e V_e = 0 \quad (7.272)$$

where the sum runs over all elements e of volume V_e . This can be rewritten

$$\vec{L} \cdot \vec{\mathcal{P}} = 0$$

and it is a constraint on the pressure solution which couples *all* pressure dofs. We can associate to it a Lagrange multiplier λ so that we must solve the modified Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} & 0 \\ \mathbb{G}^T & 0 & \vec{L}^T \\ 0 & \vec{L} & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \\ \lambda \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \\ 0 \end{pmatrix}$$

When higher order spaces are used for pressure (continuous or discontinuous) one must then carry out the above integration numerically by means of (usually) a Gauss-Legendre quadrature.

Although valid, this approach has one main disadvantage: it makes the Stokes matrix larger (although marginally so – only one row and column are added), but more importantly it prevents the use of some of the solving strategies of Section 7.11.

Implementation – the real deal

Here is what Bochev and Gunzburger [103, Section 7.6.4] have to say about this: "[...] practical implementations cheat by substituting enforcement of the true zero mean constraint by using procedures collectively known as setting the pressure datum. These procedures essentially amount to removing one degree of freedom from the pressure space. Setting the pressure datum can be accomplished in many different ways, ranging from specifying a pressure value at a grid node to more complicated approaches in which a boundary traction is specified at a single node in lieu of the velocity condition; see [16, 24, 191] and the references cited therein for more details. Not surprisingly, in practice, the simplest approach of fixing the pressure value at a node also happens to be the most widely used in practice."

The idea is actually quite simple and requires two steps:

1. remove the null space by prescribing the pressure at one location and solve the system;

2. post-process the pressure so as to arrive at a pressure field which fulfils the required normalisation (surface, volume, ...)

The reason why it works is as follows: a constant pressure value lies in the null space, so that one can add or delete any value to the pressure field without consequence. As such I can choose said constant such that the pressure at a given node/element is zero. All other computed pressures are then relative to that one. The post-processing step will redistribute a constant value to all pressures (it will shift them up or down) so that the normalising condition is respected.



Relevant Literature

In <https://scicomp.stackexchange.com/questions/27645/pressure-boundary-condition-in-lid-> we find: *Zero mean pressure space is used for convenience when one is interested in FEA theory (basically, we cannot enforce $p(x_0) = p_0$ for $p \in L^2$ since it does not make sense); from the computational point of view, it is easier to fix one of the pressure DOFs (although you can subtract mean value at the post-processing step if you want to). When you are working w/ polynomial spaces—and this is exactly what you do in FEM—it is perfectly fine to enforce $p(x_0) = p_0$. Handle this constraint like you usually handle Dirichlet BCs (e.g., via modifying your matrix). It is also fine to ignore this constraint in some cases (e.g., Krylov solvers can do fine with this).*

<https://scicomp.stackexchange.com/questions/25134/mixed-finite-element-method-for-the-s>

7.14.2 Removing rotational nullspace

nullspace.tex

When free slip boundary conditions are prescribed in an annulus or hollow sphere geometry there exists a rotational nullspace, or in other words there exists a tangential velocity field ('pure rotation') which, if added or subtracted to the solution, generates a solution which is still the solution of the PDEs.

As in the pressure normalisation case (see section 7.14.1), the solution is simple:

1. fix the tangential velocity at *one* node on a boundary, and solve the system (the nullspace has been removed)³⁹
2. post-process the solution to have the velocity field fulfill the required conditions, i.e. either a zero net angular momentum or a zero net angular velocity of the domain.

Remark. In ASPECT this is available under the option "Remove nullspace = angular momentum" and "Remove nullspace = net rotation". The "angular momentum" option removes a rotation such that the net angular momentum is zero. The "net rotation" option removes the net rotation of the domain.

Angular momentum approach In physics, velocity is not a conserved quantity, but the momentum is. In order to remove the angular momentum, we search for a rotation vector $\vec{\omega}$ such that

$$\int_{\Omega} \rho [\vec{r} \times (\vec{v} - \vec{\omega} \times \vec{r})] dV = \vec{0} \quad (7.273)$$

The angular momentum of a rigid body can be obtained from the sum of the angular momentums

³⁹<https://scicomp.stackexchange.com/questions/3531/how-to-remove-rigid-body-motions-in-linear-elasticity>

of the particles forming the body⁴⁰:

$$\vec{H} = \sum_i \vec{L}_i \quad (7.274)$$

$$= \sum_i \vec{r}_i \times m_i \vec{v}_i \quad (7.275)$$

$$= \sum_i \vec{r}_i \times m_i (\vec{\omega} \times \vec{r}_i) \quad (7.276)$$

$$= \sum_i m_i \begin{pmatrix} \sum_i m_i (y_i^2 + z_i^2) & -\sum_i m_i x_i y_i & -\sum_i m_i x_i z_i \\ -\sum_i m_i x_i y_i & \sum_i m_i (x_i^2 + z_i^2) & -\sum_i m_i y_i z_i \\ -\sum_i m_i x_i z_i & -\sum_i m_i y_i z_i & \sum_i m_i (x_i^2 + y_i^2) \end{pmatrix} \cdot \begin{pmatrix} \omega_x \\ \omega_y \\ \omega_z \end{pmatrix} \quad (7.277)$$

In the continuum limit, we have:

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \vec{r} \times \vec{v} dV \quad (7.278)$$

and the 3×3 moment of inertia tensor \mathbf{I} (also called inertia tensor) is given by⁴¹

$$\mathbf{I} = \int_{\Omega} \rho(\vec{r}) [\vec{r} \cdot \vec{r} \mathbf{1} - \vec{r} \times \vec{r}] dV \quad (7.279)$$

so that the above equation writes: $\vec{H} = \mathbf{I} \cdot \vec{\omega}$ and then $\vec{\omega} = \mathbf{I}^{-1} \cdot \vec{H}$.

Ultimately, at each velocity node a rotation about the rotation vector $\vec{\omega}$ is then subtracted from the velocity solution [1412, eq. 26]:

$$\vec{v}_{new} = \vec{v}_{old} - \vec{\omega} \times \vec{r} \quad (7.280)$$

For the special case of a solid sphere of constant density and radius R the tensor \mathbf{I} becomes diagonal and we have

$$I = \frac{2}{5} m R^2$$

where m is the mass of the sphere.

The case of a hollow sphere is explained in Section 2.4.1 of Zhong *et al.* (2008) [1412].

Three dimensions

The angular momentum vector is given by:

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} yw - zv \\ zu - xw \\ xv - yu \end{pmatrix} d\vec{r} = \begin{pmatrix} \int_{\Omega} \rho(\vec{r}) (yw - zv) d\vec{r} \\ \int_{\Omega} \rho(\vec{r}) (zu - xw) d\vec{r} \\ \int_{\Omega} \rho(\vec{r}) (xv - yu) d\vec{r} \end{pmatrix} = \begin{pmatrix} H_x \\ H_y \\ H_z \end{pmatrix} \quad (7.281)$$

⁴⁰<http://www.kwon3d.com/theory/moi/iten.html>

⁴¹https://en.wikipedia.org/wiki/Moment_of_inertia

while the inertia tensor for a continuous body is given by

$$\mathbf{I} = \int_{\Omega} \rho(\vec{r}) [\vec{r} \cdot \vec{r} \mathbf{1} - \vec{r} \times \vec{r}] d\vec{r} \quad (7.282)$$

$$= \int_{\Omega} \rho(\vec{r}) \left[\begin{pmatrix} x^2 + y^2 + z^2 & 0 & 0 \\ 0 & x^2 + y^2 + z^2 & 0 \\ 0 & 0 & x^2 + y^2 + z^2 \end{pmatrix} - \begin{pmatrix} xx & xy & xz \\ yx & yy & yz \\ zx & zy & zz \end{pmatrix} \right] d\vec{r} \quad (7.283)$$

$$= \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} y^2 + z^2 & -xy & -xz \\ -yx & x^2 + z^2 & -yz \\ -zx & -zy & x^2 + y^2 \end{pmatrix} d\vec{r} \quad (7.284)$$

$$= \begin{pmatrix} \int_{\Omega} \rho(\vec{r})(y^2 + z^2) d\vec{r} & -\int_{\Omega} \rho(\vec{r})xy d\vec{r} & -\int_{\Omega} \rho(\vec{r})xz d\vec{r} \\ -\int_{\Omega} \rho(\vec{r})yx d\vec{r} & \int_{\Omega} \rho(\vec{r})(x^2 + z^2) d\vec{r} & -\int_{\Omega} \rho(\vec{r})yz d\vec{r} \\ -\int_{\Omega} \rho(\vec{r})zx d\vec{r} & -\int_{\Omega} \rho(\vec{r})zy d\vec{r} & \int_{\Omega} \rho(\vec{r})(x^2 + y^2) d\vec{r} \end{pmatrix} \quad (7.285)$$

$$= \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} \quad (7.286)$$

Two dimensions

In two dimensions, flow is taking place in the (x, y) plane. This means that $\vec{r} = (x, y, 0)$ and $\vec{v} = (u, v, 0)$ are coplanar, and therefore that $\vec{\omega}$ is perpendicular to the plane. We have then

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} 0 \\ 0 \\ xv - yu \end{pmatrix} d\vec{r} = \begin{pmatrix} 0 \\ 0 \\ \int_{\Omega} \rho(\vec{r})(xv - yu) d\vec{r} \end{pmatrix} \quad (7.287)$$

and

$$\mathbf{I} = \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} = \begin{pmatrix} I_{xx} & I_{xy} & 0 \\ I_{yx} & I_{yy} & 0 \\ 0 & 0 & I_{zz} \end{pmatrix} \quad (7.288)$$

since $I_{xz} = I_{yz} = 0$ as $z = 0$, and with $I_{xx} = \int_{\Omega} \rho(\vec{r})y^2 d\vec{r}$ and $I_{yy} = \int_{\Omega} \rho(\vec{r})x^2 d\vec{r}$. The solution to $\mathbf{I} \cdot \vec{\omega} = \vec{H}$ can be easily obtained (see Appendix D.0.2):

$$\omega_x = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} 0 & I_{xy} & 0 \\ 0 & I_{yy} & 0 \\ H_z & 0 & I_{zz} \end{vmatrix} = 0 \quad (7.289)$$

$$\omega_y = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} I_{xx} & 0 & 0 \\ I_{yx} & 0 & 0 \\ 0 & H_z & I_{zz} \end{vmatrix} = 0 \quad (7.290)$$

$$\omega_z = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} I_{xx} & I_{xy} & 0 \\ I_{yx} & I_{yy} & 0 \\ 0 & 0 & H_z \end{vmatrix} \quad (7.291)$$

$$= \frac{1}{\det(\mathbf{I})} (I_{xx}I_{yy}H_z - I_{yx}I_{xy}H_z) \quad (7.292)$$

$$= \frac{1}{\det(\mathbf{I})} (I_{xx}I_{yy} - I_{yx}I_{xy}) H_z \quad (7.293)$$

with $\det(\mathbf{I}) = I_{xx}I_{yy}I_{zz} - I_{yx}I_{xy}I_{zz} = (I_{xx}I_{yy} - I_{yx}I_{xy})I_{zz}$ and then

$$\omega_z = \frac{(I_{xx}I_{yy} - I_{yx}I_{xy})H_z}{(I_{xx}I_{yy} - I_{yx}I_{xy})I_{zz}} = \frac{H_z}{I_{zz}} = \frac{\int_{\Omega} \rho(\vec{r})(xv - yu)d\vec{r}}{\int_{\Omega} \rho(\vec{r})(x^2 + y^2)d\vec{r}}$$

Concretely, this means that in 2D one does not need to solve the system $\mathbf{I} \cdot \vec{\omega} = \vec{H}$ since only ω_z is not zero.

Then, since $\vec{r} = (x, y, z)$ and $\vec{\omega} = (0, 0, \omega_z)$:

$$\vec{v}_{new}(\vec{r}) = \vec{v}_{old} - \vec{\omega} \times \vec{r} = \begin{pmatrix} u_{old} - (-\omega_z y) \\ v_{old} - (\omega_z x) \\ 0 \end{pmatrix} \quad (7.294)$$

Chapter 8

The Discontinuous Galerkin Finite Element Method (DG-FEM)

chapter7.tex

dgintro.tex

What is DG?

- it is a variant of the SG ("Standard Galerkin FEM")¹
- SG-FEM requires continuity of the solution along element interfaces (edges).
- DG-FEM does not require continuity of the solution along edges.
- DG methods have more degrees of freedom than SG methods.
- DG-FEM shares some properties with FVM

Various books about DG-FEM

- *Discontinuous Galerkin Methods. Theory, Computation and Applications* by Cockburn, Karniadakis and Shu [267]
- *Mathematical Aspects of Discontinuous Galerkin Methods* by Di Pietro and Ern [999]
- *Discontinuous Galerkin Methods. Analysis and Applications to Compressible Flow* by Dolejsi and Feistauer [337]
- *Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations* by Rivi re [1075]
- *Discontinuous finite elements in fluid dynamics and heat transfer* by Li [779]
- *Nodal Discontinuous Galerkin Methods. Algorithms, Analysis, and Applications* by Hesthaven & Warbuton [568]

¹Some authors use the acronym CG for Continuous Galerkin but since the Conjugate Gradient solver acronym CG is very much present in FE codes it can be confusing so we use here SG instead.

DG flavors There are many different flavours of the Discontinuous Galerkin Finite Element Method:

- **HDG**: Hybridizable DG [263, 268, 937, 938, 936]
- **IPG**: Interior Penalty G [893, 894]
- **IIPG**: Incomplete Interior Penalty G [338]
- **SIPG**: Symmetric Interior Penalty G [118, 1138]
- **LDG**: Local DG [215, 265, 211, 264]

Pro and cons for DG-FEM versus SG-FEM

- Assembly of stiffness matrix is easier to implement (ref?).
- Refinement of triangles is easier to implement (ref?).
- DG methods can easily handle adaptivity strategies since refinement or unrefinement of the grid can be achieved without taking into account the continuity restrictions typical of conforming finite element methods. Moreover, the degree of the approximating polynomial can be easily changed from one element to the other. [266]
- DG methods can support high order local approximations that can vary nonuniformly over the mesh.
- DG methods are readily parallelizable. Since the elements are discontinuous, the mass matrix is block diagonal and since the size of the blocks is equal to the number of degrees of freedom inside the corresponding elements, the blocks can be inverted by hand once and for all.[266]

The DG-FEM in geodynamics This method has not been used extensively in geodynamics with (so-far) two noticeable exceptions:

- Lehmann *et al.* , Comparison of continuous and discontinuous Galerkin approaches for variable-viscosity Stokes flow (2015) [761]
- He *et al.* , A discontinuous Galerkin method with a bound preserving limiter for the advection of non-diffusive fields in solid Earth geodynamics (2017) [555]
- Puckett *et al.* , New numerical approaches for modeling thermochemical convection in a compositionally stratified fluid (2018) [1020]

8.1 First-order advection ODE in 1D

dgfem1D.tex

What follows is borrowed from the book *Discontinuous finite elements in fluid dynamics and heat transfer* by Ben Q. Li [779].

To illustrate the basic ideas of the discontinuous finite element method, we consider a simple, one-dimensional, first order differential equation with u specified at one of the boundaries:

$$\frac{du}{dx} + g = 0 \quad x \in [a, b] \quad \text{and} \quad u(x = a) = u_a \quad (8.1)$$

where g is a constant (for simplicity). The domain is discretized such that : $\Omega_j = [x_j, x_{j+1}]$ with $j = 1, 2, \dots, nel$. Then, integrating the above equation over the element j with respect to a weighting function $f(x)$

$$\int_{x_j}^{x_{j+1}} \left(\frac{du}{dx} + g \right) f(x) dx = 0 \quad (8.2)$$

Remembering that $\int_c^d u(x)v'(x)dx = [u(x)v(x)]_c^d - \int_c^d u'(x)v(x)dx$, we can now perform an integration by parts on the differential operator and we obtain:

$$[u(x)f(x)]_{x_j}^{x_{j+1}} - \int_{x_j}^{x_{j+1}} \left(u \frac{df}{dx} - gf(x) \right) dx = 0 \quad (8.3)$$

or,

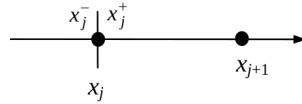
$$u(x_{j+1})f(x_{j+1}) - u(x_j)f(x_j) - \int_{x_j}^{x_{j+1}} \left(u \frac{df}{dx} - gf(x) \right) dx = 0 \quad (8.4)$$

On the element Ω_j the function u is approximated by $u_h \in H$, H being an appropriate function space of finite dimension, and f by f_h taken from the same function space as u_h . Upon substituting (u_h, f_h) for (u, f) in the equation above, we have the discontinuous Galerkin finite element formulation:

$$u_h(x_{j+1})f_h(x_{j+1}) - u_h(x_j)f_h(x_j) - \int_{x_j}^{x_{j+1}} \left(u_h \frac{df_h}{dx} - gf_h(x) \right) dx = 0 \quad (8.5)$$

In the continuous finite element approach, the field variable u_h is forced to be continuous across the boundary. The essential idea for the discontinuous method is that u_h is allowed to be discontinuous across the boundary. Therefore, across the element, the following two different values are defined at the two sides of the boundary:

$$u_j^+ = \lim_{x \searrow x_j^+} u_h(x) \quad u_j^- = \lim_{x \nearrow x_j^-} u_h(x) \quad (8.6)$$



An illustration of the jump across x_j of element j : x_j and x_{j+1} mark the boundaries of the element

Conversely, we also have:

$$u_{j+1}^+ = \lim_{x \searrow x_{j+1}^+} u_h(x) \quad u_{j+1}^- = \lim_{x \nearrow x_{j+1}^-} u_h(x) \quad (8.7)$$

It is key to remember that 1) u_h is discontinuous only at the element boundaries; 2) the solution u is smooth within (but excluding) the boundary. By this definition, the above equation contains the variables only within the integral limits of Ω_j . As a consequence, there is no direct coupling with other intervals or other elements. *The field values at a node, or the interface between two elements, are not unique.* They are calculated using the two limiting values approaching the interface from the two adjacent elements. This feature is certainly desirable for problems with internal discontinuities.

We can finally write for a single element:

$$u_{j+1}^- f_h(x_{j+1}) - u_j^+ f_h(x_j) - \int_{x_j}^{x_{j+1}} \left(u_h \frac{df_h}{dx} - gf_h(x) \right) dx = 0 \quad (8.8)$$

8.2 Steady state diffusion in 1D

dgfem1D_ssdiff.tex

Let us start simple with the 1D steady state heat conduction problem in 1D, given by the following equation:

$$\frac{d^2 T}{dx^2} = 0 \quad T(x=0) = 0 \quad T(x=1) = 1 \quad \text{on } x \in [0, 1] \quad (8.9)$$

Although this equation is usually solved as is with its second-order derivative, it can also be written in a mixed form, using the heat flux q (a scalar in 1D):

$$\begin{aligned} -\frac{dq}{dx} &= 0 \\ q - \frac{dT}{dx} &= 0 \quad x \in [0, 1] \end{aligned} \quad (8.10)$$

and the boundary conditions remain unchanged.

We apply the standard approach to establish the weak forms of these two first-order ODEs, and we do so on an element e bound by nodes k and $k+1$ with coordinates x_k and x_{k+1}

$$-\int_{x_k}^{x_{k+1}} \frac{dq}{dx} \tilde{f}(x) dx = -[q\tilde{f}]_{x_k}^{x_{k+1}} + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} q(x) dx = 0 \quad (8.11)$$

$$\int_{x_k}^{x_{k+1}} \left(q - \frac{dT}{dx} \right) \bar{f}(x) dx = \int_{x_k}^{x_{k+1}} q(x) \bar{f}(x) dx - [T\bar{f}]_{x_k}^{x_{k+1}} + \int_{x_k}^{x_{k+1}} \frac{d\bar{f}}{dx} T(x) dx = 0 \quad (8.12)$$

where \tilde{f} and \bar{f} are test functions. We now must examine the term between square brackets. Inside the element, the test functions \tilde{f} and \bar{f} are well defined polynomials and we coin:

$$\tilde{f}_k^+ = \tilde{f}(x_k^+) \quad (8.13)$$

$$\tilde{f}_{k+1}^- = \tilde{f}(x_{k+1}^-) \quad (8.14)$$

$$\bar{f}_k^+ = \bar{f}(x_k^+) \quad (8.15)$$

$$\bar{f}_{k+1}^- = \bar{f}(x_{k+1}^-) \quad (8.16)$$

Concerning q and T , we will for now give them values \hat{q}_k and \hat{T}_k at node k and \hat{q}_{k+1} and \hat{T}_{k+1} at node $k+1$, and we will specify the hat quantities as follows:

$$\begin{aligned} \hat{T}_k &= \begin{cases} T_k^- & k = 1 \\ \frac{1}{2}(T_k^- + T_k^+) + \mathcal{C}(T_k^- - T_k^+) & k = 2, \dots, N-1 \\ T_k^+ & k = N \end{cases} \\ \hat{q}_k &= \begin{cases} q_k^+ - \mathcal{E}(T_k^- - T_k^+) & k = 1 \\ \frac{1}{2}(q_k^+ + q_k^-) - \mathcal{E}(T_k^- - T_k^+) - \mathcal{C}(q_k^- - q_k^+) & k = 2, \dots, N-1 \\ q_k^- - \mathcal{E}(T_k^- - T_k^+) & k = N \end{cases} \end{aligned} \quad (8.17)$$

where N is the number of nodes and where \mathcal{C} and \mathcal{E} are two constants.

Discuss the meaning/values of these!

Remark. Note that $\hat{T}_k = T_1^-$ on the left boundary is consistent with $\hat{T}_k = \frac{1}{2}(T_k^- + T_k^+) + \mathcal{C}(T_k^- - T_k^+)$ provided $T_1^- = T_1^+$. The same goes for the right boundary, and the same reasoning applies for the heat flux terms \hat{q}_k .

Inside an element bounded by nodes k and $k+1$, the temperature T and heat flux q are interpolated over an isoparametric linear element:

$$\begin{aligned} T_h(x) &= \mathcal{N}_k(x)T_k^+ + \mathcal{N}_{k+1}(x)T_{k+1}^- \\ q_h(x) &= \mathcal{N}_k(x)q_k^+ + \mathcal{N}_{k+1}(x)q_{k+1}^- \end{aligned}$$

As in the (Continuous/Standard) Galerkin case of section 6.1, the test functions are taken to be the basis functions, and in this case for both temperature and flux.

There are four unknowns q_k^+ , q_{k+1}^- , T_k^+ and T_{k+1}^- per element. All other q and T quantities in the above/following equations will need to find their way to the rhs.

- Eq. 8.11 becomes:

$$\begin{aligned} 0 &= -\hat{q}_{k+1}\tilde{f}(x_{k+1}^-) + \hat{q}_k\tilde{f}(x_k^+) + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} q_h(x) dx \\ &= -\hat{q}_{k+1}\tilde{f}_{k+1}^- + \hat{q}_k\tilde{f}_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} (\mathcal{N}_k(x)q_k^+ + \mathcal{N}_{k+1}(x)q_{k+1}^-) dx \\ &= -\hat{q}_{k+1}\tilde{f}_{k+1}^- + \hat{q}_k\tilde{f}_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} \mathcal{N}_k(x) dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} \mathcal{N}_{k+1}(x) dx \cdot q_{k+1}^- \quad (8.18) \end{aligned}$$

- We take $\tilde{f} = \mathcal{N}_k$ and by virtue of the properties of basis functions \mathcal{N} we have:

$$\begin{aligned} \tilde{f}_k^+ &= \tilde{f}(x_k^+) = \mathcal{N}_k(x_k^+) = 1 \\ \tilde{f}_{k+1}^- &= \tilde{f}(x_{k+1}^-) = \mathcal{N}_k(x_{k+1}^-) = 0 \end{aligned}$$

so that

$$\begin{aligned} 0 &= \hat{q}_k + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k(x) dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1}(x) dx \cdot q_{k+1}^- \\ &= \frac{1}{2}(q_k^+ + q_k^-) - \mathcal{E}(T_k^- - T_k^+) - \mathcal{C}(q_k^- - q_k^+) \\ &\quad + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1} dx \cdot q_{k+1}^- \quad (8.19) \end{aligned}$$

- We take $\tilde{f} = \mathcal{N}_{k+1}$ and likewise:

$$\begin{aligned} \tilde{f}_k^+ &= \tilde{f}(x_k^+) = \mathcal{N}_{k+1}(x_k^+) = 0 \\ \tilde{f}_{k+1}^- &= \tilde{f}(x_{k+1}^-) = \mathcal{N}_{k+1}(x_{k+1}^-) = 1 \end{aligned}$$

so that

$$\begin{aligned} 0 &= -\hat{q}_{k+1} + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k(x) dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1}(x) dx \cdot q_{k+1}^- \\ &= -\left[\frac{1}{2}(q_{k+1}^+ + q_{k+1}^-) - \mathcal{E}(T_{k+1}^- - T_{k+1}^+) - \mathcal{C}(q_{k+1}^- - q_{k+1}^+) \right] \\ &\quad + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1} dx \cdot q_{k+1}^- \quad (8.20) \end{aligned}$$

and finally

$$\begin{aligned} &\int_{x_k}^{x_{k+1}} \left(\frac{\frac{d\mathcal{N}_k}{dx} \mathcal{N}_k}{\frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k} \quad \frac{\frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1}}{\frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1}} \right) dx \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\mathcal{C} + \frac{1}{2})q_k^+ \\ (\mathcal{C} - \frac{1}{2})q_{k+1}^- \end{pmatrix} + \begin{pmatrix} \mathcal{E}T_k^+ \\ \mathcal{E}T_{k+1}^- \end{pmatrix} \\ &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2})q_k^- \\ (\mathcal{C} + \frac{1}{2})q_{k+1}^+ \end{pmatrix} + \begin{pmatrix} \mathcal{E}T_k^- \\ \mathcal{E}T_{k+1}^+ \end{pmatrix} \quad (8.21) \end{aligned}$$

- Eq. 8.12 becomes:

$$\begin{aligned}
0 &= -[T\bar{f}]_{x_k}^{x_{k+1}} + \int_{x_k}^{x_{k+1}} q_h(x)\bar{f}(x)dx + \int_{x_k}^{x_{k+1}} \frac{d\bar{f}}{dx}T_h(x)dx \\
&= -\hat{T}_{k+1}\bar{f}_{k+1}^- + \hat{T}_k\bar{f}_k^+ + \int_{x_k}^{x_{k+1}} q_h(x)\bar{f}(x)dx + \int_{x_k}^{x_{k+1}} \frac{d\bar{f}}{dx}T_h(x)dx
\end{aligned}$$

- We take $\bar{f} = \mathcal{N}_k$:

$$\begin{aligned}
\bar{f}_k^+ &= \bar{f}(x_k^+) = \mathcal{N}_k(x_k^+) = 1 \\
\bar{f}_{k+1}^- &= \bar{f}(x_{k+1}^-) = \mathcal{N}_k(x_{k+1}^-) = 0
\end{aligned}$$

so that

$$\begin{aligned}
0 &= \hat{T}_k + \int_{x_k}^{x_{k+1}} q_h(x)\mathcal{N}_k dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx}T_h(x)dx \\
&= \frac{1}{2}(T_k^- + \textcolor{red}{T}_k^+) + \mathcal{C}(T_k^- - \textcolor{red}{T}_k^+) \\
&+ \int_{x_k}^{x_{k+1}} (\mathcal{N}_k(x)\textcolor{red}{q}_k^+ + \mathcal{N}_{k+1}(x)\textcolor{red}{q}_{k+1}^-)\mathcal{N}_k dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx}(\mathcal{N}_k(x)\textcolor{red}{T}_k^+ + \mathcal{N}_{k+1}(x)\textcolor{red}{T}_{k+1}^-)dx
\end{aligned}$$

- We take $\bar{f} = \mathcal{N}_{k+1}$:

$$\begin{aligned}
\bar{f}_k^+ &= \bar{f}(x_k^+) = \mathcal{N}_{k+1}(x_k^+) = 0 \\
\bar{f}_{k+1}^- &= \bar{f}(x_{k+1}^-) = \mathcal{N}_{k+1}(x_{k+1}^-) = 1
\end{aligned}$$

so that

$$\begin{aligned}
0 &= -\hat{T}_{k+1} + \int_{x_k}^{x_{k+1}} q_h(x)\mathcal{N}_{k+1}dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx}T_h(x)dx \\
&= -\left[\frac{1}{2}(\textcolor{red}{T}_{k+1}^- + T_{k+1}^+) + \mathcal{C}(\textcolor{red}{T}_{k+1}^- - T_{k+1}^+)\right] \\
&+ \int_{x_k}^{x_{k+1}} (\mathcal{N}_k(x)\textcolor{red}{q}_k^+ + \mathcal{N}_{k+1}(x)\textcolor{red}{q}_{k+1}^-)\mathcal{N}_{k+1}dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx}(\mathcal{N}_k(x)\textcolor{red}{T}_k^+ + \mathcal{N}_{k+1}(x)\textcolor{red}{T}_{k+1}^-)dx
\end{aligned} \tag{8.23}$$

and finally

$$\begin{aligned}
&\int_{x_k}^{x_{k+1}} \begin{pmatrix} \mathcal{N}_k\mathcal{N}_k & \mathcal{N}_k\mathcal{N}_{k+1} \\ \mathcal{N}_{k+1}\mathcal{N}_k & \mathcal{N}_{k+1}\mathcal{N}_{k+1} \end{pmatrix} dx \begin{pmatrix} \textcolor{red}{q}_k^+ \\ \textcolor{red}{q}_{k+1}^- \end{pmatrix} + \int_{x_k}^{x_{k+1}} \begin{pmatrix} \frac{d\mathcal{N}_k}{dx}\mathcal{N}_k & \frac{d\mathcal{N}_k}{dx}\mathcal{N}_{k+1} \\ \frac{d\mathcal{N}_{k+1}}{dx}\mathcal{N}_k & \frac{d\mathcal{N}_{k+1}}{dx}\mathcal{N}_{k+1} \end{pmatrix} dx \begin{pmatrix} \textcolor{red}{T}_k^+ \\ \textcolor{red}{T}_{k+1}^- \end{pmatrix} \\
&+ \begin{pmatrix} (\frac{1}{2} - \mathcal{C})\textcolor{red}{T}_k^+ \\ -(\mathcal{C} + \frac{1}{2})\textcolor{red}{T}_{k+1}^- \end{pmatrix} = \begin{pmatrix} -(\mathcal{C} + \frac{1}{2})T_k^- \\ (\frac{1}{2} - \mathcal{C})T_{k+1}^+ \end{pmatrix}
\end{aligned} \tag{8.25}$$

We will also use the results obtained in Appendix E for 1D linear elements:

$$\begin{aligned} \mathbf{M}^e &= \int_{\Omega_e} \vec{\mathcal{N}}^T \vec{\mathcal{N}} dV = \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_k \mathcal{N}_k & \mathcal{N}_k \mathcal{N}_{k+1} \\ \mathcal{N}_{k+1} \mathcal{N}_k & \mathcal{N}_{k+1} \mathcal{N}_{k+1} \end{pmatrix} dV = \frac{h}{2} \frac{1}{3} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = \frac{h}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \\ \mathbf{K}^e &= \int_{\Omega_e} \begin{pmatrix} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k & \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1} \\ \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k & \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1} \end{pmatrix} dV = \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix} \end{aligned} \quad (8.26)$$

Filling this into equations (8.21) and (8.25), gives

$$\begin{aligned} \mathbf{K}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\mathcal{C} + \frac{1}{2}) q_k^+ \\ (\mathcal{C} - \frac{1}{2}) q_{k+1}^- \end{pmatrix} + \begin{pmatrix} \mathcal{E} T_k^+ \\ \mathcal{E} T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2}) q_k^- \\ (\mathcal{C} + \frac{1}{2}) q_{k+1}^+ \end{pmatrix} + \begin{pmatrix} \mathcal{E} T_k^- \\ \mathcal{E} T_{k+1}^+ \end{pmatrix} \\ \mathbf{M}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \mathbf{K}^e \cdot \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\frac{1}{2} - \mathcal{C}) T_k^+ \\ -(\mathcal{C} + \frac{1}{2}) T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\mathcal{C} + \frac{1}{2}) T_k^- \\ (\frac{1}{2} - \mathcal{C}) T_{k+1}^+ \end{pmatrix} \end{aligned} \quad (8.27)$$

which becomes

$$\begin{aligned} \begin{pmatrix} \mathcal{C} & -\frac{1}{2} \\ \frac{1}{2} & \mathcal{C} \end{pmatrix} \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} \mathcal{E} & 0 \\ 0 & \mathcal{E} \end{pmatrix} \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2}) q_k^- + \mathcal{E} T_k^- \\ (\mathcal{C} + \frac{1}{2}) q_{k+1}^+ + \mathcal{E} T_{k+1}^+ \end{pmatrix} \\ \frac{h}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} -\mathcal{C} & -\frac{1}{2} \\ \frac{1}{2} & -\mathcal{C} \end{pmatrix} \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\frac{1}{2} + \mathcal{C}) T_k^- \\ (\frac{1}{2} - \mathcal{C}) T_{k+1}^+ \end{pmatrix} \end{aligned}$$

Combining these equations gives the expression for the linear system for the element under consideration:

$$\left(\begin{pmatrix} \frac{h}{3} & \frac{h}{6} & -\mathcal{C} & -\frac{1}{2} \\ \frac{h}{6} & \frac{h}{3} & \frac{1}{2} & -\mathcal{C} \\ \mathcal{C} & -\frac{1}{2} & \mathcal{E} & 0 \\ \frac{1}{2} & \mathcal{C} & 0 & \mathcal{E} \end{pmatrix} \begin{pmatrix} q_k^+ \\ q_{k+1}^- \\ T_k^+ \\ T_{k+1}^- \end{pmatrix} = \begin{pmatrix} -(\frac{1}{2} + \mathcal{C}) T_k^- \\ (\frac{1}{2} - \mathcal{C}) T_{k+1}^+ \\ -(\frac{1}{2} - \mathcal{C}) q_k^- + \mathcal{E} T_k^- \\ (\frac{1}{2} + \mathcal{C}) q_{k+1}^+ + \mathcal{E} T_{k+1}^+ \end{pmatrix} \right) \quad (8.28)$$

Left boundary Special care must be taken for the two elements on the boundaries of the domain. On the left, we have

$$\begin{aligned} \hat{q}_1 &= q_1^+ - \mathcal{E}(T_1^- - T_1^+) \\ \hat{T}_1 &= T_1^- \end{aligned}$$

Eq. (8.19) becomes:

$$\begin{aligned} 0 &= \hat{q}_k + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k(x) dx \cdot q_k^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1}(x) dx \cdot q_{k+1}^- \\ &= q_1^+ - \mathcal{E}(T_k^- - T_1^+) + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k dx \cdot q_1^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1} dx \cdot q_2^- \end{aligned} \quad (8.29)$$

Eq. (8.22) becomes:

$$\begin{aligned} 0 &= \hat{T}_1 + \int_{x_k}^{x_{k+1}} q_h(x) \mathcal{N}_k dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} T_h(x) dx \\ &= T_1^- + \int_{x_k}^{x_{k+1}} (\mathcal{N}_k(x) q_1^+ + \mathcal{N}_{k+1}(x) q_2^-) \mathcal{N}_k dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_k}{dx} (\mathcal{N}_k(x) T_1^+ + \mathcal{N}_{k+1}(x) T_2^-) dx \end{aligned} \quad (8.30)$$

Eq. (8.28) then becomes:

$$\left(\begin{pmatrix} \frac{h}{3} & \frac{h}{6} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{h}{6} & \frac{h}{3} & \frac{1}{2} & -\mathcal{C} \\ \frac{1}{2} & -\frac{1}{2} & \mathcal{E} & 0 \\ \frac{1}{2} & \mathcal{C} & 0 & \mathcal{E} \end{pmatrix} \begin{pmatrix} q_1^+ \\ q_2^- \\ T_1^+ \\ T_2^- \end{pmatrix} = \begin{pmatrix} -T_1^- \\ (\frac{1}{2} - \mathcal{C}) T_{k+1}^+ \\ \mathcal{E} T_k^- \\ (\frac{1}{2} + \mathcal{C}) q_{k+1}^+ + \mathcal{E} T_{k+1}^+ \end{pmatrix} \right) \quad (8.31)$$

Right boundary The element is composed of nodes $N - 1$ and N . The fluxes are

$$\begin{aligned}\hat{q}_N &= q_N^+ - \mathcal{E}(T_N^- - T_N^+) \\ \hat{T}_N &= T_N^-\end{aligned}$$

Eq. (8.20) becomes:

$$\begin{aligned}0 &= -\hat{q}_N + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k(x) dx \cdot q_{N-1}^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1}(x) dx \cdot q_N^- \\ &= -[q_N^+ - \mathcal{E}(T_N^- - T_N^+)] + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k dx \cdot \textcolor{red}{q}_{N-1}^+ + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1} dx \cdot \textcolor{red}{q}_N^- \quad (8.32)\end{aligned}$$

Eq. (8.24) becomes:

$$\begin{aligned}0 &= -\hat{T}_N + \int_{x_k}^{x_{k+1}} q_h(x) \mathcal{N}_{k+1} dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} T_h(x) dx \\ &= -T_N^- + \int_{x_k}^{x_{k+1}} (\mathcal{N}_k(x) \textcolor{red}{q}_{N-1}^+ + \mathcal{N}_{k+1}(x) \textcolor{red}{q}_N^-) \mathcal{N}_{k+1} dx + \int_{x_k}^{x_{k+1}} \frac{d\mathcal{N}_{k+1}}{dx} (\mathcal{N}_k(x) \textcolor{red}{T}_{N-1}^+ + \mathcal{N}_{k+1}(x) \textcolor{red}{T}_N^-) dx\end{aligned}$$

Eq. (8.28) then becomes:

$$\begin{pmatrix} \frac{h}{3} & \frac{h}{6} & -\mathcal{C} & -\frac{1}{2} \\ \frac{h}{6} & \frac{h}{3} & \frac{1}{2} & \frac{1}{2} \\ \mathcal{C} & -\frac{1}{2} & \mathcal{E} & 0 \\ \frac{1}{2} & -\frac{1}{2} & 0 & \mathcal{E} \end{pmatrix} \begin{pmatrix} \textcolor{red}{q}_{N-1}^+ \\ \textcolor{red}{q}_N^- \\ \textcolor{red}{T}_{N-1}^+ \\ \textcolor{red}{T}_N^- \end{pmatrix} = \begin{pmatrix} -(\frac{1}{2} + \mathcal{C})T_{N-1}^- \\ T_N^+ \\ -(\frac{1}{2} - \mathcal{C})q_{N-1}^- + \mathcal{E}T_{N-1}^- \\ \mathcal{E}T_N^+ \end{pmatrix} \quad (8.33)$$

Solving strategies Following Li [779], there are three main strategies:

- Successive substitution: all the variables are initialized to zero. Eq. (8.31) is solved to obtain the data for the first element, where boundary conditions are specified. Then Eq. (8.28) is used for all interior elements. Finally Eq. (8.33) is used for the last element. This procedure is carried out until all fields have converged.
- Global assembly: this approach is identical to the one we have taken so far with the continuous Galerkin Finite Element method. We form a large global matrix and then solve the linear system to obtain the solution. The disadvantage of this approach lies in the size of the generated matrix: each node counts 4 dofs so the assembled matrix will be about 4 times as big as in the standard FE case.
- Elimination then assembly: one can first eliminate the variable q and solve for the temperature T only. This speeds up the calculations, but also increases the bandwidth of the element matrix. Li [779] states: "Further comparison shows that the saving in CPU time for solving T alone is less significant than the $q - T$ iterative solution, in particular, for 3-D problems."

Eq. (8.28) can be rewritten:

$$\begin{pmatrix} \mathbf{M}_e & \mathbf{C}_1 \\ \mathbf{C}_2 & \mathbf{E} \end{pmatrix} \cdot \begin{pmatrix} \vec{q} \\ \vec{T} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{g} \end{pmatrix} \quad (8.34)$$

The unknown of the original ODE is temperature so this is the quantity we are after. The first line of the matrix can be written

$$\mathbf{M}_e \cdot \vec{q} + \mathbf{C}_1 \cdot \vec{T} = \vec{f}$$

or,

$$\vec{q} = \mathbf{M}_e^{-1} \cdot (\vec{f} - \mathbf{C}_1 \cdot \vec{T})$$

The second line of the matrix is

$$\mathbf{C}_2 \cdot \vec{q} + \mathbf{E} \cdot \vec{T} = \vec{g}$$

and we then replace \vec{q} by the expression above:

$$\mathbf{C}_2 \cdot [\mathbf{M}_e^{-1} \cdot (\vec{f} - \mathbf{C}_1 \cdot \vec{T})] + \mathbf{E} \cdot \vec{T} = \vec{g}$$

or,

$$-\mathbf{C}_2 \cdot \mathbf{M}_e^{-1} \cdot \mathbf{C}_1 \cdot \vec{T} + \mathbf{E} \cdot \vec{T} = \vec{g} - \mathbf{C}_2 \cdot \mathbf{M}_e^{-1} \cdot \vec{f}$$

and finally

$$[\mathbf{E} - \mathbf{C}_2 \cdot \mathbf{M}_e^{-1} \cdot \mathbf{C}_1] \cdot \vec{T} = \vec{g} - \mathbf{C}_2 \cdot \mathbf{M}_e^{-1} \cdot \vec{f}$$

Note that the matrix will still be twice as big than in the standard FEM case since each node counts two temperature dofs.

Choice of \mathcal{C} and \mathcal{E} Li [779] shows satisfying results for $\mathcal{E} = 4$ and $\mathcal{C} = -1/2, 0, 1/2$ or $\mathcal{E} = 0$ and $\mathcal{C} = 1, 4, 10$.

Remark. *Aside from the sheer complexity of the above derivations as compared to those for the SG method, the fact that we have two tuning parameters \mathcal{E} and \mathcal{C} is a real drawback.*

8.3 Time-dependent diffusion PDE in 1D

dgfem1D-diff.tex

Starting from the simple transient 1-D heat conduction problem similar to the steady state heat conduction problem only with added time dependence:

$$\frac{\partial T}{\partial t} = \frac{\partial^2 T}{\partial x^2} \quad T(x=0) = 0 \quad T(x=1) = 1 \quad \text{on } x \in [0, 1] \quad (8.35)$$

Once again we split this system into two separate first order equations:

$$\begin{aligned} \frac{\partial T}{\partial t} - \frac{\partial q}{\partial x} &= 0 \\ \frac{\partial T}{\partial x} - q &= 0 \end{aligned} \quad (8.36)$$

We apply the standard approach to establish the weak forms of these two first-order PDEs, and we do so on an element e bound by nodes k and $k+1$ with coordinates x_k and x_{k+1}

$$-\int_{x_k}^{x_{k+1}} \left(\frac{\partial T}{\partial t} - \frac{\partial q}{\partial x} \right) \tilde{f}(x) dx = \int_{x_k}^{x_{k+1}} \frac{\partial T}{\partial t} \tilde{f}(x) dx - [q \tilde{f}]_{x_k}^{x_{k+1}} + \int_{x_k}^{x_{k+1}} \frac{d\tilde{f}}{dx} q(x) dx = 0 \quad (8.37)$$

$$\int_{x_k}^{x_{k+1}} \left(q - \frac{\partial T}{\partial x} \right) \bar{f}(x) dx = \int_{x_k}^{x_{k+1}} q(x) \bar{f}(x) dx - [T \bar{f}]_{x_k}^{x_{k+1}} + \int_{x_k}^{x_{k+1}} \frac{d\bar{f}}{dx} T(x) dx = 0 \quad (8.38)$$

where \tilde{f} and \bar{f} are test functions.

In what follows we coin $\dot{T} = \partial T / \partial t$ (for convenience of notation). We once again recover Equations (8.21) and (8.25), although with an additional time derivative term.

Filling this into equations (8.21) and (8.25), gives

$$\begin{aligned} \mathbf{K}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \mathbf{M}^e \cdot \begin{pmatrix} \dot{T}_k^+ \\ \dot{T}_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\mathcal{C} + \frac{1}{2})q_k^+ \\ (\mathcal{C} - \frac{1}{2})q_{k+1}^- \end{pmatrix} + \begin{pmatrix} \mathcal{E}T_k^+ \\ \mathcal{E}T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2})q_k^- \\ (\mathcal{C} + \frac{1}{2})q_{k+1}^+ \end{pmatrix} + \begin{pmatrix} \mathcal{E}T_k^- \\ \mathcal{E}T_{k+1}^+ \end{pmatrix} \\ \mathbf{M}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \mathbf{K}^e \cdot \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\frac{1}{2} - \mathcal{C})T_k^+ \\ -(\mathcal{C} + \frac{1}{2})T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\mathcal{C} + \frac{1}{2})T_k^- \\ (\frac{1}{2} - \mathcal{C})T_{k+1}^+ \end{pmatrix} \end{aligned} \quad (8.39)$$

In what follows we set $\mathcal{E} = 0$ so that we have

$$\begin{aligned} \mathbf{K}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \mathbf{M}^e \cdot \begin{pmatrix} \dot{T}_k^+ \\ \dot{T}_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\mathcal{C} + \frac{1}{2})q_k^+ \\ (\mathcal{C} - \frac{1}{2})q_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2})q_k^- \\ (\mathcal{C} + \frac{1}{2})q_{k+1}^+ \end{pmatrix} \\ \mathbf{M}^e \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \mathbf{K}^e \cdot \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\frac{1}{2} - \mathcal{C})T_k^+ \\ -(\mathcal{C} + \frac{1}{2})T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\mathcal{C} + \frac{1}{2})T_k^- \\ (\frac{1}{2} - \mathcal{C})T_{k+1}^+ \end{pmatrix} \end{aligned} \quad (8.40)$$

Using the expressions for \mathbf{M}^e and \mathbf{K}^e obtained in Appendix E for 1D linear elements we arrive at

$$\begin{aligned} \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \frac{h}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \cdot \begin{pmatrix} \dot{T}_k^+ \\ \dot{T}_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\mathcal{C} + \frac{1}{2})q_k^+ \\ (\mathcal{C} - \frac{1}{2})q_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2})q_k^- \\ (\mathcal{C} + \frac{1}{2})q_{k+1}^+ \end{pmatrix} \\ \frac{h}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} + \begin{pmatrix} (\frac{1}{2} - \mathcal{C})T_k^+ \\ -(\mathcal{C} + \frac{1}{2})T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\mathcal{C} + \frac{1}{2})T_k^- \\ (\frac{1}{2} - \mathcal{C})T_{k+1}^+ \end{pmatrix} \end{aligned} \quad (8.41)$$

which simplifies to

$$\begin{aligned} \begin{pmatrix} C & -1/2 \\ 1/2 & C \end{pmatrix} \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} h/3 & h/6 \\ h/6 & h/3 \end{pmatrix} \cdot \begin{pmatrix} \dot{T}_k^+ \\ \dot{T}_{k+1}^- \end{pmatrix} &= \begin{pmatrix} (\mathcal{C} - \frac{1}{2})q_k^- \\ (\mathcal{C} + \frac{1}{2})q_{k+1}^+ \end{pmatrix} \\ \begin{pmatrix} h/3 & h/6 \\ h/6 & h/3 \end{pmatrix} \cdot \begin{pmatrix} q_k^+ \\ q_{k+1}^- \end{pmatrix} + \begin{pmatrix} -C & -1/2 \\ 1/2 & -C \end{pmatrix} \cdot \begin{pmatrix} T_k^+ \\ T_{k+1}^- \end{pmatrix} &= \begin{pmatrix} -(\mathcal{C} + \frac{1}{2})T_k^- \\ (\frac{1}{2} - \mathcal{C})T_{k+1}^+ \end{pmatrix} \end{aligned} \quad (8.42)$$

or,

$$\begin{aligned} \mathbf{C}_1 \vec{q} + \mathbf{M} \vec{T} &= \vec{f} \\ \mathbf{M} \vec{q} + \mathbf{C}_2 \vec{T} &= \vec{g} \end{aligned}$$

so

$$\vec{q} = \mathbf{M}^{-1}(\vec{g} - \mathbf{C}_2 \vec{T})$$

and then

$$\mathbf{C}_1[\mathbf{M}^{-1}(\vec{g} - \mathbf{C}_2 \vec{T})] + \mathbf{M} \vec{T} = \vec{f}$$

NOT REALLY FINISHED...

8.4 Time-dependent advection PDE in 1D

Starting from the 1-D advection equation:

$$\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} = 0 \quad (8.43)$$

where T is the temperature and u the velocity. As shown before we start by discretizing the domain into a collection of elements. Then the above equation can be integrated over the element which is bounded by nodes x_k and x_{k+1} .

$$\int_{x_k}^{x_{k+1}} \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} \right) \tilde{f}(x) dx = \int_{x_k}^{x_{k+1}} \tilde{f}(x) \frac{\partial T}{\partial t} dx + \left[u T \tilde{f} \right]_{x_k}^{x_{k+1}} - \int_{x_k}^{x_{k+1}} \frac{\partial \tilde{f}}{\partial x} u T dx = 0$$

with \tilde{f} the test function. Inside the elements the test functions are defined by well defined polynomials. We once again define

$$\begin{aligned} \tilde{f}_k^+ &= \tilde{f}(x_k^+) \\ \tilde{f}_{k+1}^- &= \tilde{f}(x_{k+1}^-) \end{aligned}$$

$$\int_{x_k}^{x_{k+1}} \left(\tilde{f}(x) \frac{\partial T_h}{\partial t} - \frac{\partial \tilde{f}}{\partial x} u T_h \right) dx + \tilde{f}(x_{k+1}) \widehat{uT}(T_{k+1}^-, T_{k+1}^+) - \tilde{f}(x_k) \widehat{uT}(T_k^-, T_k^+) = 0 \quad (8.44)$$

For a constant u or a linear problem, an effective numerical flux is the Lax-Friedrichs flux:

$$\widehat{uT}(a, b) = u \frac{(a+b)}{2} - |u| \frac{(b-a)}{2} \quad (8.45)$$

when $u > 0$ this flux then simply becomes:

$$uT(a, b) = ua \quad (8.46)$$

which is in essence an upwinding scheme. Filling this into equation (8.44) gives²:

$$\int_{x_k}^{x_{k+1}} \left(\tilde{f}(x) \frac{\partial T_h}{\partial t} - \frac{\partial \tilde{f}}{\partial x} u T_h \right) dx + \tilde{f}_{k+1}^- u T_{k+1}^- - \tilde{f}_k^+ u T_k^- = 0 \quad (8.47)$$

The function T_h inside the element can be approximated as follows:

$$T_h(x) = \sum_{i=1}^m \mathcal{N}_i(x) T_i = \mathcal{N}_k(x) T_k^+ + \mathcal{N}_{k+1}(x) T_{k+1}^- \quad (8.48)$$

where the red quantities are the unknown dofs. In what follows we coin $\dot{T} = \partial T / \partial t$ so

$$\dot{T}_h(x) = \sum_{i=1}^m \mathcal{N}_i(x) \dot{T}_i = \mathcal{N}_k(x) \dot{T}_k^+ + \mathcal{N}_{k+1}(x) \dot{T}_{k+1}^- \quad (8.49)$$

Taking $\tilde{f}(x) = \mathcal{N}_k(x)$ and then $\tilde{f}(x) = \mathcal{N}_{k+1}(x)$ we arrive at

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_k(x) \dot{T}_h - \frac{\partial \mathcal{N}_k}{\partial x} u T_h \right) dx + \underbrace{\mathcal{N}_k(x_{k+1}^-)}_{=0} u T_{k+1}^- - \underbrace{\mathcal{N}_k(x_k^-)}_{=1} u T_k^- = 0 \quad (8.50)$$

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_{k+1}(x) \dot{T}_h - \frac{\partial \mathcal{N}_{k+1}}{\partial x} u T_h \right) dx + \underbrace{\mathcal{N}_{k+1}(x_{k+1}^-)}_{=1} u T_{k+1}^- - \underbrace{\mathcal{N}_{k+1}(x_k^-)}_{=0} u T_k^- = 0. \quad (8.51)$$

²why do we only keep T^- ? because of the velocity direction?

The underbraced terms are either 0 or 1 due to the properties of the basis functions $\mathcal{N}_i(\vec{r}_j) = \delta_{ij}$. Finally:

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_k(x) \dot{T}_h - \frac{\partial \mathcal{N}_k}{\partial x} u T_h \right) dx - u T_k^- = 0 \quad (8.52)$$

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_{k+1}(x) \dot{T}_h - \frac{\partial \mathcal{N}_{k+1}}{\partial x} u T_h \right) dx + u \underline{T}_{k+1}^- = 0 \quad (8.53)$$

We now use Eqs. (8.48) and (8.49) in Eq. (8.52) and Eq. (8.53):

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_k [\mathcal{N}_k \dot{T}_k^+ + \mathcal{N}_{k+1} \underline{\dot{T}}_{k+1}^-] - \frac{\partial \mathcal{N}_k}{\partial x} u [\mathcal{N}_k T_k^+ + \mathcal{N}_{k+1} \underline{T}_{k+1}^-] \right) dx - u T_k^- = 0 \quad (8.54)$$

$$\int_{x_k}^{x_{k+1}} \left(\mathcal{N}_{k+1} [\mathcal{N}_k \dot{T}_k^+ + \mathcal{N}_{k+1} \underline{\dot{T}}_{k+1}^-] - \frac{\partial \mathcal{N}_{k+1}}{\partial x} u [\mathcal{N}_k T_k^+ + \mathcal{N}_{k+1} \underline{T}_{k+1}^-] \right) dx + u \underline{T}_{k+1}^- = 0 \quad (8.55)$$

Defining again (see Appendix E)

$$\mathbf{M}_e = \int_{x_k}^{x_{k+1}} \begin{pmatrix} \mathcal{N}_k \mathcal{N}_k & \mathcal{N}_k \mathcal{N}_{k+1} \\ \mathcal{N}_{k+1} \mathcal{N}_k & \mathcal{N}_{k+1} \mathcal{N}_{k+1} \end{pmatrix} dx = \frac{h}{6} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

and

$$\mathbf{K}^e = \int_{\Omega_e} \begin{pmatrix} \frac{d\mathcal{N}_k}{dx} \mathcal{N}_k & \frac{d\mathcal{N}_k}{dx} \mathcal{N}_{k+1} \\ \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_k & \frac{d\mathcal{N}_{k+1}}{dx} \mathcal{N}_{k+1} \end{pmatrix} dV = \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix}$$

This results in:

$$\mathbf{M}_e \cdot \begin{pmatrix} \dot{T}_k^+ \\ \underline{\dot{T}}_{k+1}^- \end{pmatrix} - u \mathbf{K}_e \cdot \begin{pmatrix} T_k^+ \\ \underline{T}_{k+1}^- \end{pmatrix} + u \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ \underline{T}_{k+1}^- \end{pmatrix} - \begin{pmatrix} u T_k^- \\ 0 \end{pmatrix} = 0 \quad (8.56)$$

or,

$$\mathbf{M}_e \cdot \begin{pmatrix} \dot{T}_k^+ \\ \underline{\dot{T}}_{k+1}^- \end{pmatrix} = u \left[\mathbf{K}_e - \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \right] \cdot \begin{pmatrix} T_k^+ \\ \underline{T}_{k+1}^- \end{pmatrix} + u \begin{pmatrix} T_k^- \\ 0 \end{pmatrix} \quad (8.57)$$

$$\mathbf{M}_e \cdot \begin{pmatrix} \dot{T}_k^+ \\ \underline{\dot{T}}_{k+1}^- \end{pmatrix} = u \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & -1 \end{pmatrix} \cdot \begin{pmatrix} T_k^+ \\ \underline{T}_{k+1}^- \end{pmatrix} + u \begin{pmatrix} T_k^- \\ 0 \end{pmatrix} \quad (8.58)$$

$$\begin{pmatrix} \dot{T}_k^+ \\ \underline{\dot{T}}_{k+1}^- \end{pmatrix} = u \mathbf{M}_e^{-1} \cdot \frac{1}{2} \begin{pmatrix} -1 & -1 \\ 1 & -1 \end{pmatrix} \cdot \begin{pmatrix} T_k^+ \\ \underline{T}_{k+1}^- \end{pmatrix} + u \mathbf{M}_e^{-1} \cdot \begin{pmatrix} T_k^- \\ 0 \end{pmatrix} \quad (8.59)$$

We have already established that

$$\mathbf{M}_e^{-1} = \frac{2}{h} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$$

so

$$\boxed{\begin{pmatrix} \dot{T}_k^+ \\ \underline{\dot{T}}_{k+1}^- \end{pmatrix} = \frac{u}{h} \begin{pmatrix} -3 & -1 \\ 3 & -1 \end{pmatrix} \begin{pmatrix} T_k^+ \\ \underline{T}_{k+1}^- \end{pmatrix} + \frac{u}{h} \begin{pmatrix} 4 \\ -2 \end{pmatrix} T_k^-} \quad (8.60)$$

Using the same first order Runge-Kutta method as in the previous section,

$$T^k(t + \delta t) = T_k(t) + \delta t \dot{T}_k \quad (8.61)$$

we multiply the equation by δt and we obtain

$$\begin{pmatrix} T_k^+(t + \delta t) \\ \underline{T}_{k+1}^-(t + \delta t) \end{pmatrix} = \begin{pmatrix} T_k^+(t) \\ \underline{T}_{k+1}^-(t) \end{pmatrix} + \frac{u \delta t}{h} \begin{pmatrix} -3 & -1 \\ 3 & -1 \end{pmatrix} \begin{pmatrix} T_k^+(t) \\ \underline{T}_{k+1}^-(t) \end{pmatrix} + \frac{u \delta t}{h} \begin{pmatrix} 4 \\ -2 \end{pmatrix} T_k^-(t + \delta t) \quad (8.62)$$

and finally

$$\begin{pmatrix} T_k^+(t + \delta t) \\ T_{k+1}^-(t + \delta t) \end{pmatrix} = \left[\mathbf{1} + \frac{u\delta t}{h} \begin{pmatrix} -3 & -1 \\ 3 & -1 \end{pmatrix} \right] \cdot \begin{pmatrix} T_k^+(t) \\ T_{k+1}^-(t) \end{pmatrix} + \frac{u\delta t}{h} \begin{pmatrix} 4 \\ -2 \end{pmatrix} T_k^-(t + \delta t) \quad (8.63)$$

Also, since $C = u\delta t/h$, then the equation above can also be written

$$\begin{pmatrix} T_k^+(t + \delta t) \\ T_{k+1}^-(t + \delta t) \end{pmatrix} = \left[\mathbf{1} + C \begin{pmatrix} -3 & -1 \\ 3 & -1 \end{pmatrix} \right] \cdot \begin{pmatrix} T_k^+(t) \\ T_{k+1}^-(t) \end{pmatrix} + C \begin{pmatrix} 4 \\ -2 \end{pmatrix} T_k^-(t + \delta t) \quad (8.64)$$

This problem can be solved starting from the left boundary and sweeping through all the elements. The updated values of adjacent elements are used in the calculation of the next element as soon as this element becomes available which is why the last term T_k^- is taken at $t + \delta t$.

8.5 Steady-state diffusion in 2D

Let us start from the 2D steady state heat diffusion equation:

$$\vec{\nabla} \cdot k \vec{\nabla} T + H = 0 \quad (8.65)$$

Just as in the 1D case this equation can be split in two separate first order differential equations:

$$\underbrace{-\vec{\nabla} \cdot \vec{q} + H = 0}_{\text{ODE 1}} \quad ; \quad \underbrace{\vec{q} = -k \vec{\nabla} T}_{\text{ODE 2}} \quad (8.66)$$

Let N_i^θ be the temperature basis functions so that the temperature inside an element is given by

$$T_h(\vec{r}) = \sum_{i=1}^m N_i^\theta(\vec{r}) T_i = \vec{N}^\theta \cdot \vec{T} \quad (8.67)$$

where \vec{T} is a vector of length m , the number of nodes per element. Similarly we let the basis function for the heat flux be such that

$$q_h(\vec{r}) = \sum_{i=1}^m N_i^q(\vec{r}) q_i = \vec{N}^q \cdot \vec{q} \quad (8.68)$$

$$q_h(\vec{r}) = \sum_{i=1}^m N_i^q(\vec{r}) q_i = \vec{N}^q \cdot \vec{q} \quad (8.69)$$

where \vec{q} , \vec{q} and \vec{N}^q are vectors of length m too. Implicitly if m is the same for temperature and heat flux, then $N^\theta = N^q$. Let us establish the weak forms of the 1st order ODEs.

ODE 1 This results in:

$$\int_{\Omega} N_i^\theta \vec{\nabla} \cdot \vec{q} dV = \int_{\Omega} N_i^\theta H dV \quad (8.70)$$

Using the product rule $\vec{\nabla} \cdot (N_i^\theta \vec{q}) = N_i^\theta \vec{\nabla} \cdot \vec{q} + \vec{\nabla} N_i^\theta \cdot \vec{q}$ then we can write

$$N_i^\theta \vec{\nabla} \cdot \vec{q} = \vec{\nabla} \cdot (N_i^\theta \vec{q}) - \vec{\nabla} N_i^\theta \cdot \vec{q}$$

which we insert in Eq. (8.70) results in

$$\int_{\Omega} \vec{\nabla} \cdot (N_i^\theta \vec{q}) dV - \int_{\Omega} \vec{\nabla} N_i^\theta \cdot \vec{q} dV = \int_{\Omega} N_i^\theta H dV \quad (8.71)$$

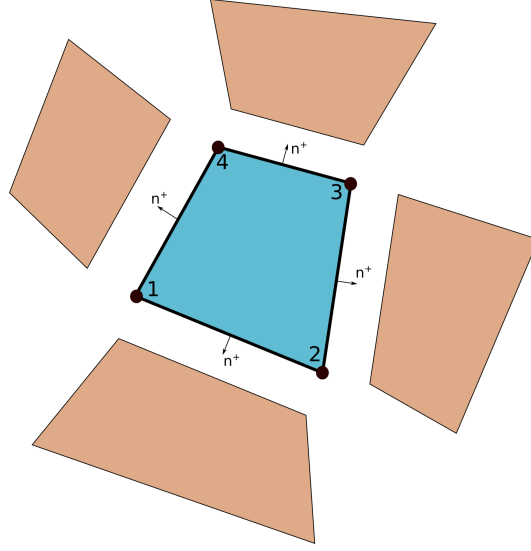
Using the divergence theorem $\int_{\Omega} \vec{\nabla} \cdot \vec{F} d\Omega = \int_{d\Omega} (\vec{F} \cdot \vec{n}) dS$ applied to Eq. (8.71) leads to

$$\int_{d\Omega} N_i^\theta \vec{q} \cdot \vec{n} dS - \int_{\Omega} \vec{\nabla} N_i^\theta \cdot \vec{q} dV = \int_{\Omega} N_i^\theta H dV \quad (8.72)$$

Here \vec{n} is the outward vector everywhere on the boundary. The exact solution $\vec{q} = (q, q)$ can be approximated with \vec{q}_h in a finite element space, same for the approximation of the flux at the boundary $\hat{q} = \hat{q}_h$ which takes a special form in the context of the DG methods (as reflected by the presence of the hat). We decompose the resulting equation in 3 terms A , B and C :

$$\underbrace{\int_{d\Omega} N_i^\theta \vec{q}_h \cdot \vec{n} dS}_A - \underbrace{\int_{\Omega} \vec{\nabla} N_i^\theta \cdot \vec{q}_h dV}_B = \underbrace{\int_{\Omega} N_i^\theta H dV}_C \quad (8.73)$$

Note that terms B and C are explained in Section 6.4. In all what follows blue symbols belong the the element under consideration and brown symbols belong to its neighbour(s).



\vec{n}^+ indicates the outward vector at the boundary and \vec{n}^- is the outward vector of the neighbouring element.

Let us turn to Li's book for useful definitions:

- Definition of jump operators The square brackets denote the jump operator:

$$\begin{aligned} [\vec{q}_h] &= \vec{q}_h \cdot \vec{n}^+ + \vec{q}_h \cdot \vec{n}^- & \text{or} & & [\vec{q}_h] &= (\vec{q}_h - \vec{q}_h) \cdot \vec{n}^+ \\ [T_h] &= T_h \vec{n}^+ + T_h \vec{n}^- & \text{or} & & [T_h] &= (T_h - T_h) \vec{n}^+ \end{aligned} \quad (8.74)$$

Note that $[\vec{q}_h]$ is a scalar function which involves the normal components only; while $[T_h]$ is a vector function.

- Definition of average operators The curly brackets indicate the average operator

$$\begin{aligned} \{\vec{q}_h\} &= \frac{1}{2}(\vec{q}_h + \vec{q}_h) & \text{so} & & \{q_h\}_x &= \frac{1}{2}(q_h + q_h) & \text{and} & & \{q_h\}_y &= \frac{1}{2}(q_h + q_h) \\ [T_h] &= \frac{1}{2}(T_h + T_h) \end{aligned}$$

Note that $\{\vec{q}_h\}$ is a vector, while $\{T_h\}$ is a scalar.

- Definitions of fluxes In the LDG method the boundary flux \vec{q}_h is defined as

$$\vec{q}_h = \{\vec{q}_h\} - \mathcal{E}[T_h] - \vec{\mathcal{C}}[\vec{q}_h]$$

where \mathcal{E} is a scalar (since $[T_h]$ is a vector) and $\vec{\mathcal{C}}$ is a vector (since $[\vec{q}_h]$ is a scalar). To be once again very explicit, the above equation writes as follows for a 2D Cartesian space:

$$\begin{aligned} \hat{q}_h_x &= \{q_h\}_x - \mathcal{E}[T_h]_x - \mathcal{C}_x[\vec{q}_h] \\ &= \frac{1}{2}(q_h + q_h)_x - \mathcal{E}(T_h n_x^+ + T_h n_x^-) - \mathcal{C}_x(\vec{q}_h \cdot \vec{n}^+ + \vec{q}_h \cdot \vec{n}^-) \end{aligned} \quad (8.75)$$

$$\begin{aligned} \hat{q}_h_y &= \{q_h\}_y - \mathcal{E}[T_h]_y - \mathcal{C}_y[\vec{q}_h] \\ &= \frac{1}{2}(q_h + q_h)_y - \mathcal{E}(T_h n_y^+ + T_h n_y^-) - \mathcal{C}_y(\vec{q}_h \cdot \vec{n}^+ + \vec{q}_h \cdot \vec{n}^-) \end{aligned} \quad (8.76)$$

Remark. In the book the note under table 4.1 states that the C_{ij} coefficients are constant matrices, which is quite misleading since some are actually scalars and others vectors.

Filling Eqs. (8.75) and (8.76) into A of Eq. (8.73) leads to

$$\begin{aligned}
A &= \int_{\partial\Omega} N_i^\theta \vec{q}_h \cdot \vec{n} \, dS \\
&= \int_{\partial\Omega} N_i^\theta \left[\{\vec{q}_h\} - \mathcal{E}[T_h] - \vec{\mathcal{C}}[\vec{q}_h] \right] \cdot \vec{n}^+ \, dS \\
&= \int_{\partial\Omega} N_i^\theta \left[\frac{1}{2}(\vec{q}_h + \vec{q}_h) - \mathcal{E}(\textcolor{blue}{T}_h \vec{n}^+ + \textcolor{brown}{T}_h \vec{n}^-) - \vec{\mathcal{C}}[\vec{q}_h] \right] \cdot \vec{n}^+ \, dS \\
&= \int_{\partial\Omega} N_i^\theta \left[\frac{1}{2}(\vec{q}_h + \vec{q}_h) \cdot \vec{n}^+ - \mathcal{E}(\textcolor{blue}{T}_h \vec{n}^+ + \textcolor{brown}{T}_h \vec{n}^-) \cdot \vec{n}^+ - \vec{\mathcal{C}} \cdot \vec{n}^+ [\vec{q}_h] \right] \, dS \\
&= \int_{\partial\Omega} N_i^\theta \left[\frac{1}{2}(\vec{q}_h + \vec{q}_h) \cdot \vec{n}^+ - \mathcal{E}(\textcolor{blue}{T}_h - \textcolor{brown}{T}_h) - (\vec{\mathcal{C}} \cdot \vec{n}^+) (\vec{q}_h - \vec{q}_h) \cdot \vec{n}^+ \right] \, dS \quad \text{since } \vec{n}^+ \cdot \vec{n}^+ = 1 \quad \vec{n}^+ \cdot \vec{n}^- = 0 \\
&= \int_{\partial\Omega} N_i^\theta \left[\left(\frac{1}{2} - \vec{\mathcal{C}} \cdot \vec{n}^+ \right) \vec{q}_h \cdot \vec{n}^+ - \mathcal{E} \textcolor{blue}{T}_h \right] \, dS + \int_{\partial\Omega} N_i^\theta \left[\left(\frac{1}{2} + \vec{\mathcal{C}} \cdot \vec{n}^+ \right) \vec{q}_h \cdot \vec{n}^+ + \mathcal{E} \textcolor{brown}{T}_h \right] \, dS \\
&= \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} - \vec{\mathcal{C}} \cdot \vec{n}^+ \right) \vec{q}_h \cdot \vec{n}^+ \, dS - \int_{\partial\Omega} N_i^\theta \mathcal{E} \textcolor{blue}{T}_h \, dS + \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} + \vec{\mathcal{C}} \cdot \vec{n}^+ \right) \vec{q}_h \cdot \vec{n}^+ \, dS + \int_{\partial\Omega} N_i^\theta \mathcal{E} \textcolor{brown}{T}_h \, dS \\
&= A_1 + A_2 + A_3 + A_4 \tag{8.}
\end{aligned}$$

In order to simplify notations we choose $N^q = N^\theta = N$ and drop the h subscripts.

$$\begin{aligned}
A_1 &= \int_{\partial\Omega} N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{q} \cdot \vec{n}^+ dS \\
&= \int_{\partial\Omega} N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) (q_x n_x^+ + q_y n_y^+) dS \\
&= \int_{\partial\Omega} N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) q_x n_x^+ dS + \int_{\partial\Omega} N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) q_y n_y^+ dS \\
\Rightarrow \vec{A}_1 &= \left(\int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \right) \cdot \vec{q}_y \quad (8.78)
\end{aligned}$$

$$\begin{aligned}
A_2 &= - \int_{\partial\Omega} N_i \mathcal{E} T dS \\
\Rightarrow \vec{A}_2 &= - \left(\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \right) \cdot \vec{T} \quad (8.79)
\end{aligned}$$

$$\begin{aligned}
A_3 &= \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{q} \cdot \vec{n}^+ dS \\
&= \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) (q_x n_x^+ + q_y n_y^+) dS \\
&= \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) q_x n_x^+ dS + \int_{\partial\Omega} N_i^\theta \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) q_y n_y^+ dS \\
\Rightarrow \vec{A}_3 &= \left(\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \right) \cdot \vec{q}_y \quad (8.80)
\end{aligned}$$

$$\begin{aligned}
A_4 &= \int_{\partial\Omega} N_i \mathcal{E} T dS \\
\Rightarrow \vec{A}_4 &= \left(\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \right) \cdot \vec{T} \quad (8.81)
\end{aligned}$$

$$\begin{aligned}
B &= \int_{\Omega} \vec{\nabla} N_i \cdot \vec{q} dV \\
&= \int_{\Omega} (\partial_x N_i q_x + \partial_y N_i q_y) dV \\
&= \int_{\Omega} \partial_x N_i q_x dV + \int_{\Omega} \partial_y N_i q_y dV \\
\Rightarrow \vec{B} &= \left(\int_{\Omega} \partial_x \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_x + \left(\int_{\Omega} \partial_y \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_y \quad (8.82)
\end{aligned}$$

$$\begin{aligned}
C &= \int_{\Omega} N_i H dV \\
\Rightarrow \vec{C} &= \int_{\Omega} \vec{N}^T H dV \quad (8.83)
\end{aligned}$$

The expressions above find their equivalent in the book (NB stands for neighbour):

$$\begin{aligned}
&\mathbf{B} \left(\int_{\Omega_j} (\nabla \Phi) \Phi^T dV \right) \cdot \underline{\mathbf{q}} - \left(\int_{\partial\Omega_j} (1/2 - C_{12}) \Phi \Phi^T \mathbf{n} dS \right) \cdot \underline{\mathbf{q}} \quad \mathbf{A}_1 \\
&\mathbf{A}_3 \left(- \int_{\partial\Omega_j} (1/2 + C_{12}) \Phi \Phi^T \mathbf{n} dS \right) \cdot \underline{\mathbf{q}}_{(NB)} - \left(\int_{\partial\Omega_j} C_{11} \Phi \Phi^T dS \right) \underline{\mathbf{T}} \quad \mathbf{A}_2 \\
&\mathbf{A}_4 \left(- \int_{\partial\Omega_j} C_{11} \Phi \Phi^T dS \right) \underline{\mathbf{T}}_{(NB)} - \left(\int_{\Omega_j} \Phi Q dV \right) \underline{\mathbf{c}} \quad (4.23b)
\end{aligned}$$

check minus signs

ODE 2 Its weak form writes:

$$\int_{\Omega} N_i^q (\vec{q} + k \vec{\nabla} T) dV = 0$$

or, (once again we drop the superscript on the basis functions and the h):

$$\begin{aligned} 0 &= \int_{\Omega} N_i \vec{q} dV + \int_{\Omega} N_i k \vec{\nabla} T dV \\ &= \int_{\Omega} N_i \vec{q} dV + \int_{\Omega} \vec{\nabla} (N_i k T) dV - \int_{\Omega} \vec{\nabla} (k N_i) T dV \end{aligned} \quad (8.84)$$

We then use $\int_{\Omega} \vec{\nabla} f dV = \int_{\partial\Omega} f \vec{S} dS$ and as before the temperature on the edge integral should be \hat{T} :

$$\int_{\Omega} N_i \vec{q} dV + \int_{\partial\Omega} k N_i \hat{T} \vec{n}^+ dS - \int_{\Omega} \vec{\nabla} (k N_i) T dV = 0$$

or, if decomposed in a 2D Cartesian axis system

$$\begin{aligned} 0 &= \underbrace{\int_{\Omega} N_i q_x dV}_D + \underbrace{\int_{\partial\Omega} k N_i \hat{T} n_x^+ dS}_E - \underbrace{\int_{\Omega} \partial_x (k N_i) T dV}_F & (0 = D + E - F) \\ 0 &= \underbrace{\int_{\Omega} N_i q_y dV}_G + \underbrace{\int_{\partial\Omega} k N_i \hat{T} n_y^+ dS}_H - \underbrace{\int_{\Omega} \partial_y (k N_i) T dV}_I & (0 = G + H - I) \end{aligned} \quad (8.85)$$

which (aside from a minus sign coming from a different definition of the heat flux) is 'identical' to the book (although the notations in the book are hella confusing):

$$\int_{\Omega_j} \mathbf{q}_h \cdot \mathbf{w} dV = - \int_{\Omega_j} T_h \nabla \cdot (\kappa \mathbf{w}) dV + \int_{\partial\Omega_j} \kappa \hat{T}_h \mathbf{n}_j \cdot \mathbf{w} dS \quad (\text{with } \kappa \rightarrow k, \mathbf{w} \rightarrow \vec{N})$$

The temperature flux \hat{T} is chosen to be (table 4.1 in book):

$$\hat{T}_h = \{T_h\} + \vec{C}[T_h] - \mathcal{F}[\vec{q}_h] \quad (8.86)$$

$$= \frac{1}{2} (T_h + T_h) + \vec{C} (T_h \vec{n}^+ + T_h \vec{n}^-) - \mathcal{F} (\vec{q}_h \cdot \vec{n}^+ + \vec{q}_h \cdot \vec{n}^-) \quad (8.87)$$

$$= \frac{1}{2} (T_h + T_h) + (T_h - T_h) \vec{C} \cdot \vec{n}^+ - \mathcal{F} (\vec{q}_h - \vec{q}_h) \cdot \vec{n}^+ \quad (8.88)$$

$$= \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) T_h + \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) T_h - \mathcal{F} (\vec{q}_h - \vec{q}_h) \cdot \vec{n}^+ \quad (8.89)$$

In what follows we assume k to be constant within each element so that $\partial_x (k N_i) = k \partial_x N_i$.

$$D = \int_{\Omega} N_i q_x dV \Rightarrow \vec{D} = \left(\int_{\Omega} \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_x \quad (8.90)$$

$$F = \int_{\Omega} k \partial_x N_i T dV \Rightarrow \vec{F} = \left(\int_{\Omega} k \partial_x \vec{N}^T \vec{N} dV \right) \cdot \vec{T} \quad (8.91)$$

$$G = \int_{\Omega} N_i q_y^h dV \Rightarrow \vec{G} = \left(\int_{\Omega} \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_y \quad (8.92)$$

$$I = \int_{\Omega} k \partial_y N_i T dV \Rightarrow \vec{I} = \left(\int_{\Omega} k \partial_y \vec{N}^T \vec{N} dV \right) \cdot \vec{T} \quad (8.93)$$

$$\begin{aligned} E &= \int_{\partial\Omega} k N_i \hat{T} n_x^+ dS \\ &= \int_{\partial\Omega} k N_i \left[\left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) T_h + \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) T_h - \mathcal{F}(\vec{q}_h - \vec{q}_h) \cdot \vec{n}^+ \right] n_x^+ dS \\ &= \int_{\partial\Omega} k N_i \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) T_h n_x^+ dS + \int_{\partial\Omega} k N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) T_h n_x^+ dS \\ &\quad - \int_{\partial\Omega} k N_i \mathcal{F}(\vec{q}_h \cdot \vec{n}^+) n_x^+ dS + \int_{\partial\Omega} k N_i \mathcal{F}(\vec{q}_h \cdot \vec{n}^+) n_x^+ dS \\ &= E_1 + E_2 + E_3 + E_4 \\ H &= \int_{\partial\Omega} k N_i \hat{T} n_y^+ dS \\ &= H_1 + H_2 + H_3 + H_4 \end{aligned} \quad (8.94)$$

Then

$$\begin{aligned} E_1 &= \int_{\partial\Omega} k N_i \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) T_h n_x^+ dS \Rightarrow \left(\int_{\partial\Omega} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \right) \cdot \vec{T} \\ E_2 &= \int_{\partial\Omega} k N_i \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) T_h n_x^+ dS \Rightarrow \left(\int_{\partial\Omega} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \right) \cdot \vec{T} \\ E_3 &= - \int_{\partial\Omega} k N_i \mathcal{F}(\vec{q}_h \cdot \vec{n}^+) n_x^+ dS \\ &= - \int_{\partial\Omega} k N_i \mathcal{F}(q_x n_x^+ + q_y n_y^+) n_x^+ dS \\ &= - \int_{\partial\Omega} k N_i \mathcal{F} q_x n_x^+ n_x^+ dS - \int_{\partial\Omega} k N_i \mathcal{F} q_y n_y^+ n_x^+ dS \\ &\Rightarrow - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \right) \cdot \vec{q}_x - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \right) \cdot \vec{q}_y \\ E_4 &= \int_{\partial\Omega} k N_i \mathcal{F}(\vec{q}_h \cdot \vec{n}^+) n_x^+ dS \\ &= \int_{\partial\Omega} k N_i \mathcal{F}(q_x n_x^+ + q_y n_y^+) n_x^+ dS \\ &= \int_{\partial\Omega} k N_i \mathcal{F} q_x n_x^+ n_x^+ dS + \int_{\partial\Omega} k N_i \mathcal{F} q_y n_y^+ n_x^+ dS \\ &\Rightarrow \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \right) \cdot \vec{q}_y \end{aligned} \quad (8.95)$$

The H_i terms are so similar to the E_i terms that there is not need to write them out explicitly.

We have seen that ODE #2 and #1 write

$$D + E - F = D + E_1 + E_2 + E_3 + E_4 - F = 0 \quad (8.96)$$

$$G + H - I = G + H_1 + H_2 + H_3 + H_4 - I = 0 \quad (8.97)$$

$$A - B = A_1 + A_2 + A_3 + A_4 - B = \vec{C} \quad (8.98)$$

so that

$$\begin{aligned}
& \left(\int_{\Omega} \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \right) \cdot \vec{T} + \left(\int_{\partial\Omega} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \right) \cdot \vec{T} \\
& - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \right) \cdot \vec{q}_x - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \right) \cdot \vec{q}_y + \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \right) \cdot \vec{q}_y - \left(\int_{\Omega} k \partial_x \vec{N}^T \vec{N} dV \right) \cdot \vec{T} = \\
& \left(\int_{\Omega} \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_y + \left(\int_{\partial\Omega} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \right) \cdot \vec{T} + \left(\int_{\partial\Omega} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \right) \cdot \vec{T} \\
& - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS \right) \cdot \vec{q}_x - \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS \right) \cdot \vec{q}_y + \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS \right) \cdot \vec{q}_y - \left(\int_{\Omega} k \partial_y \vec{N}^T \vec{N} dV \right) \cdot \vec{T} = \\
& \left(\int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \right) \cdot \vec{q}_y - \left(\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \right) \cdot \vec{T} \\
& + \left(\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \right) \cdot \vec{q}_x + \left(\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \right) \cdot \vec{q}_y + \left(\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \right) \cdot \vec{T} - \left(\int_{\Omega} \partial_x \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_x - \left(\int_{\Omega} \partial_y \vec{N}^T \vec{N} dV \right) \cdot \vec{q}_y = \\
& \left[\underbrace{\begin{pmatrix} \int_{\Omega} \vec{N}^T \vec{N} dV & 0 & -\int_{\Omega} k \partial_x \vec{N}^T \vec{N} dV \\ 0 & \int_{\Omega} \vec{N}^T \vec{N} dV & -\int_{\Omega} k \partial_y \vec{N}^T \vec{N} dV \\ -\int_{\Omega} \partial_x \vec{N}^T \vec{N} dV & -\int_{\Omega} \partial_y \vec{N}^T \vec{N} dV & 0 \end{pmatrix}}_{\mathcal{A}_{\Omega}} + \underbrace{\begin{pmatrix} -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS & -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS & \int_{\partial\Omega} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \\ -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS & -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS & \int_{\partial\Omega} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \\ \int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS & \int_{\partial\Omega} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS & -\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \end{pmatrix}}_{\mathcal{A}_{\partial\Omega}} \right] \cdot \begin{pmatrix} \vec{q}_x \\ \vec{q}_y \\ \vec{T} \end{pmatrix} \\
& = \begin{pmatrix} -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS & -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS & -\int_{\partial\Omega} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \\ -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS & -\int_{\partial\Omega} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS & -\int_{\partial\Omega} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \\ -\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS & -\int_{\partial\Omega} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS & -\int_{\partial\Omega} \mathcal{E} \vec{N}^T \vec{N} dS \end{pmatrix} \cdot \begin{pmatrix} \vec{q}_x \\ \vec{q}_y \\ \vec{T} \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ \int_{\Omega} \vec{N}^T \vec{N} dV \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
& \left[\underbrace{\begin{pmatrix} \int_{\Omega} \vec{N}^T \vec{N} dV & 0 & -\int_{\Omega} k \partial_x \vec{N}^T \vec{N} dV \\ 0 & \int_{\Omega} \vec{N}^T \vec{N} dV & -\int_{\Omega} k \partial_y \vec{N}^T \vec{N} dV \\ -\int_{\Omega} \partial_x \vec{N}^T \vec{N} dV & -\int_{\Omega} \partial_y \vec{N}^T \vec{N} dV & 0 \end{pmatrix}}_{\mathcal{A}_{\Omega}} + \sum_{i=1}^{nedges} \underbrace{\begin{pmatrix} -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS & -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS & \int_{\partial\Omega_i} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \\ -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS & -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS & \int_{\partial\Omega_i} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \\ \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS & \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS & -\int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS \end{pmatrix}}_{\mathcal{A}_{\partial\Omega}} \right] \cdot \left(\begin{pmatrix} \vec{q}_x \\ \vec{q}_y \\ \vec{T} \end{pmatrix} \right)_i + \begin{pmatrix} 0 \\ 0 \\ \int_{\Omega} \vec{N}^T H dV \end{pmatrix} \\
& = \sum_{i=1}^{nedges} \begin{pmatrix} -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS & -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS & -\int_{\partial\Omega_i} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \\ -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS & -\int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS & -\int_{\partial\Omega_i} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \\ -\int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS & -\int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS & -\int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS \end{pmatrix} \cdot \begin{pmatrix} \vec{q}_x \\ \vec{q}_y \\ \vec{T} \end{pmatrix}_i + \begin{pmatrix} 0 \\ 0 \\ \int_{\Omega} \vec{N}^T H dV \end{pmatrix} \\
& \left[\begin{pmatrix} \mathbf{E} & \mathbf{0} & \mathbf{H}_x \\ \mathbf{0} & \mathbf{E} & \mathbf{H}_y \\ \mathbf{J}_x & \mathbf{J}_y & \mathbf{0} \end{pmatrix} + \sum_{i=1}^{nedges} \begin{pmatrix} \mathbf{E}_{xx,i} & \mathbf{E}_{xy,i} & \mathbf{H}_{x,i} \\ \mathbf{E}_{yx,i} & \mathbf{E}_{yy,i} & \mathbf{H}_{y,i} \\ \mathbf{J}_{x,i} & \mathbf{J}_{y,i} & \mathbf{G}_{T,i} \end{pmatrix} \right] \cdot \begin{pmatrix} \vec{q}_x \\ \vec{q}_y \\ \vec{T} \end{pmatrix} =
\end{aligned}$$

which is identical to the equation 4.24 in Li's book (if the terms related to the third dimension are disregarded):

$$\begin{aligned}
& \begin{bmatrix} \mathbf{E} & \mathbf{0} & \mathbf{0} & \mathbf{H}_x \\ \mathbf{0} & \mathbf{E} & \mathbf{0} & \mathbf{H}_y \\ \mathbf{0} & \mathbf{0} & \mathbf{E} & \mathbf{H}_z \\ \mathbf{J}_x & \mathbf{J}_y & \mathbf{J}_z & \mathbf{0} \end{bmatrix} \begin{pmatrix} \underline{\mathbf{q}}_x \\ \underline{\mathbf{q}}_y \\ \underline{\mathbf{q}}_z \\ \underline{\mathbf{T}} \end{pmatrix} + \sum_{i=1}^{NS} \begin{bmatrix} \mathbf{E}_{xx,i} & \mathbf{E}_{xy,i} & \mathbf{E}_{xz,i} & \mathbf{H}_{x,i} \\ \mathbf{E}_{yx,i} & \mathbf{E}_{yy,i} & \mathbf{E}_{yz,i} & \mathbf{H}_{y,i} \\ \mathbf{E}_{zx,i} & \mathbf{E}_{zy,i} & \mathbf{E}_{zz,i} & \mathbf{H}_{z,i} \\ \mathbf{J}_{x,i} & \mathbf{J}_{y,i} & \mathbf{J}_{z,i} & \mathbf{G}_{T,i} \end{bmatrix} \begin{pmatrix} \underline{\mathbf{q}}_x \\ \underline{\mathbf{q}}_y \\ \underline{\mathbf{q}}_z \\ \underline{\mathbf{T}} \end{pmatrix} \\
& + \sum_{i=1}^{NS} \begin{bmatrix} \mathbf{E}_{xx,B,i} & \mathbf{E}_{xy,B,i} & \mathbf{E}_{xz,B,i} & \mathbf{H}_{x,B,i} \\ \mathbf{E}_{yx,B,i} & \mathbf{E}_{yy,B,i} & \mathbf{E}_{yz,B,i} & \mathbf{H}_{y,B,i} \\ \mathbf{E}_{zx,B,i} & \mathbf{E}_{zy,B,i} & \mathbf{E}_{zz,B,i} & \mathbf{H}_{z,B,i} \\ \mathbf{J}_{x,B,i} & \mathbf{J}_{y,B,i} & \mathbf{J}_{z,B,i} & \mathbf{G}_{T,B,i} \end{bmatrix} \begin{pmatrix} \underline{\mathbf{q}}_x \\ \underline{\mathbf{q}}_y \\ \underline{\mathbf{q}}_z \\ \underline{\mathbf{T}} \end{pmatrix}_{(NB,i)} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{S}_T \end{pmatrix}
\end{aligned}$$

$$\mathbf{E} = \int_{\Omega} \vec{N}^T \vec{N} dV \quad (8.99)$$

$$\mathbf{H}_x = - \int_{\Omega} k \partial_x \vec{N}^T \vec{N} dV \quad (8.100)$$

$$\mathbf{H}_y = - \int_{\Omega} k \partial_y \vec{N}^T \vec{N} dV \quad (8.101)$$

$$\mathbf{J}_x = - \int_{\Omega} \partial_x \vec{N}^T \vec{N} dV \quad (8.102)$$

$$\mathbf{J}_y = - \int_{\Omega} \partial_y \vec{N}^T \vec{N} dV \quad (8.103)$$

$$\mathbf{E}_{xx,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \quad (8.104)$$

$$\mathbf{E}_{xy,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \quad (8.105)$$

$$\mathbf{E}_{yx,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS \quad (8.106)$$

$$\mathbf{E}_{yy,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS \quad (8.107)$$

$$\mathbf{H}_{x,i} = \int_{\partial\Omega_i} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \quad (8.108)$$

$$\mathbf{H}_{y,i} = \int_{\partial\Omega_i} k \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \quad (8.109)$$

$$\mathbf{J}_{x,i} = \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \quad (8.110)$$

$$\mathbf{J}_{y,i} = \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \quad (8.111)$$

$$\mathbf{G}_{T,i} = - \int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS \quad (8.112)$$

$$\mathbf{E}_{xx,B,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \quad (8.113)$$

$$\mathbf{E}_{xy,B,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \quad (8.114)$$

$$\mathbf{E}_{yx,B,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS \quad (8.115)$$

$$\mathbf{E}_{yy,B,i} = - \int_{\partial\Omega_i} k \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS \quad (8.116)$$

$$\mathbf{H}_{x,B,i} = - \int_{\partial\Omega_i} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS \quad (8.117)$$

$$\mathbf{H}_{y,B,i} = - \int_{\partial\Omega_i} k \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS \quad (8.118)$$

$$\mathbf{J}_{x,B,i} = - \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \quad (8.119)$$

$$\mathbf{J}_{y,B,i} = - \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \quad (8.120)$$

$$\mathbf{G}_{T,B,i} = - \int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS \quad (8.121)$$

$$\mathbf{S}_T = \int_{\Omega} \vec{N}^T H dV \quad (8.122)$$

Note that $\mathcal{E} = C_{11}$ and $\mathcal{F} = C_{22}$ in the book.

If k is constant per element, then:

$$\mathbf{E} = \int_{\Omega} \vec{N}^T \vec{N} dV \quad (8.123)$$

$$\mathbf{H}_x = -k \int_{\Omega} \partial_x \vec{N}^T \vec{N} dV \quad (8.124)$$

$$\mathbf{H}_y = -k \int_{\Omega} \partial_y \vec{N}^T \vec{N} dV \quad (8.125)$$

$$\mathbf{J}_x = - \int_{\Omega} \partial_x \vec{N}^T \vec{N} dV \quad (8.126)$$

$$\mathbf{J}_y = - \int_{\Omega} \partial_y \vec{N}^T \vec{N} dV \quad (8.127)$$

$$\mathbf{E}_{xx,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS \quad (8.128)$$

$$\mathbf{E}_{xy,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS \quad (8.129)$$

$$\mathbf{E}_{yx,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS \quad (8.130)$$

$$\mathbf{E}_{yy,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS \quad (8.131)$$

$$\mathbf{H}_{x,i} = k \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS = -k \mathbf{J}_{x,B,i} \quad (8.132)$$

$$\mathbf{H}_{y,i} = k \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS = -k \mathbf{J}_{y,B,i} \quad (8.133)$$

$$\mathbf{J}_{x,i} = \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \quad (8.134)$$

$$\mathbf{J}_{y,i} = \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \quad (8.135)$$

$$\mathbf{G}_{T,i} = - \int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS \quad (8.136)$$

$$\mathbf{E}_{xx,B,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_x^+ dS = \mathbf{E}_{xx,i} \quad (8.137)$$

$$\mathbf{E}_{xy,B,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_x^+ dS = \mathbf{E}_{xy,i} \quad (8.138)$$

$$\mathbf{E}_{yx,B,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_x^+ n_y^+ dS = \mathbf{E}_{yx,i} \quad (8.139)$$

$$\mathbf{E}_{yy,B,i} = -k \int_{\partial\Omega_i} \vec{N}^T \vec{N} \mathcal{F} n_y^+ n_y^+ dS = \mathbf{E}_{yy,i} \quad (8.140)$$

$$\mathbf{H}_{x,B,i} = -k \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_x^+ dS = -k \mathbf{J}_{x,i} \quad (8.141)$$

$$\mathbf{H}_{y,B,i} = -k \int_{\partial\Omega_i} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) N^T \vec{N} n_y^+ dS = -k \mathbf{J}_{y,i} \quad (8.142)$$

$$\mathbf{J}_{x,B,i} = - \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS \quad (8.143)$$

$$\mathbf{J}_{y,B,i} = - \int_{\partial\Omega_i} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS \quad (8.144)$$

$$\mathbf{G}_{T,B,i} = - \int_{\partial\Omega_i} \mathcal{E} \vec{N}^T \vec{N} dS = \mathbf{G}_{T,i} \quad (8.145)$$

$$\mathbf{S}_T = \int_{\Omega} \vec{N}^T H dV \quad (8.146)$$

8.5.1 The special case of linear rectangular elements

Let us start with \mathcal{A}_e where we assume that k is constant within an element:

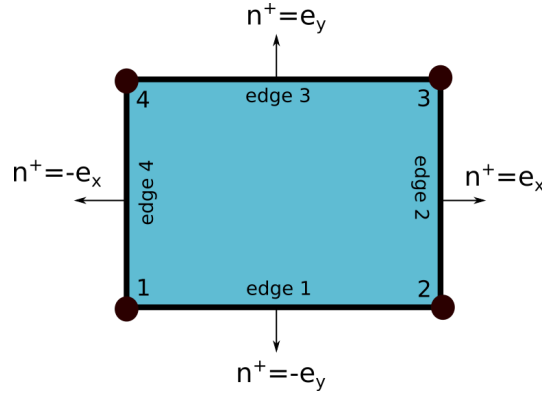
$$\mathcal{A}_{\Omega_e} = \begin{pmatrix} \mathbf{E} & \mathbf{0} & \mathbf{H}_x \\ \mathbf{0} & \mathbf{E} & \mathbf{H}_y \\ \mathbf{J}_x & \mathbf{J}_y & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{E} & \mathbf{0} & -k_e \mathbf{J}_x \\ \mathbf{0} & \mathbf{E} & -k_e \mathbf{J}_y \\ \mathbf{J}_x & \mathbf{J}_y & \mathbf{0} \end{pmatrix} \quad (8.147)$$

The matrices \mathbf{E} , \mathbf{H}_x , and \mathbf{H}_y have been analytically derived in Appendix E.0.2:

$$\mathbf{E} = \frac{h_x h_y}{9} \begin{pmatrix} 1 & 1/2 & 1/4 & 1/2 \\ 1/2 & 1 & 1/2 & 1/4 \\ 1/4 & 1/2 & 1 & 1/2 \\ 1/2 & 1/4 & 1/2 & 1 \end{pmatrix} \quad \mathbf{J}_x = \frac{h_y}{12} \begin{pmatrix} -2 & -2 & -1 & -1 \\ 2 & 2 & 1 & 1 \\ 1 & 1 & 2 & 2 \\ -1 & -1 & -2 & -2 \end{pmatrix} \quad \mathbf{J}_y = \frac{h_x}{12} \begin{pmatrix} -2 & -1 & -1 & -2 \\ -1 & -2 & -2 & -1 \\ 1 & 2 & 2 & 1 \\ 2 & 1 & 1 & 2 \end{pmatrix}$$

The matrix \mathcal{A}_{Ω_e} is therefore trivial to implement.

Let us now turn to $\mathcal{A}_{\partial\Omega}$ which is specific to the DG method. Because elements are rectangles then $n_i^+ n_j^+ = 0$ if $i \neq j$ (where i and j are edge indices). Also, if $i = j$ then $n_i^+ n_j^+ = 1$.



Assuming here again that heat conductivities are constant inside an element it then follows that

$$\mathcal{A}_{\partial\Omega_e} = \sum_{i=1}^{nedges} \begin{pmatrix} \mathbf{E}_{xx,i} & \mathbf{E}_{xy,i} & \mathbf{H}_{x,i} \\ \mathbf{E}_{yx,i} & \mathbf{E}_{yy,i} & \mathbf{H}_{y,i} \\ \mathbf{J}_{x,i} & \mathbf{J}_{y,i} & \mathbf{G}_{T,i} \end{pmatrix} = \sum_{i=1}^{nedges} \begin{pmatrix} \mathbf{E}_{xx,i} & \mathbf{0} & \mathbf{H}_{x,i} \\ \mathbf{0} & \mathbf{E}_{yy,i} & \mathbf{H}_{y,i} \\ \mathbf{J}_{x,i} & \mathbf{J}_{y,i} & \mathbf{G}_{T,i} \end{pmatrix} \quad (8.148)$$

Note that

- $i = 1$: bottom edge, i.e. $s = -1$ and then $N_4 = N_3 = 0$; Also $\vec{n}_x^+ = 0$, $\vec{n}_y^+ = -1$
- $i = 2$: right edge, i.e. $r = +1$ and then $N_1 = N_4 = 0$; Also $\vec{n}_x^+ = 1$, $\vec{n}_y^+ = 0$
- $i = 3$: top edge, i.e. $s = +1$ and then $N_1 = N_2 = 0$; Also $\vec{n}_x^+ = 0$, $\vec{n}_y^+ = 1$
- $i = 4$: left edge, i.e. $r = -1$ and then $N_2 = N_3 = 0$; Also $\vec{n}_x^+ = -1$, $\vec{n}_y^+ = 0$

Then

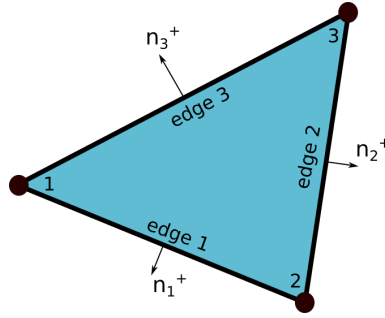
$$\begin{aligned} \mathbf{G}_{T,1} &= -\mathcal{E} \int_{\partial\Omega_1} \vec{N}^T \vec{N} dS = -\mathcal{E} \mathbf{C}_1 \\ \mathbf{G}_{T,2} &= -\mathcal{E} \int_{\partial\Omega_2} \vec{N}^T \vec{N} dS = -\mathcal{E} \mathbf{C}_2 \\ \mathbf{G}_{T,3} &= -\mathcal{E} \int_{\partial\Omega_3} \vec{N}^T \vec{N} dS = -\mathcal{E} \mathbf{C}_3 \\ \mathbf{G}_{T,4} &= -\mathcal{E} \int_{\partial\Omega_4} \vec{N}^T \vec{N} dS = -\mathcal{E} \mathbf{C}_4 \end{aligned}$$

where the matrices \mathbf{C}_i have been worked out in detail in appendix [E.0.2](#):

$$\begin{aligned}\mathbf{C}_1 &= \int_{1 \rightarrow 2} \vec{N}^T \vec{N} dS = \frac{h_x}{6} \begin{pmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \mathbf{C}_2 &= \int_{2 \rightarrow 3} \vec{N}^T \vec{N} dS = \frac{h_y}{6} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 2 & 1 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \mathbf{C}_3 &= \int_{3 \rightarrow 4} \vec{N}^T \vec{N} dS = \frac{h_x}{6} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{pmatrix} \\ \mathbf{C}_4 &= \int_{4 \rightarrow 1} \vec{N}^T \vec{N} dS = \frac{h_y}{6} \begin{pmatrix} 2 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 2 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}
\mathbf{E}_{xx,1} &= -k_e \mathcal{F} \int_{\partial\Omega_1} \vec{N}^T \vec{N} n_x^+ n_x^+ dS = 0 \\
\mathbf{E}_{xx,2} &= -k_e \mathcal{F} \int_{\partial\Omega_2} \vec{N}^T \vec{N} n_x^+ n_x^+ dS = -k_e \mathcal{F} \mathbf{C}_2 \\
\mathbf{E}_{xx,3} &= -k_e \mathcal{F} \int_{\partial\Omega_3} \vec{N}^T \vec{N} n_x^+ n_x^+ dS = 0 \\
\mathbf{E}_{xx,4} &= -k_e \mathcal{F} \int_{\partial\Omega_4} \vec{N}^T \vec{N} n_x^+ n_x^+ dS = -k_e \mathcal{F} \mathbf{C}_4 \\
\mathbf{E}_{yy,1} &= -k_e \mathcal{F} \int_{\partial\Omega_1} \vec{N}^T \vec{N} n_y^+ n_y^+ dS = -k_e \mathcal{F} \mathbf{C}_1 \\
\mathbf{E}_{yy,2} &= -k_e \mathcal{F} \int_{\partial\Omega_2} \vec{N}^T \vec{N} n_y^+ n_y^+ dS = 0 \\
\mathbf{E}_{yy,3} &= -k_e \mathcal{F} \int_{\partial\Omega_3} \vec{N}^T \vec{N} n_y^+ n_y^+ dS = -k_e \mathcal{F} \mathbf{C}_3 \\
\mathbf{E}_{yy,4} &= -k_e \mathcal{F} \int_{\partial\Omega_4} \vec{N}^T \vec{N} n_y^+ n_y^+ dS = 0 \\
\mathbf{H}_{x,1} &= k_e \int_{\partial\Omega_1} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = 0 \\
\mathbf{H}_{x,2} &= k_e \int_{\partial\Omega_2} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = k_e \left(\frac{1}{2} + \mathcal{C}_x \right) \mathbf{C}_2 \\
\mathbf{H}_{x,3} &= k_e \int_{\partial\Omega_3} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = 0 \\
\mathbf{H}_{x,4} &= k_e \int_{\partial\Omega_4} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = -k_e \left(\frac{1}{2} - \mathcal{C}_x \right) \mathbf{C}_4 \\
\mathbf{H}_{y,1} &= k_e \int_{\partial\Omega_1} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = -k_e \left(\frac{1}{2} - \mathcal{C}_y \right) \mathbf{C}_1 \\
\mathbf{H}_{y,2} &= k_e \int_{\partial\Omega_2} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = 0 \\
\mathbf{H}_{y,3} &= k_e \int_{\partial\Omega_3} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = k_e \left(\frac{1}{2} + \mathcal{C}_y \right) \mathbf{C}_3 \\
\mathbf{H}_{y,4} &= k_e \int_{\partial\Omega_4} \left(\frac{1}{2} + \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = 0 \\
\mathbf{J}_{x,1} &= \int_{\partial\Omega_1} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = 0 \\
\mathbf{J}_{x,2} &= \int_{\partial\Omega_2} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = \left(\frac{1}{2} - \mathcal{C}_x \right) \mathbf{C}_2 \\
\mathbf{J}_{x,3} &= \int_{\partial\Omega_3} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = 0 \\
\mathbf{J}_{x,4} &= \int_{\partial\Omega_4} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_x^+ dS = - \left(\frac{1}{2} + \mathcal{C}_x \right) \mathbf{C}_4 \\
\mathbf{J}_{y,1} &= \int_{\partial\Omega_1} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = - \left(\frac{1}{2} + \mathcal{C}_y \right) \mathbf{C}_1 \\
\mathbf{J}_{y,2} &= \int_{\partial\Omega_2} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = 0 \\
\mathbf{J}_{y,3} &= \int_{\partial\Omega_3} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = \left(\frac{1}{2} - \mathcal{C}_y \right) \mathbf{C}_3 \\
\mathbf{J}_{y,4} &= \int_{\partial\Omega_4} \left(\frac{1}{2} - \vec{C} \cdot \vec{n}^+ \right) \vec{N}^T \vec{N} n_y^+ dS = 0
\end{aligned}$$

8.5.2 The special case of linear triangular elements



The linear basis functions in the triangle are

$$\begin{aligned} N_1(x, y) &= \frac{1}{2S}(x_2y_3 - x_3y_2 + (y_2 - y_3)x + (x_3 - x_2)y) \\ N_2(x, y) &= \frac{1}{2S}(x_3y_1 - x_1y_3 + (y_3 - y_1)x + (x_1 - x_3)y) \\ N_3(x, y) &= \frac{1}{2S}(x_1y_2 - x_2y_1 + (y_1 - y_2)x + (x_2 - x_1)y) \end{aligned}$$

where S is the area of the element:

$$S = \frac{1}{2}[(x_1 - x_3)(y_2 - y_3) - (x_2 - x_3)(y_1 - y_3)]$$

We can easily verify that $N_i(x_j, y_j) = \delta_{ij}$. We then have

$$\begin{aligned} \partial_x N_1(x, y) &= \frac{1}{2S}(y_2 - y_3) \\ \partial_x N_2(x, y) &= \frac{1}{2S}(y_3 - y_1) \\ \partial_x N_3(x, y) &= \frac{1}{2S}(y_1 - y_2) \end{aligned}$$

and

$$\begin{aligned} \partial_y N_1(x, y) &= \frac{1}{2S}(x_3 - x_2) \\ \partial_y N_2(x, y) &= \frac{1}{2S}(x_1 - x_3) \\ \partial_y N_3(x, y) &= \frac{1}{2S}(x_2 - x_1) \end{aligned}$$

Then, as shown in Section E.0.5, the mass matrix³ and the \mathbf{J}_x and \mathbf{J}_y matrices are:

$$\mathbf{E} = \int_{\Omega} \vec{N}^T \vec{N} dV = \frac{S}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} \quad (8.149)$$

$$\begin{aligned} \mathbf{J}_x &= - \int_{\Omega} \partial_x \vec{N}^T \vec{N} dV = -\frac{1}{6} \begin{pmatrix} y_2 - y_3 & y_2 - y_3 & y_2 - y_3 \\ y_3 - y_1 & y_3 - y_1 & y_3 - y_1 \\ y_1 - y_2 & y_1 - y_2 & y_1 - y_2 \end{pmatrix} \\ \mathbf{J}_y &= - \int_{\Omega} \partial_y \vec{N}^T \vec{N} dV = -\frac{1}{6} \begin{pmatrix} x_3 - x_2 & x_3 - x_2 & x_3 - x_2 \\ x_1 - x_3 & x_1 - x_3 & x_1 - x_3 \\ x_2 - x_1 & x_2 - x_1 & x_2 - x_1 \end{pmatrix} \end{aligned} \quad (8.150)$$

³The mass matrix is commonly called \mathbf{M} but I use here the same notations as in Li's book.

In the same appendix we show that

$$\mathbf{C}_1 = \int_{\partial\Omega_3} \vec{N}^T \vec{N} dS = \frac{L_1}{6} \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (8.151)$$

$$\mathbf{C}_2 = \int_{\partial\Omega_1} \vec{N}^T \vec{N} dS = \frac{L_2}{6} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad (8.152)$$

$$\mathbf{C}_3 = \int_{\partial\Omega_2} \vec{N}^T \vec{N} dS = \frac{L_3}{6} \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 2 \end{pmatrix} \quad (8.153)$$

Testing the waters - constant temperature field

Let us assume that the temperature is constant in space. It then follows that the heat flux is identically zero.

Testing the waters - linear temperature field

If the temperature field is given by $T(x, y) = T_0 - ax - by$ then $q_x = a$ and $q_y = b$.

8.6 Time-dependent diffusion PDE in 2D

8.7 Stokes equations

Two relevant papers:

- Cockburn *et al.* (2002) [265] - LDG
- Cockburn *et al.* (2010) [268] - HDG

Let us start with the dimensionless Stokes system [265]:

$$-\eta\Delta\vec{v} + \vec{\nabla}p = \vec{f} \quad \text{in } \Omega \quad (8.154)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad \text{in } \Omega \quad (8.155)$$

$$\vec{v} = \vec{v}_D \quad \text{on } \Gamma \quad (8.156)$$

where Ω is a bounded domain of \mathbb{R}^d and the Dirichlet boundary conditions are such that they satisfy the compatibility condition

$$\int_{\Gamma} \vec{v}_D \cdot \vec{n} = 0$$

where \vec{n} is the outward unit normal.

Gradient-based formulation In order to obtain the LDG methods we first rewrite this system as the following collection of conservation laws [265]:

$$\mathbf{L} = \vec{\nabla}\vec{v} \quad \text{in } \Omega \quad (8.157)$$

$$\vec{\nabla} \cdot (-2\eta\mathbf{L} + p\mathbf{1}) = \vec{f} \quad \text{in } \Omega \quad (8.158)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad \text{in } \Omega \quad (8.159)$$

$$\vec{v} = \vec{v}_D \quad \text{on } \Gamma \quad (8.160)$$

supplemented by

$$\int_{\Omega} p = 0$$

where \mathbf{L} is the gradient tensor, $\mathbf{1}$ is the unit tensor.

Remark. *It may appear counter-intuitive at first to define \mathbf{L} as being the gradient of the velocity instead of the strain rate tensor but under the assumption of incompressibility $\partial_x u + \partial_y v = 0$ (and constant viscosity) we can write:*

$$\vec{\nabla} \cdot (2\eta\mathbf{L}) = 2\eta \begin{pmatrix} \partial_x^2 u + \frac{1}{2}\partial_x\partial_y v + \frac{1}{2}\partial_y^2 u \\ \frac{1}{2}\partial_x^2 v + \frac{1}{2}\partial_y\partial_x u + \partial_y^2 v \end{pmatrix} = 2\eta \begin{pmatrix} \partial_x^2 u + \frac{1}{2}\partial_x(-\partial_x u) + \frac{1}{2}\partial_y^2 u \\ \frac{1}{2}\partial_x^2 v + \frac{1}{2}\partial_y(-\partial_y v) + \partial_y^2 v \end{pmatrix} = \eta \begin{pmatrix} \partial_x^2 u + \partial_y^2 u \\ \partial_x^2 v + \partial_y^2 v \end{pmatrix} = \eta\Delta\vec{v}$$

Remark. *Cockburn et al. (2010) [268] also introduce the vorticity-based formulation and the stress-based formulation but we will not explore these in what follows.*

RETYPE section 2.1 of [265]

Chapter 9

Additional techniques, features, measurements

chapter8.tex

Solving the Stokes equations and the energy equations is one thing. Doing it in a geodynamical context requires a lot of additional techniques.

9.1 Dealing with a free surface (and mesh deformation)

When carrying out global models, typically mantle convection, the effect of the free surface is often neglected/negligible: topography ranges from $\sim 10\text{km}$ depth to $\sim 10\text{km}$ height, which is very small compared to the depth of the mantle ($\sim 3000\text{km}$).

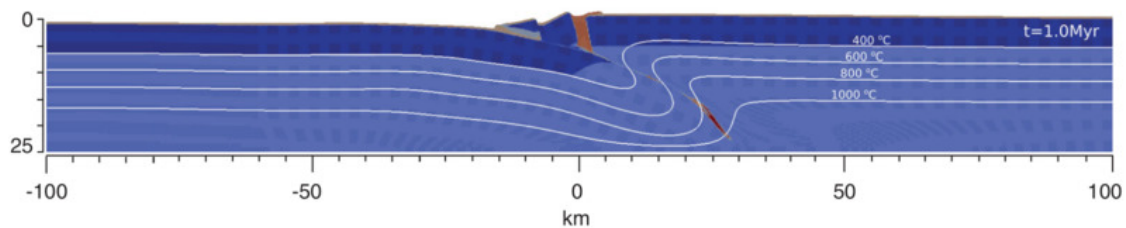
However, it has long been recognised that there is a feedback between topography and crust/lithosphere deformation: the surface of the Earth reflects the deeper processes, from orogeny, back-arc basins, rifts, mid-ocean ridges, etc ... (see for instance Braun (2010) [138]).

Remark. *Free surface flows are not unique to Earth sciences, and their modelling has given rise to many studies and textbooks. A typical free-surface flow problem in the CFD literature is the so-called 'dam break' problem [888, 42, 788, 757, 584, 17]. Other occurrences involve sea waves, flow over structures, flow around ships, mould filling, flow with bubbles [788].*

Remark. *Free surface flows have also been extensively studied and even benchmarked in the ice-sheet modelling community, see [XXX].*

What distinguishes geodynamics free surface modelling from its engineering counterpart is (i) the absence of surface tension, (ii) the fact that the fluids under consideration are Stokesian, (iii) their rheology is complex (the elastic and plastic components can be dominant at the surface).

The problem of dealing with a free surface can be deceptively simple at first glance: as mentioned before the amplitude of surface movement is often less than 1% of the domain size. Isostasy-driven movements are easy to deal with since the movement is vertical (and often characterised by long wavelength). However, computational problems quickly arise for example in subduction modelling: the downgoing lithosphere subducts below the overriding plate and the relative convergence of the two is likely to generate a cusp at the trench. The presence of shear bands intersecting the surface accentuates the problem:

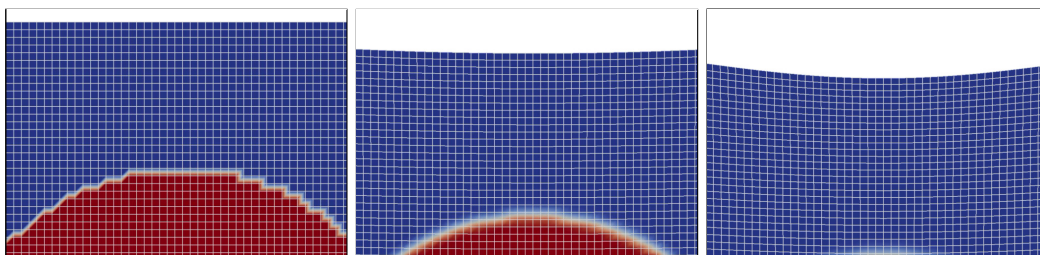


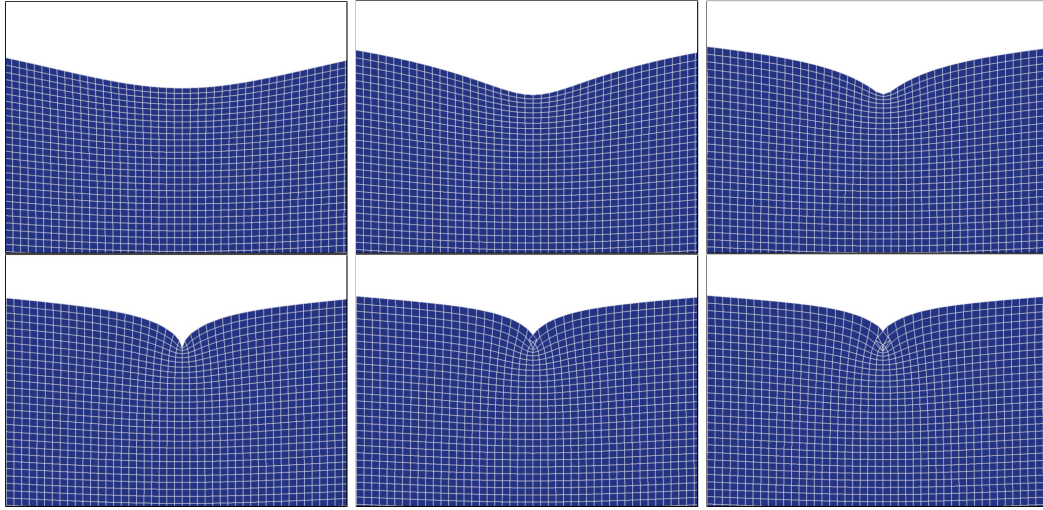
Taken from Maffione *et al.* [823]. Example of free surface deformation above intra-oceanic subduction initiation

Remark. *It is difficult to talk about free surface without including the underlying mesh. What follows should be read alongside Section 9.6.*

The fully Lagrangian approach

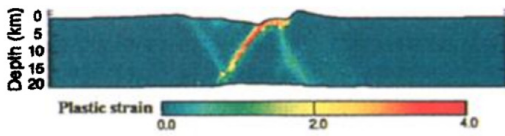
In this case the mesh is deformed with the velocity (or displacement) computed on its nodes. It is sometimes called 'body fitting' [285] or 'boundary fitted'. In the case when large deformation occurs (which is rather frequent in geodynamics – think about subduction or rifting processes where materials end up moving 100's or 1000's of km, horizontally and/or vertically) – it leads to highly deformed elements, and in some case even bow-tied:



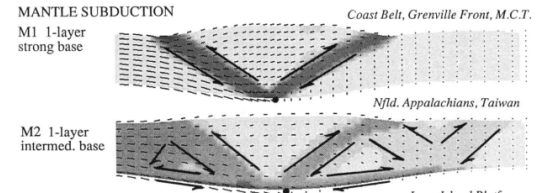


Example of a free surface evolution above a sinking sphere. The isostatic rebound above the sphere generates a cusp which, if no special measure is taken, ultimately leads to a bow-tied element. Once this occurs the simulation stops since the mapping of the bow-tied element to the reference element yields to wrong elemental matrix. Courtesy of M. Fraters

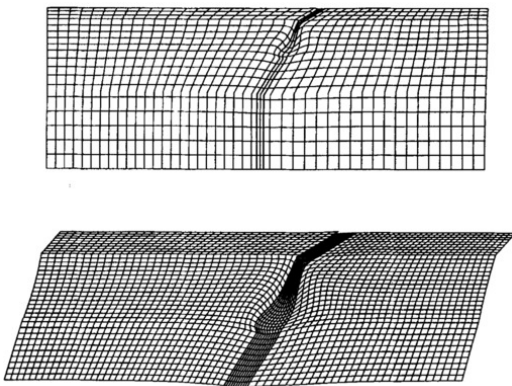
In the mildest cases this does not occur but it has long been established that large mesh deformation yields low accuracy calculations, especially when angles between edges become small or large. One way to overcome this problem is to remesh, i.e. generate a better mesh based on the available information on the deformed one. In 2D this is routinely done, especially when triangular elements are used. In 3D, multiple remeshing are very costly and it is generally avoided. Note also that re-meshing often involves some form of interpolation and therefore some unwanted numerical diffusion. When deformation is reasonably small, fully lagrangian methods work and have been used in geodynamics, see for example Hassani & Chéry (1996) [553], Melosh & Raefsky (1980) [862], or Lavier *et al.* (2000) [754].



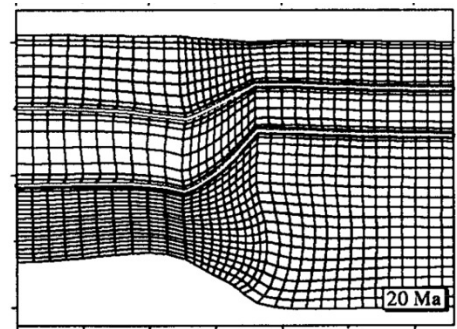
Taken from Lavier *et al.* (2000) [754]. Upper-crustal faulting, note that the bottom and the top surface are deformed.



Taken from Beaumon *et al.* (1994) [60]. Strain rate and velocity field for crustal S-point models.



Taken from Gurnis *et al.* (1996) [516]. Subduction model, topographic expression is shown without vertical exaggeration.



Taken from Govers & Wortel [479]. Asymmetric lithospheric extension.

The Eulerian approach: using sticky air

'Sticky air' is the default option for numerical methods which mesh cannot be deformed (typically the finite difference method). In this case, the air above the crust/sediments is modelled as a zero-

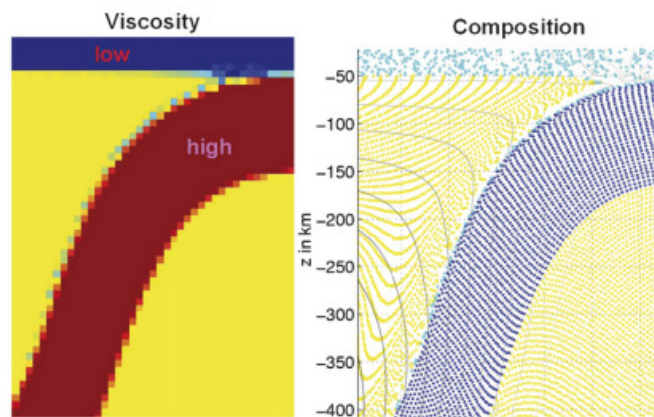
density¹ fluid with very low viscosity (see for instance the early article by Zaleski and Julien [1400]). One problem quickly arises when one realises that the viscosity of the air ($\sim 18.5 \cdot 10^{-6} \text{ Pa}\cdot\text{s}$ ²) is almost 25-30 orders of magnitude lower than the (effective) viscosity of Earth materials. Real air viscosity cannot therefore be used because of 1) round-off errors, 2) extremely poorly-conditioned matrices. Low viscosities around $10^{16} - 10^{19} \text{ Pa}\cdot\text{s}$ are then commonly used as they are still negligible next to those of the (plastic) crust, and the flow of air parallel to Earth materials only generates extremely small shear and normal stress values (thereby approaching the true nature of a free surface). This approach is the one employed in all the papers based on the I2/I3(EL)VIS code (see Appendix ??) and has been benchmarked in Cramer *et al.* (2012) [285].

This approach has a few advantages:

1. it is simple to implement
2. it is compatible with all the standard numerical methods (FEM, FDM, FVM)
3. it avoids (potentially complicated or costly) remeshing

and quite a few drawbacks:

1. it increases the size of the computational domain, thereby adding more unknowns to the linear system: in Schmeling *et al.* (2008) [1124] the air layer is set to 50 km while the lithospheric domain underneath is 700 km thick;
2. it requires the use of averaging all along the free-surface where very large viscosity contrasts are present. Here is what Poliakov and Podlachikov [1008] say about the sticky air method: "Zaleski & Julien [1400] used a top layer with a very low viscosity and density to represent air or water above the surface. This allows a simple representation of the free surface. However, due to the very high viscosity and density contrast and diffusion between the top layer and the underlying layers, calculations sometimes become unstable and give significant errors."
3. it can yield air entrainment in the mantle:



Taken from Schmeling *et al.* (2008) [1124]. Details of the entrainment and lubrication of the soft surface layer. Light blue particles are sticky air particle and are found to greatly alter the viscosity of the subduction channel.

4. it is not clear how thick the air layer must be
5. it often requires to ascribe thermal parameters to the air;

¹Sometimes a zero value is problematic and one then resorts to a very low value instead.

²<https://en.wikipedia.org/wiki/Viscosity>

6. it makes the implementation of Dirichlet or Neuman boundary conditions for temperature at the surface less obvious.
7. it makes the coupling with surface processes codes less straightforward.
8. its accuracy depends on the method used to track materials in the rest of the code (markers, level sets, ...). If markers are used, the free surface position is then known up to the average distance between markers.
9. it negatively impacts the condition number of the matrix.

The term 'sticky water' is sometimes employed too. The dynamic viscosity of water is about 10^{-3} Pas so that it is also negligible compared to the viscosity of Earth materials and the same reasoning as air applies. However, in such a case a density of about 1000 kg m^{-3} is then assigned to the layer (e.g. Gerya & Burg (2007) [451]).

In conclusion, as stated in Cramer *et al.* (2012) [285]: "the sticky air method is a good way to simulate a free surface for Eulerian approaches, provided that its parameters are chosen carefully."

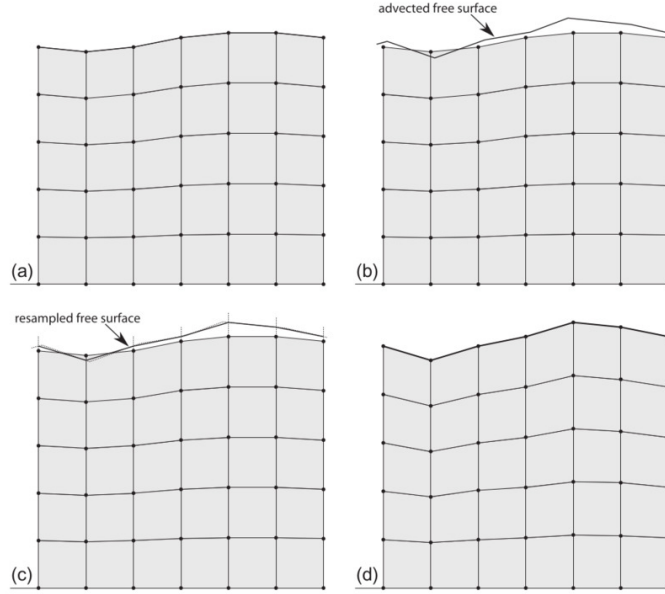
The Arbitrary Lagrangian Eulerian (ALE) approach

It is a very widely used approach in FEM-based geodynamics codes but originates in the field of CFD (Hirt *et al.* (1974)[575], Hughes *et al.* (1981) [609]) and is described at length in Souli & Zolesio (2001) [1182], Donea *et al.* (2004) [339], Donea & Huerta [341]. To put it very simply, the key idea in the ALE formulation is the introduction of a computational mesh which can move and deform with a velocity independent of the velocity carried by the material particles.

The simple approach in Thieulot (2011) [1258]. What follows is written with a 2D Cartesian model in mind ($Q_1 \times P_0$ elements are used). The computational domain is a rectangle of size $L_x \times L_y$ and a $nnx \times nny$ rectangular grid spanning the simulation domain is generated. The grid points constituting the top row of the grid define the discrete free surface of the domain. Once the Eulerian velocity field has been computed on these, their position is first updated using a simple Eulerian advection step (see a,b on figure hereunder):

$$\vec{r}'_i(t + \delta t) = \vec{r}_i(t) + \vec{v}_i \cdot \delta t \quad i = 1, \dots, nnx$$

The other boundaries of the system remain fixed at locations $x = 0$, $x = L_x$ and $y = 0$. Even though the Eulerian grid must conform to the current domain shape, only vertical motion of grid nodes is allowed. It is therefore necessary to resample the predicted free surface given by \vec{r}'_i at equidistant positions between $x = 0$ and $x = L_x$. The resampling is carried out either with Spline functions or a moving least square algorithm. Finally, the vertical position of all the nodes corresponding to column $i \in [1, nnx]$ is recalculated so that they are equidistant, as sketched in Figure d. This has the advantage of keeping the mesh distortion to a minimum in the case of large deformation.



The ALE algorithm of [1258] in 2D. (a) Grid and free surface at a given time t ; (b) advection of the free surface; (c) resampling of the free surface at equidistant abscissae; (d) vertical adjustment of grid nodes in each column at equidistant ordinates.

The ALE method is used in the SOPALE , SULEC , FANTOM , ELEFANT , and ASPECT codes to name a few (see Appendix ??) although each code has its own specific implementation, as detailed in what follows for ASPECT .

The not-so-simple but rather elegant approach of Aspect What follows is mostly borrowed from Rose *et al.* [1085]. Their approach has the advantage that it does not presuppose a geometry (Cartesian, Spherical, ...) nor a number of dimensions. It is also designed to work in parallel and on octree-based meshes, and with various combinations of boundary conditions. Note that the authors specify that "for moderate mesh deformation, the mesh stays smooth and well conditioned, though it breaks down for large deformations".

This approach is obtained by simply imposing the obvious condition that no particle (fluid parcel) can cross the free surface (because it is a material surface). This can be imposed in a straightforward manner by using a Lagrangian description along this surface. However, this condition may be relaxed by imposing only the necessary condition: \vec{v} equal to zero along the normal to the boundary (ie. $\vec{n} \cdot \vec{v} = 0$, where \vec{n} is the outward unit normal to the fluid domain). The mesh position, normal to the free surface, is determined from the normal component of the particle velocity and remeshing can be performed along the tangent; see, for instance Huerta and Liu (1988) [603] or Braess and Wriggers (2000) [130].

As mentioned above, the mesh velocity in normal direction at the free surface (with unit normal \vec{n}) has to be consistent with the velocity of the Stokes velocity solution $\vec{v}(t)$:

$$\vec{v}_{\text{mesh}}(t) \cdot \vec{n} = \vec{v}(t) \cdot \vec{n} \quad \text{on } \Gamma_F \quad (9.1)$$

In ALE calculations the internal mesh velocity is usually undetermined, but one wants to smoothly deform the mesh so as to preserve its regularity, avoiding inverted or otherwise poorly conditioned cells. The mesh deformation can be calculated in many different ways, including algebraic (as mentioned in the previous paragraph) and PDE based approaches. The latter is chosen here. The Laplace equation is solved where the unknown is the mesh velocity, i.e. one must solve:

$$\Delta \vec{v}_{\text{mesh}} = 0 \quad (9.2)$$

subjected to the following boundary conditions:

$$\begin{aligned}\vec{v}_{\text{mesh}} &= \vec{0} & \text{on } \Gamma_0 \\ \vec{v}_{\text{mesh}} &= (\vec{v} \cdot \vec{n})\vec{n} & \text{on } \Gamma_F \\ \vec{v}_{\text{mesh}} \cdot \vec{n} &= 0 & \text{on } \Gamma_{FS}\end{aligned}\tag{9.3}$$

where Γ_{FS} is the part of the boundary with free slip boundary conditions, Γ_0 is the no-slip part and Γ_{FS} is the free slip part.

Once the mesh velocity has been obtained for all mesh points, these can be moved with said velocity. However, it must be noted that the multiple occurrences of the normal vector in the above equations is not without problem as the normal vectors are not well defined on the mesh vertices, which is where the mesh velocity is defined.

INSERT FIGURE

This yields what the authors coin the 'quasi-implicit' scheme (we have so far neglected any kind of stabilisation):

1. Solve the Stokes system;
2. Solve for the surface mesh velocity using Equation 9.4;
3. Solve for the internal mesh velocity using Equations 9.2, 9.3;
4. Advect the mesh forward in time using displacements determined by the forward Euler scheme:
 $\vec{x}(t^{n+1}) = \vec{x}(t^n) + \vec{v}_{\text{mesh}}\delta t.$

The authors list two simple methods of computing the normals:

- one can take \vec{n} as the direction of the local vertical,
- one could compute \vec{n} as some weighted average of the cell normals adjacent to a given vertex

but conclude that they have found that these schemes do not necessarily have good mass conservation properties.

A better approach is proposed in the form of an L_2 projection of the normal velocity $\vec{v} \cdot \vec{n}$ onto the free surface Γ_F . Multiplying the boundary conditions $\vec{v}_{\text{mesh}} = (\vec{v} \cdot \vec{n})\vec{n}$ by a test function \vec{w} and integrating over the free surface part of the boundary, we find:

$$\int_{\Gamma_F} \vec{w} \cdot \vec{v}_{\text{mesh}} d\Gamma = \int_{\Gamma_F} \vec{w} \cdot (\vec{v} \cdot \vec{n})\vec{n} d\Gamma = \int_{\Gamma_F} (\vec{w} \cdot \vec{n})(\vec{v} \cdot \vec{n}) d\Gamma\tag{9.4}$$

When discretized, this forms a linear system which can be solved for the mesh velocity \vec{v}_{mesh} at the free surface. This system, being nonzero over only the free surface, is relatively computationally inexpensive to solve. The authors unfortunately fail to mention that this approach is particularly interesting since the numerical quadrature used to compute the above integrals require the normal \vec{n} between the nodes and these normals are well defined over each segment joining two nodes!³

In what follows I present in some detail how to carry out the L_2 projection to arrive at the surface velocity for both Q_1 and Q_2 elements.

I start from the following integral over a Q_1 element:

$$\int_{\Gamma_e} N_i \vec{v}_{\text{mesh}} d\Gamma = \int N_i \begin{pmatrix} u_{\text{mesh}} \\ v_{\text{mesh}} \end{pmatrix} d\Gamma\tag{9.5}$$

$$= \int N_i \begin{pmatrix} N_1 & 0 & N_2 & 0 \\ 0 & N_1 & 0 & N_2 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \end{pmatrix} d\Gamma\tag{9.6}$$

$$\tag{9.7}$$

³what if Q_k with $k > 1$ elements are used and nodes on the surface no more form a line?

Writing this equation alternatively for $N_i = N_1, N_2$ yields:

$$\int_{\Gamma_e} \begin{pmatrix} N_1 N_1 & 0 & N_1 N_2 & 0 \\ 0 & N_1 N_1 & 0 & N_1 N_2 \\ N_2 N_1 & 0 & N_2 N_2 & 0 \\ 0 & N_2 N_1 & 0 & N_2 N_2 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \end{pmatrix} d\Gamma = \int_{\Gamma_e} \begin{pmatrix} N_1 N_1 & 0 & N_1 N_2 & 0 \\ 0 & N_1 N_1 & 0 & N_1 N_2 \\ N_2 N_1 & 0 & N_2 N_2 & 0 \\ 0 & N_2 N_1 & 0 & N_2 N_2 \end{pmatrix} d\Gamma \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \end{pmatrix}$$

Turning now to the right hand side $\int_{\Gamma_e} N_i (\vec{\nu} \cdot \vec{n}_e) \vec{n}_e d\Gamma$, it yields the following rhs:

$$\int_{\Gamma_e} (\vec{\nu} \cdot \vec{n}_e) \begin{pmatrix} N_1 n_x \\ N_1 n_y \\ N_2 n_x \\ N_2 n_y \end{pmatrix} d\Gamma$$

The elemental matrix and rhs must be built for each element and assembled in a global matrix and rhs. The solution is the mesh velocity vector at all surface nodes. the same approach can be taken for Q_2 elements:

$$\int_{\Gamma_e} N_i \vec{\nu}_{mesh} d\Gamma = \int N_i \begin{pmatrix} u_{mesh} \\ v_{mesh} \end{pmatrix} d\Gamma \quad (9.8)$$

$$= \int N_i \begin{pmatrix} N_1 & 0 & N_2 & 0 & N_3 & 0 \\ 0 & N_1 & 0 & N_2 & 0 & N_3 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \end{pmatrix} d\Gamma \quad (9.9)$$

Writing this equation alternatively for $N_i = N_1, N_2, N_3$ yields:

$$\begin{aligned} & \int_{\Gamma_e} \begin{pmatrix} N_1 N_1 & 0 & N_1 N_2 & 0 & N_1 N_3 & 0 \\ 0 & N_1 N_1 & 0 & N_1 N_2 & 0 & N_1 N_3 \\ N_2 N_1 & 0 & N_2 N_2 & 0 & N_2 N_3 & 0 \\ 0 & N_2 N_1 & 0 & N_2 N_2 & 0 & N_2 N_3 \\ N_3 N_1 & 0 & N_3 N_2 & 0 & N_3 N_3 & 0 \\ 0 & N_3 N_1 & 0 & N_3 N_2 & 0 & N_3 N_3 \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \end{pmatrix} d\Gamma \\ &= \int_{\Gamma_e} \begin{pmatrix} N_1 N_1 & 0 & N_1 N_2 & 0 & N_1 N_3 & 0 \\ 0 & N_1 N_1 & 0 & N_1 N_2 & 0 & N_1 N_3 \\ N_2 N_1 & 0 & N_2 N_2 & 0 & N_2 N_3 & 0 \\ 0 & N_2 N_1 & 0 & N_2 N_2 & 0 & N_2 N_3 \\ N_3 N_1 & 0 & N_3 N_2 & 0 & N_3 N_3 & 0 \\ 0 & N_3 N_1 & 0 & N_3 N_2 & 0 & N_3 N_3 \end{pmatrix} d\Gamma \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \end{pmatrix} \end{aligned} \quad (9.10)$$

The right hand side is then

$$\int_{\Gamma_e} (\vec{\nu} \cdot \vec{n}_e) \begin{pmatrix} N_1 n_x \\ N_1 n_y \\ N_2 n_x \\ N_2 n_y \\ N_3 n_x \\ N_3 n_y \end{pmatrix} d\Gamma$$

Having obtained the boundary condition velocity for the Laplace equation, we can now turn our attention to solving this ODE.

In what follows I omit the subscript 'mesh' and focus on the 2D case. The components of the (mesh) velocity are given by

$$u^h = \sum_{i=1}^{m_v} N_i^v u_i \quad v^h = \sum_{i=1}^{m_v} N_i^v v_i \quad \vec{v}^h = \begin{pmatrix} u^h \\ v^h \end{pmatrix}$$

We start from the ODE to solve in its strong form:

$$\Delta \vec{v}^h = \vec{0}$$

We multiply it by a velocity test function N_i^v and integrate over an element:

$$\begin{aligned} & \vec{0} \\ &= \int_{\Omega_e} N_i^v \Delta \vec{v}^h \\ &= \int_{\Omega_e} N_i^v \Delta \vec{v}^h dV \\ &= \int_{\Omega_e} \begin{pmatrix} N_i^v \Delta u^h \\ N_i^v \Delta v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} N_i^v \vec{\nabla} \cdot \vec{\nabla} u^h \\ N_i^v \vec{\nabla} \cdot \vec{\nabla} v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \vec{\nabla} N_i^v \cdot \vec{\nabla} u^h \\ \vec{\nabla} N_i^v \cdot \vec{\nabla} v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \partial_x N_i^v \partial_x u^h + \partial_y N_i^v \partial_y u^h \\ \partial_x N_i^v \partial_x v^h + \partial_y N_i^v \partial_y v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \partial_x N_i^v & \partial_y N_i^v & 0 & 0 \\ 0 & 0 & \partial_x N_i^v & \partial_y N_i^v \end{pmatrix} \cdot \begin{pmatrix} \partial_x u^h \\ \partial_y u^h \\ \partial_x v^h \\ \partial_y v^h \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \frac{\partial N_1^v}{\partial x} & \frac{\partial N_1^v}{\partial y} & 0 & 0 \\ 0 & 0 & \frac{\partial N_1^v}{\partial x} & \frac{\partial N_1^v}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial N_1^v}{\partial x} & 0 & \frac{\partial N_2^v}{\partial x} & 0 & \dots & \frac{\partial N_{m_v}^v}{\partial x} & 0 \\ \frac{\partial N_1^v}{\partial y} & 0 & \frac{\partial N_2^v}{\partial y} & 0 & \dots & \frac{\partial N_{m_v}^v}{\partial y} & 0 \\ 0 & \frac{\partial N_1^v}{\partial x} & 0 & \frac{\partial N_2^v}{\partial x} & \dots & 0 & \frac{\partial N_{m_v}^v}{\partial x} \\ 0 & \frac{\partial N_1^v}{\partial y} & 0 & \frac{\partial N_2^v}{\partial y} & \dots & 0 & \frac{\partial N_{m_v}^v}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix} dV \end{aligned}$$

Writing this equation for $i = 1, \dots, m_v$, we obtain:

$$\int \begin{pmatrix} \frac{\partial N_1^v}{\partial x} & \frac{\partial N_1^v}{\partial y} & 0 & 0 \\ 0 & 0 & \frac{\partial N_1^v}{\partial x} & \frac{\partial N_1^v}{\partial y} \\ \frac{\partial N_2^v}{\partial x} & \frac{\partial N_2^v}{\partial y} & 0 & 0 \\ 0 & 0 & \frac{\partial N_2^v}{\partial x} & \frac{\partial N_2^v}{\partial y} \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial N_{m_v}^v}{\partial x} & \frac{\partial N_{m_v}^v}{\partial y} & 0 & 0 \\ 0 & 0 & \frac{\partial N_{m_v}^v}{\partial x} & \frac{\partial N_{m_v}^v}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial N_1^v}{\partial x} & 0 & \frac{\partial N_2^v}{\partial x} & 0 & \dots & \frac{\partial N_{m_v}^v}{\partial x} & 0 \\ \frac{\partial N_1^v}{\partial y} & 0 & \frac{\partial N_2^v}{\partial y} & 0 & \dots & \frac{\partial N_{m_v}^v}{\partial y} & 0 \\ 0 & \frac{\partial N_1^v}{\partial x} & 0 & \frac{\partial N_2^v}{\partial x} & \dots & 0 & \frac{\partial N_{m_v}^v}{\partial x} \\ 0 & \frac{\partial N_1^v}{\partial y} & 0 & \frac{\partial N_2^v}{\partial y} & \dots & 0 & \frac{\partial N_{m_v}^v}{\partial y} \end{pmatrix} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix}}_{\vec{v}} dV = \vec{0}$$

or,

$$\left(\int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{B} \, dV \right) \cdot \vec{V} = \vec{0}$$

where \mathbf{B} is a $(ndim * ndim) \times (m_v * ndofV)$ matrix. This is implemented in Stone 54 ??.

Remark. *The integration by parts should have a minus appear but since the left hand side is 0, it is not taken into account.*

Note that Rose *et al.* (2017) [1085] go further than this and propose a 'nonstandard finite difference scheme' and make a link with the stabilisation presented in Kaus *et al.* (2010) [681].

surface terms arising from the integration by parts are neglected. EXPLAIN WHY!

Yet another approach [339] The unknown position of free surfaces can be computed using the following approach: for the simple case of a single-valued function $h = h(x, y, t)$, a hyperbolic equation must be solved,

$$\frac{\partial h}{\partial t} + (\vec{v} \cdot \vec{\nabla})h = 0 \quad (9.11)$$

This is the kinematic equation of the surface and has been used, for instance, by Ramaswamy and Kawahara (1987), Huerta and Liu, 1988b, 1990; Souli and Zolesio (2001).

Eq. (9.11) is a simple advection equation. One could also add a diffusion operator with a diffusion coefficient D . Low values of D could be used to stabilise the surface while higher values (possibly nonlinear ones) could be used to account for simple surface processes.

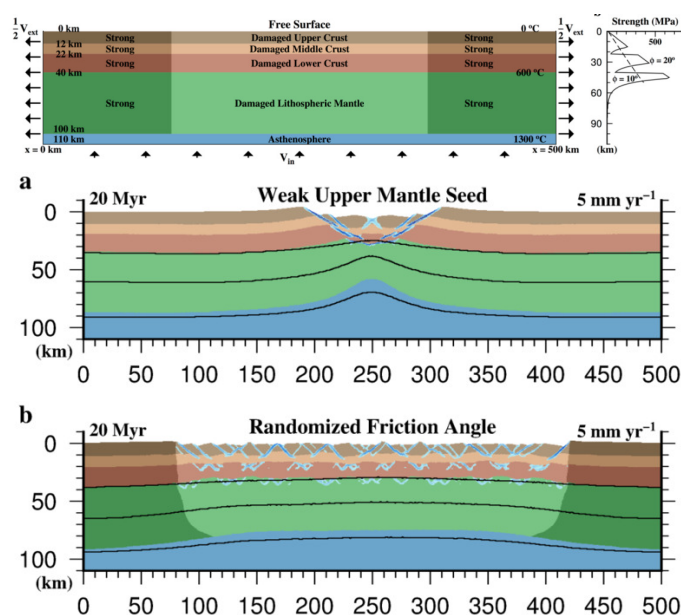
$$\frac{\partial h}{\partial t} + (\vec{v} \cdot \vec{\nabla})h = D\Delta h \quad (9.12)$$

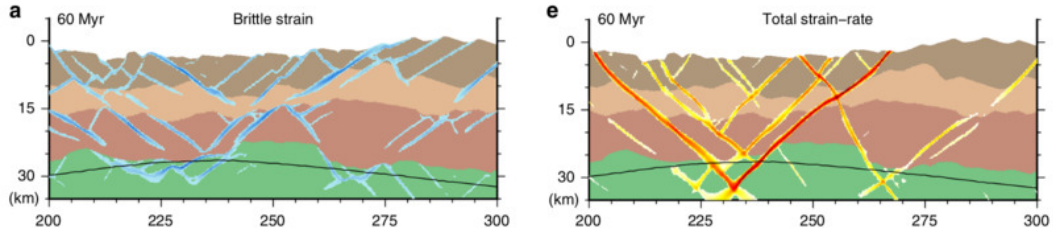
Also, Hansen and Nielsen [532, 531] write: “During the entire model evolution surface processes act to re-distribute sediments. These processes are modelled by a diffusion equation with a source term enabling the transport of sediments to and from the model profile. The transport equation is written

$$\dot{h} = \nabla \cdot (\kappa \nabla h) + \dot{s}(w)$$

where $\kappa = 200 km^2/Ma$ is the diffusivity of topography and $\dot{s}(w)$ is a linear function of water depth.”

The following pictures are taken from Naliboff *et al.* (2017) [927] on the topic of how complex fault interaction controls continental rifting. It is a beautiful example (among many) of the importance of free surface geodynamical expression and large deformation:





Taken from [927]

On the topic of moving internal nodes

Braess & Wriggers [130] propose the following interesting algorithm: "A measure of the quality of a triangular mesh is the quotient of the outer radius r_{out} and the inner radius r_{in} of each element. This quotient is important because it plays a certain role in a priori error estimates. If an element degenerates this quotient will approach infinity. Another important feature of good mesh is that no element becomes very large. With these considerations in mind the penalty function W is defined:

$$W = \sum_{elts} \left(\frac{r_{out}}{r_{in}} \right)^m \left(\frac{r_{out}}{r_0} \right)^n \quad (9.13)$$

where m , n and r_0 are positive constants. For our calculations we chose $m = 3$, $n = 1$ and $r_0 = 1$, but the results seem to depend only slightly on this choice. Whenever a triangle is distorted or very large, this function becomes very large. A similar penalty function was presented in [657] for four-node elements. In that case the angles of the elements are used to construct the penalty function. In order to regularize a distorted mesh the coordinates of the internal nodes will be chosen such that W is minimized. It is not necessary to reach the global minimum, a rough approximation is sufficient. Therefore the minimization of the potential can be done efficiently with standard procedures and will not be discussed in any detail. This algorithm can also be applied to h -adaptive mesh-generation by choosing appropriate constants r_0 for each triangle."⁴

This is still WORK IN PROGRESS. I Need to look at those papers:

ALE: Ramaswamy & Kawahara (1987) [1034], ALE: Huerta & Liu (1988) [603], ALE: Ponthot & Belytschko (1998)[1010], ALE: Benson (1989) [71], ALE: Sung *et al.* (2000) [1219], ALE: Ramaswamy (1990) [1033], Tezduyar *et al.* (1992) [1248](moving pulse) Andres-Martinez *et al.* (2015) [24] (read !! extract benchmarks ?) Kramer *et al.* (2012) [731] (read !! extract benchmarks ?) Steer *et al.* (2011) [1195] Maierova [825] Zhong *et al.* [1410], Gurnis *et al.* [516] Ellis *et al.* (2004) [369] Braess & Wriggers (2000) [130] Parsons & Daly (1983) [980]

Free surface treatment in FDM: [353].

Free surface stabilization algorithm (FSSA)

To start with, Duretz *et al.* (2011) [352] is concerned with the Finite Difference method so I will not expose this one further (see Gerya's book [456] for all things FDM). Note that their approach is similar to what follows.

⁴Indeed, if r_0 is the same for all elements this parameter will not play any role at all in the minimisation process.

There are then four articles left: Kaus *et al.* (2010) [681], Quinquis & Buiter (2011) [1028], Rose *et al.* (2017) [1085], and Schuh-Senlis *et al.* (2020) [1144].

The first one is rather technical while the second is more to the point with some simple reasoning. Apparently both author groups arrived at about the same formulation/algorithm at about the same time. The third and fourth paper implement the algorithm in ASPECT and the FAIStokes code respectively.

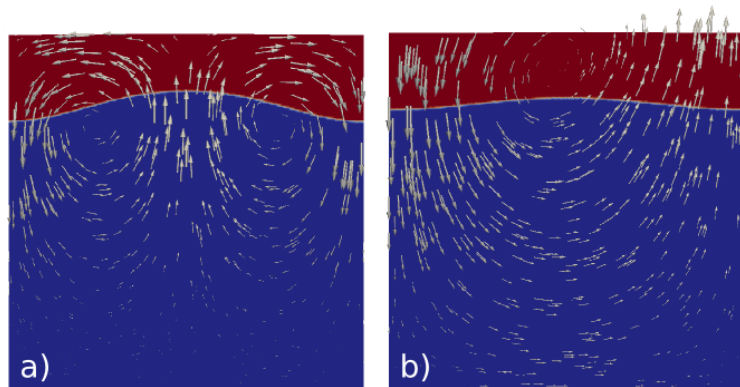
The basic idea is rather simple. As stated in Quinquis *et al.* (2011) [1028]: "In numerical subduction models which include a free surface or other interfaces with large density contrasts, an instability can occur as a result of numerical overshoot when computing restoring forces".

In essence when the timestep δt exceeds (or is comparable) to a characteristic relaxation time t_s , then body and surface forces will be out of balance with the updated free surface. As a consequence oscillations will occur and get amplified with sometimes a sloshing effect, also called 'drunken sailor'. The solution is to add a term in the FE matrix which takes into account the incremental change in normal forces across density interfaces during each timestep:

$$\Delta F_y = \Delta \rho g_y \delta t v_y$$

where ΔF_y is the extra vertical force term across the interface, $\Delta \rho$ is the change in density across the interface, and v_y is the vertical velocity on the interface. This extra term is applied to the free surface, and at any other interface with a density contrast. Unfortunately, although this is easy to understand, it is not very clear how to implement it, which is when Kaus *et al.* (2010) becomes useful.

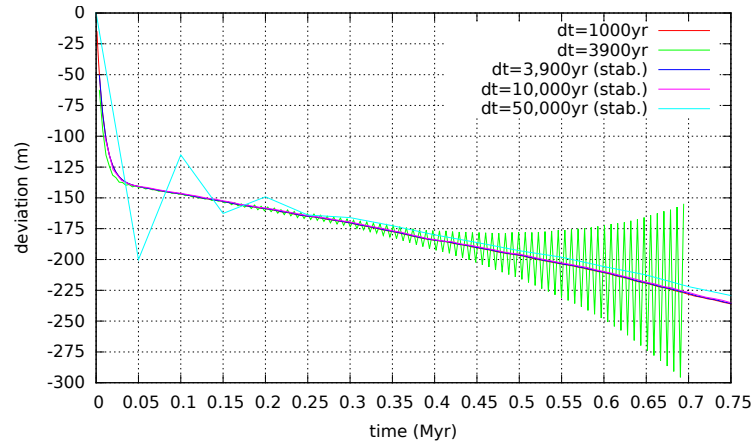
I first present results from my doomed 2014 ELEFANT paper ⁵ [1257] in which I carried out the Rayleigh-Taylor instability experiment described in [681] where a dense fluid overlays another fluid in a square domain with their initial interface being of sinusoidal shape.



Taken from Thieulot (2014) [1257] When the time step δt is small, the simulation evolves smoothly, but too large a time step leads to a sloshing instability, also called 'drunken sailor effect' in the community.

The position $y(t)$ of the free surface point situated at $x = L_x$ is monitored over time:

⁵This paper is about my ELEFANT code which relied at the time solely on $Q_1 \times P_0$ elements with the penalty formulation.



The stabilisation algorithm is first switched off and up to a time step of $\delta t \sim 3800\text{yr}$, the free surface does not develop any instability (all curves are superimposed). For a time step of 3900yr and above, a clear oscillation develops (green seesaw curve on the figure). Turning the algorithm on, time steps larger than $dt \sim 10,000\text{yr}$ can be used and no oscillation is observed.

I found that: a) the presence of the stabilisation algorithm with small time steps does not introduce a significant difference in the outcome (less than a meter of deviation after 1Myr); b) the time step was increased up to $\delta t = 20,000\text{yr}$ and the simulation remained stable (the vertical deviation differed by approximately 4m from the one obtained with a very small time step, which remains a remarkable result since it only represents $\sim 0.2\%$ of the element size); c) time steps up to $\delta t = 50,000\text{yr}$ remain stable but lead to oscillations at the beginning and deviate in the end by about 5meters .

Finally, let us look at the mathematical details behind this stabilisation, as presented in Kaus *et al.* (2010) [681].

FINISH!!



Relevant Literature: Furuchi (2011) [428]

9.2 Convergence criterion for nonlinear iterations

nlconvcrit.tex

Disclaimer: the topic of nonlinear PDEs solving is vast and has received much attention from the mathematical community. In what follows I present a few key ideas which are at the core of many codes and publications in computational geodynamics.

Looking at the conservation equations that we must solve, i.e. conservation of mass, momentum and energy, we find that more often than not the coefficients of these PDEs depend on the strain rate, temperature, pressure, etc ... This makes solving the PDEs even harder. Also the advection term $\vec{v} \cdot \vec{\nabla}$ couples the two primary variables velocity and temperature.

One simple approach consists first in 'separating' the mass+momentum equations from the energy equation: one solves the first two equations assuming temperature known while the energy equation is solved assuming velocity and pressure known. Better schemes obviously exist and iterate on these equations until convergence for velocity, pressure and temperature is reached (see for instance the ASPECT manual). In what follows I focus on the mass and momentum equations assuming temperature known.

The main source of nonlinearity lies in the (effective) viscosity which often depends on strain rate and pressure (note that density can also depend on pressure in compressible cases):

$$\begin{aligned}\vec{\nabla} \cdot (2\eta_{\text{eff}}(\dot{\epsilon}, p) \dot{\epsilon}) - \vec{\nabla} p + \rho \vec{g} &= \vec{0} \\ \vec{\nabla} \cdot \vec{v} &= 0\end{aligned}$$

Simply put, in order to solve these equations and obtain the velocity and pressure fields I need to specify the density and viscosity (and of course appropriate boundary conditions!), but in order to compute the viscosity I need the strain rate and pressure fields.

The simplest approach here consists in so-called Picard iterations as explained in Section 9.19.

Let us start with the penalty-based FEM codes. In this case the mass and momentum equations are 'merged' into a single PDE where pressure has been eliminated:

$$\vec{\nabla} \cdot (2\eta_{\text{eff}}(\dot{\epsilon}, p) \dot{\epsilon}) + \lambda \vec{\nabla} (\vec{\nabla} \cdot \vec{v}) + \rho \vec{g} = \vec{0}$$

In this case the algorithm is simple:

1. start with a guess for the velocity and pressure fields, i.e. \vec{v}^{old} and \vec{p}^{old}
2. compute the effective viscosity field with \vec{v}^{old} and \vec{p}^{old}
3. solve PDE, obtain new solution \vec{v}^{new}
4. compute \vec{p}^{new} from \vec{v}^{new}
5. assess convergence, i.e. answer 'how close are the newly obtained fields from the old ones?'
6. $\vec{v}^{\text{old}} \leftarrow \vec{v}^{\text{new}}$, and $\vec{p}^{\text{old}} \leftarrow \vec{p}^{\text{new}}$
7. if not converged go back to 2, else exit

Thieulot (2011) [1258] computes the means $\langle \vec{v}^i \rangle$, $\langle \vec{v}^{i+1} \rangle$, and the variances σ_v^i and σ_v^{i+1} followed by the correlation

$$R^{i,i+1} = \frac{\langle (\vec{v}^i - \langle \vec{v}^i \rangle) \cdot (\vec{v}^{i+1} - \langle \vec{v}^{i+1} \rangle) \rangle}{\sqrt{\sigma_v^i \sigma_v^{i+1}}}$$

Since the correlation is normalised, it takes values between 0 (very dissimilar velocity fields) and 1 (very similar fields). The following convergence criterion, formulated in terms of the variable $\chi = 1 - R^{i,i+1}$ has been implemented: convergence is reached when $\chi < tol$. Since pressure is a derived quantity from velocity, if velocity is converged so is pressure⁶.

When the algorithm above is close to convergence then $\vec{\mathcal{V}}^i$ and $\vec{\mathcal{V}}^{i-1}$ are close. If these were scalar quantities we could subtract them and look at the (absolute) difference: if it is 'small enough' then the algorithm has converged. However there are two problems with this:

1. $\vec{\mathcal{V}}^i$ and $\vec{\mathcal{V}}^{i-1}$ are vector quantities (with potentially millions of values) so in order to measure a scalar difference between these we must take the norm of the difference, or $||\vec{\mathcal{V}}^i - \vec{\mathcal{V}}^{i-1}||$ and it is common to take the L^2 -norm.
2. we do not know a priori the (magnitude of the) solution so that 'small enough' is a dangerous statement. We could check for $||\vec{\mathcal{V}}^i - \vec{\mathcal{V}}^{i-1}|| < tol$ and set $tol = 10^{-6}$ for example. However in geodynamics velocities are of the order of a cm yr^{-1} which is about $3.1 \cdot 10^{-10} \text{m s}^{-1}$. Small velocity changes would then be enforced only if $tol < 10^{-12}$. This value might prove completely unpractical for other applications. In light of all this one then resorts to assessing the *relative* change in the velocity by normalising the previous quantity by the average velocity in the domain $||\vec{\mathcal{V}}^i||$.

This is the very approach taken by Spiegelman *et al.* [1187] who monitor the relative changes in the solution from iteration to iteration:

$$\frac{||\Delta \vec{\mathcal{V}}||_{L2}}{||\vec{\mathcal{V}}||_{L2}} = \left(\frac{\int_{\Omega} (\vec{\mathcal{V}}_i - \vec{\mathcal{V}}_{i-1}) \cdot (\vec{\mathcal{V}}_i - \vec{\mathcal{V}}_{i-1}) dV}{\int_{\Omega} \vec{\mathcal{V}}_i \cdot \vec{\mathcal{V}}_i dV} \right)^{1/2}$$

Of course, if a mixed formulation is used where velocity and pressure are solved for unknowns, the same monitoring can be done for pressure:

$$\frac{||\Delta \vec{\mathcal{P}}||_{L2}}{||\vec{\mathcal{P}}||_{L2}} = \left(\frac{\int_{\Omega} (\vec{\mathcal{P}}_i - \vec{\mathcal{P}}_{i-1}) \cdot (\vec{\mathcal{P}}_i - \vec{\mathcal{P}}_{i-1}) dV}{\int_{\Omega} \vec{\mathcal{P}}_i \cdot \vec{\mathcal{P}}_i dV} \right)^{1/2}$$

Convergence is reached when both are below 0.001 (as in Lemiale *et al.* (2008) [764]) or 0.0001 (as in Kaus *et al.* (2010) [679]).

The last option is via the nonlinear residual. Coming back to the penalty formulation, we can form the nonlinear residual as follows:

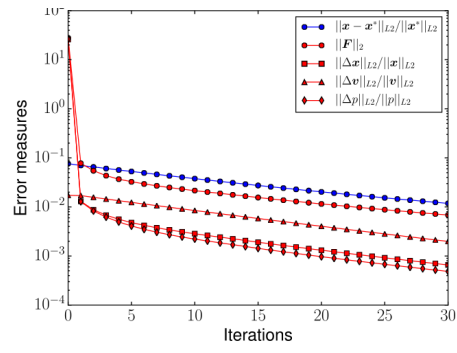
$$\vec{\mathcal{R}}^i = \mathbb{K}(\eta_{\text{eff}}(\dot{\epsilon}^i, p^i)) \cdot \vec{\mathcal{V}}^i - \vec{f}$$

where \mathbb{K} is defined in Section 7.4. Close to convergence $\vec{\mathcal{V}}^i$ and $\vec{\mathcal{V}}^{i-1}$ are very close so that we expect the residual $\vec{\mathcal{R}}$ to become smaller and smaller. In order to extract a scalar from $\vec{\mathcal{R}}$ we once again resort to the L^2 -norm and we also wish to monitor its relative change. In this case it is customary to use $\vec{\mathcal{R}}^0 = \vec{f}$ so that the convergence criterion becomes

$$\frac{||\vec{\mathcal{R}}^i||}{||\vec{\mathcal{R}}^0||} < tol.$$

FINISH: explain problem with mixed formulation!

⁶Two caveats here: the amplitude of the checkerboard mode might come into play - in the case $Q_1 \times P_0$ elements are used- and so does the applied smoothing.



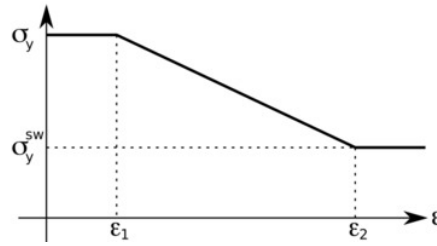
Taken from Spiegelman, May & Wilson *et al.* (2016). Example of reported nonlinear convergence.

9.3 Strain weakening

strainweakening.tex

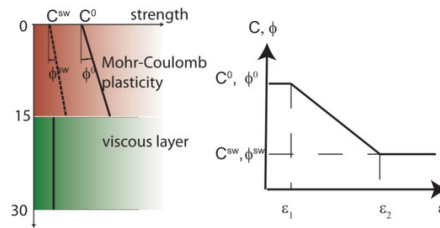
Several mechanisms may contribute to strain or strain rate dependent weakening but their relative and absolute importance is poorly constrained. Furthermore, weakening mechanisms are often crudely parameterised in geodynamical codes with simple mathematical functions and a limited number of parameters.

For example, Allken, Huismans, and Thieulot [10] (2011) authors use a von Mises plasticity formulation so that the rheology is parameterised by the cohesion c , or $c = \sigma_y$ in their notations. The yield strength σ_y starts is constant until the strain ε reaches the threshold value ε_1 . It then decreases linearly from σ_y to σ_y^{sw} between ε_1 and ε_2 . For strain values $\varepsilon > \varepsilon_2$, the yield strength remains constant at σ_y^{sw} .



Taken from Allken, Huismans, and Thieulot [10] (2011).

The same authors in a subsequent study use a Drucker-Prager rheology parameterised by cohesion c and friction angle ϕ . They use the same approach as before but now both parameters are subjected to strain weakening:



Taken from Allken, Huismans, and Thieulot [9] (2012), see also Thieulot [1258] (2011).

They further define the factor $R = C^0/C^{sw} = \phi^0/\phi^{sw} \geq 1$ which is a proxy for the ratio σ_y/σ_y^{sw} where $\sigma_y = p \sin \phi + c \cos \phi$, and carry out 3D crustal extensional models for R between 2 and 5.

- In Le Pourhiet, May, Huille, Watremez, and Leroy [756] (2017) the authors also define

$$\tau_y = p \sin(\phi(\varepsilon^p)) + c_0 \cos(\phi(\varepsilon^p))$$

but the cohesion is regarded to be constant. The angle of friction ϕ is assumed to decrease as a function of the accumulated plastic strain ε^p to

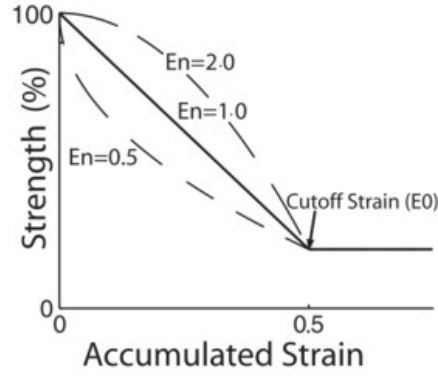
$$\phi(\varepsilon^p) = \max \left(\phi_\infty, \phi_0 - \frac{\varepsilon^p(\phi_0 - \phi_\infty)}{\varepsilon_\infty^p} \right)$$

This equation defines an empirical softening relation which reduces the friction angle linearly with accumulated plastic strain. ϕ_0 defines the initial friction angle, ε_∞^p represents the measure of plastic strain after which complete softening is achieved and internal friction angle reaches ϕ_∞ . Plastic strain represents an integrated, tensorial invariant measure of the deformation which has occurred due to plastic yielding. Thus, the quantity ε^p can be regarded as a simplified measure of material damage.

- In Dyksterhuis *et al.* Dyksterhuis, Rey, Mueller, and Moresi [356] (2007) a variant of the above formulation is used:

$$f(\varepsilon) = \begin{cases} 1 - (1 - a)(\varepsilon/\varepsilon_0)^n & \varepsilon \leq \varepsilon_0 \\ a & \varepsilon \geq \varepsilon_0 \end{cases}$$

where ε is the accumulated plastic strain, ε_0 is the saturation strain beyond which no further weakening takes place, n is an exponent that controls the shape of the function and a is a maximum value of strain weakening beyond which no further weakening occurs. This equation leads to the following plot:



Strain-softening behaviour showing strength weakening from 100 to 20% after an accumulated strain of 0.5, after which no further weakening occurs.

Dashed lines show the effect of the exponential parameter (En) on the curve. Taken from Dyksterhuis, Rey, Mueller, and Moresi [356] (2007).

Although it is not specified in Dyksterhuis, Rey, Mueller, and Moresi [356] what f is, other users of the code specify that the yield strength is given by

$$\sigma_y = (B_0 + B_1 p) f(\varepsilon)$$

where p is the pressure, B_0 is the cohesion, or yield stress at zero pressure, and B_p is the pressure dependence of the yield stress, equivalent to the friction coefficient in Byerlee's law.

In Yang, Moresi, Zhao, Sandiford, and Whittaker [1380] (2018) the authors take a different approach:

$$C = C_0 + C_1 \exp\left(-\frac{\varepsilon_{plast}}{\varepsilon_{ref}}\right)$$


$$\mu = \mu_0 + \mu_1 \exp\left(-\frac{\varepsilon_{plast}}{\varepsilon_{ref}}\right)$$

where C_0 and $C_0 + C_1$ represent the minimum and maximum cohesions, respectively; μ_0 and $\mu_0 + \mu_1$ represent the minimum and maximum frictional coefficients, respectively. ε_{plast} and ε_{ref} represent accumulated plastic strain and reference strain, respectively.

- In Leroy and Ortiz [772] (1989) the authors describe another formulation for plastic hardening. The angle of friction changes with the accumulated plastic strain:

$$\sin \phi = \sin \phi_i + \frac{2(\sin \phi_f - \sin \phi_i) \sqrt{\varepsilon_c^p \varepsilon^p}}{\varepsilon^p + \varepsilon_c^p}$$

where ϕ transitions from an initial value ϕ_i to a maximum ϕ_f attained when the effective plastic strain reaches a critical value ε_c^p . When $\varepsilon^p \rightarrow \varepsilon_c^p$ then $\phi \rightarrow \phi_f$.

 Relevant Literature: Sterpi [1211] (1999), Nijholt and Govers [942] (2015).

9.4 Assigning values to quadrature points

As we have seen in Section 7, the building of the elemental matrix and rhs requires (at least) to assign a density and viscosity value to each quadrature point inside the element. Depending on the type of modelling, this task can prove more complex than one might expect and have large consequences on the solution accuracy.

Here are several options:

- The simplest way (which is often used for benchmarks) consists in computing the 'real' coordinates (x_q, y_q, z_q) of a given quadrature point based on its reduced coordinates (r_q, s_q, t_q) , and passing these coordinates to a function which returns density and/or viscosity at this location. For instance, for the Stokes sphere:

```
def rho(x,y):
    if (x-.5)**2+(y-0.5)**2<0.123**2:
        val=2.
    else:
        val=1.
    return val

def mu(x,y):
    if (x-.5)**2+(y-0.5)**2<0.123**2:
        val=1.e2
    else:
        val=1.
    return val
```

This is very simple, but it has been shown to potentially be problematic. In essence, it can introduce very large contrasts inside a single element and perturb the quadrature. Please read section 3.3 of [560] and/or have a look at the section titled "Averaging material properties" in the ASPECT manual.

- another similar approach consists in assigning a density and viscosity value to the nodes of the FE mesh first, and then using these nodal values to assign values to the quadrature points. Very often, and quite logically, the basis functions are used to this effect. Indeed we have seen before that for any point (r, s, t) inside an element we have

$$f_h(r, s, t) = \sum_i^m f_i N_i(r, s, t)$$

where the f_i are the nodal values and the N_i the corresponding basis functions.

In the case of linear elements (Q_1 basis functions), this is straightforward. In fact, the basis functions N_i can be seen as moving weights: the closer the point is to a node, the higher the weight (basis function value).

However, this is quite another story for quadratic elements (Q_2 basis functions). In order to illustrate the problem, let us consider a 1D problem. The basis functions are

$$N_1(r) = \frac{1}{2}r(r-1) \quad N_2(r) = 1-r^2 \quad N_3(r) = \frac{1}{2}r(r+1)$$

Let us further assign: $\rho_1 = \rho_2 = 0$ and $\rho_3 = 1$. Then

$$\rho_h(r) = \sum_i^m \rho_i N_i(r) = N_3(r)$$

There lies the core of the problem: the $N_3(r)$ basis function is negative for $r \in [-1, 0]$. This means that the quadrature point in this interval will be assigned a negative density, which is nonsensical and numerically problematic!

use 2X Q1. write about it !

The above methods work fine as long as the domain contains a single material. As soon as there are multiple fluids in the domain a special technique is needed to track either the fluids themselves or their interfaces. Let us start with markers. We are then confronted to the infernal trio (a *menage a trois*?) which is present for each element, composed of its nodes, its markers and its quadrature points.

Each marker carries the material information (density and viscosity). This information must ultimately be projected onto the quadrature points. Two main options are possible: an algorithm is designed and projects the marker-based fields onto the quadrature points directly or the marker fields are first projected onto the FE nodes and then onto the quadrature points using the techniques above.

At a given time, every element e contains n^e markers. During the FE matrix building process, viscosity and density values are needed at the quadrature points. One therefore needs to project the values carried by the markers at these locations. Several approaches are currently in use in the community and the topic has been investigated by [330] and [352] for instance.

ELEFANT adopts a simple approach: viscosity and density are considered to be elemental values, i.e. all the markers within a given element contribute to assign a unique constant density and viscosity value to the element by means of an averaging scheme.

While it is common in the literature to treat the so-called arithmetic, geometric and harmonic means as separate averagings, I hereby wish to introduce the notion of generalised mean, which is a family of functions for aggregating sets of numbers that include as special cases the arithmetic, geometric and harmonic means.

If p is a non-zero real number, we can define the generalised mean (or power mean) with exponent p of the positive real numbers a_1, \dots, a_n as:

$$M_p(a_1, \dots, a_n) = \left(\frac{1}{n} \sum_{i=1}^n a_i^p \right)^{1/p} \quad (9.14)$$

and it is trivial to verify that we then have the special cases:

$$M_{-\infty} = \lim_{p \rightarrow -\infty} M_p = \min(a_1, \dots, a_n) \quad (\text{minimum}) \quad (9.15)$$

$$M_{-1} = \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \dots + \frac{1}{a_n}} \quad (\text{harm. avrg.}) \quad (9.16)$$

$$M_0 = \lim_{p \rightarrow 0} M_p = \left(\prod_{i=1}^n a_i \right)^{1/n} \quad (\text{geom. avrg.}) \quad (9.17)$$

$$M_{+1} = \frac{1}{n} \sum_{i=1}^n a_i \quad (\text{arithm. avrg.}) \quad (9.18)$$

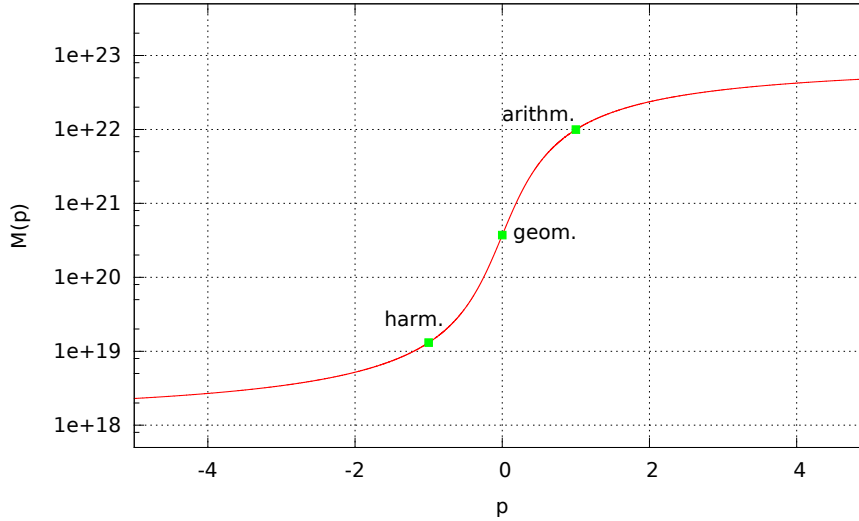
$$M_{+2} = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2} \quad (\text{root mean square}) \quad (9.19)$$

$$M_{+\infty} = \lim_{p \rightarrow +\infty} M_p = \max(a_1, \dots, a_n) \quad (\text{maximum}) \quad (9.20)$$

Note that the proofs of the limit convergence are given in [169].

An interesting property of the generalised mean is as follows: for two real values p and q , if $p < q$ then $M_p \leq M_q$. This property has for instance been illustrated in Fig. 20 of [1124].

One can then for instance look at the generalised mean of a randomly generated set of 1000 viscosity values within $10^{18} Pa.s$ and $10^{23} Pa.s$ for $-5 \leq p \leq 5$. Results are shown in the figure hereunder and the arithmetic, geometric and harmonic values are indicated too. The function M_p assumes an arctangent-like shape: very low values of p will ultimately yield the minimum viscosity in the array while very high values will yield its maximum. In between, the transition is smooth and occurs essentially for $|p| \leq 5$.



▷ `python_codes/fieldstone_markers_avrg`

9.5 Matrix (Sparse) storage

storage.tex

The FE matrix (or the blocks which compose it) is the result of the assembly process of all elemental matrices. Its size can become quite large when the resolution is being increased (from thousands of lines/columns to tens of millions).

One important property of the matrix is its sparsity. Typically much less than 1% of the matrix terms is not zero and this means that the matrix storage can and *should* be optimised. Clever storage formats were designed early on since the amount of RAM memory in computers was the limiting factor 3 or 4 decades ago [1092].

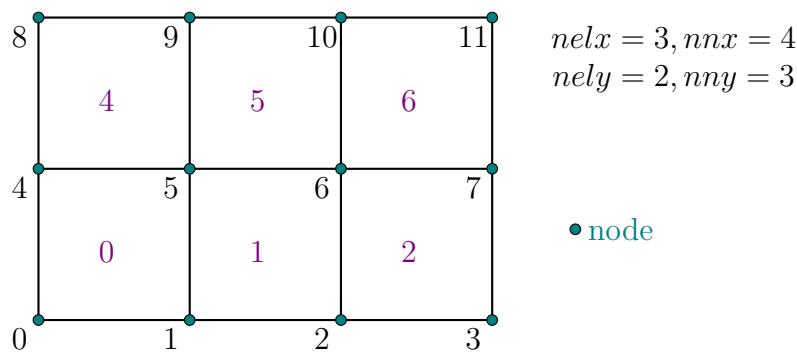
There are several standard formats⁷, e.g.:

- compressed sparse row format (CSR)
- compressed sparse column format (CSC)
- the Coordinate Format (COO)
- Skyline Storage Format

I focus on the CSR format in what follows since it is the most common format and it is the one used in ELEFANT .

9.5.1 2D domain - Q_1 - One degree of freedom per node

Let us consider again the 3×2 element grid which counts 12 nodes.



In the case there is only a single degree of freedom per node, the assembled FEM matrix \mathbf{M} will look like this:

$$\mathbf{M} = \begin{pmatrix} \square & \square & & & \square & \square & & & & & & \\ & \square & \square & \square & & \square & \square & \square & & & & \\ & & \square & \square & \square & & \square & \square & \square & & & \\ & & & \square & \square & & & \square & \square & & & \\ \square & \square & & & \square & \square & & & \square & \square & & \\ \square & \square & \square & & \square & \square & \square & & \square & \square & \square & \\ & \square & \square & \square & & \square & \square & \square & & \square & \square & \square \\ & & \square & \square & & & \square & \square & & \square & \square & \\ & & & \square & \square & & & \square & \square & & & \\ & & & & \square & \square & \square & & \square & \square & \square & \\ & & & & & \square & \square & \square & & \square & \square & \square \\ & & & & & & \square & \square & & \square & \square & \\ & & & & & & & \square & \square & & \square & \square \end{pmatrix}$$

⁷https://en.wikipedia.org/wiki/Sparse_matrix

where the \square stand for non-zero terms. This matrix structure stems from the fact that

- node 0 sees nodes 0,1,4,5 (1st line/column of the matrix)
- node 1 sees nodes 0,1,2,4,5,6 (2nd line/column of the matrix)
- node 2 sees nodes 1,2,3,5,6,7 (3rd line/column of the matrix)
- node 3 sees nodes 2,3,6,7
- node 4 sees nodes 0,1,4,5,8,9
- node 5 sees nodes 0,1,2,4,5,6,8,9,10
- node 6 sees nodes 1,2,3,5,6,7,9,10,11
- node 7 sees nodes 2,3,6,7,10,11
- node 8 sees nodes 4,5,8,9
- node 9 sees nodes 4,5,6,8,9,10
- node 10 sees nodes 5,6,7,9,10,11
- node 11 sees nodes 6,7,10,11 (last line/column of the matrix)

In light thereof, we have

- 4 corner nodes which have 4 neighbours (counting themselves)
- $2(nnx-2)$ nodes which have 6 neighbours
- $2(nny-2)$ nodes which have 6 neighbours
- $(nnx-2) \times (nny-2)$ nodes which have 9 neighbours

In total, the number of non-zero terms in the matrix above is then:

$$NZ = 4 \times 4 + 4 \times 6 + 2 \times 6 + 2 \times 9 = 70$$

and in general, we would then have:

$$NZ = 4 \times 4 + [2(nnx - 2) + 2(nny - 2)] \times 6 + (nnx - 2)(nny - 2) \times 9$$

Let us temporarily assume $nnx = nny = n$. The matrix size (total number of unknowns) is then $N = n^2$ and

$$NZ = 16 + 24(n - 2) + 9(n - 2)^2$$

A full matrix array would contain $N^2 = n^4$ terms. The ratio of NZ (the actual number of reals to store) to the full matrix size (the number of reals a full matrix contains) is then

$$R = \frac{16 + 24(n - 2) + 9(n - 2)^2}{n^4}$$

It is then obvious that when n is large enough $R \sim 1/n^2$.

CSR stores the nonzeros of the matrix row by row, in a single indexed array A of double precision numbers. Another array $COLIND$ contains the column index of each corresponding entry in the A array. A third integer array $RWPTR$ contains pointers to the beginning of each row, which an additional pointer to the first index following the nonzeros of the matrix A . A and $COLIND$ have length NZ and $RWPTR$ has length $N+1$.

In the case of the here-above matrix, the arrays $COLIND$ and $RWPTR$ will look like:

$$\begin{aligned} COLIND &= (0, 1, 4, 5, 0, 1, 2, 4, 5, 6, 1, 2, 3, 5, 6, 7, \dots, 6, 7, 10, 11) \\ RWPTR &= (0, 4, 10, 16, \dots) \end{aligned}$$

9.5.2 2D domain - Q_1 - Symmetric matrix CSR storage

If the matrix is symmetric, i.e. $\mathbf{M} = \mathbf{M}^T$, then we may wish to only store half of it, always in the interest of saving memory. Only the following remaining \square entries are relevant now:

$$M = \begin{pmatrix} \square & & & & & & & & & \\ & \square & & & & & & & & \\ & & \square & & & & & & & \\ & & & \square & & & & & & \\ & & & & \square & & & & & \\ & & & & & \square & & & & \\ & & & & & & \square & & & \\ & & & & & & & \square & & \\ & & & & & & & & \square & \\ & & & & & & & & & \square \end{pmatrix}$$

We see that the number of nonzeros is now

$$NZ_{symm} = \frac{NZ - n}{2} + n$$

and in this case $NZ_{symm} = (70 - 12)/2 + 12 = 41$. Then

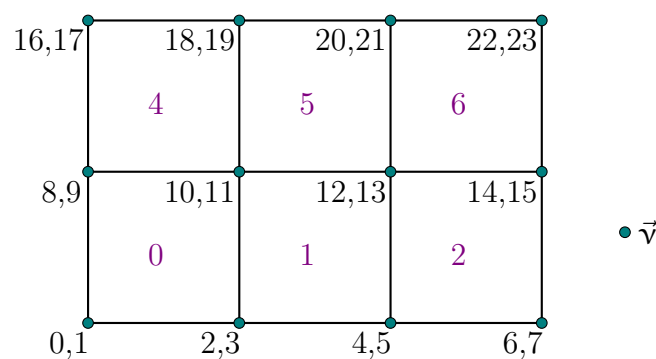
$$\begin{aligned} COLIND &= (0, 1, 4, 5, 1, 2, 4, 5, 6, 3, 5, 6, 7, \dots, 11) \\ RWPTR &= (0, 4, 9, 14, \dots) \end{aligned}$$

In case the numbering is Fortran-like, then

$$\begin{aligned} ja = COLIND &= (1, 2, 5, 6, \quad 2, 3, 5, 6, 7, \quad 3, 4, 6, 7, 8, \quad 4, 7, 8, \quad 5, 6, 9, 10, \quad 6, 7, 9, 10, 11, \\ &\quad 7, 8, 10, 11, 12, \quad 8, 11, 12, \quad 9, 10, \quad 10, 11, \quad 11, 12, \quad 12) \\ ia = RWPTR &= (1, 5, 10, 15, 18, 22, 27, 32, 35, 37, 39, 41, 42) \end{aligned}$$

9.5.3 2D domain - Q_1 - Two degrees of freedom per node

When there are now two degrees of freedom per node, such as in the case of the Stokes equation in two-dimensions, the size of the \mathbb{K} matrix is given by $NfemV = nnx * nny * ndofV$ where $NfemV$ is the total number of velocity degrees of freedom.



In the case of the small grid above, we have then $NfemV = 24$ and elemental matrices are now 8×8 in size.

We still have

- 4 corner nodes which have 4 neighbours
- $2(nnx - 2)$ nodes which have 6 neighbours
- $2(nny - 2)$ nodes which have 6 neighbours
- $(nnx - 2) \cdot (nny - 2)$ nodes which have 9 neighbours

but now each degree of freedom from a node sees the other two degrees of freedom of another node too. In that case, the number of nonzeros has been multiplied by four and the assembled FEM matrix looks like:

A large grid of 100 small squares arranged in a 10x10 pattern, with some squares missing, forming a sparse, irregular shape. The grid is enclosed in large parentheses.

Note that the degrees of freedom are organised as follows:

$$(u_0, v_0, u_1, v_1, u_2, v_2, \dots, u_{11}, v_{11})$$

In general, we would then have:

$$NZ = 4[4 \times 4 + [2(nnx - 2) + 2(nny - 2)] \times 6 + (nnx - 2)(nny - 2) \times 9]$$

and in the case of the small grid, the number of non-zero terms in the matrix is then:

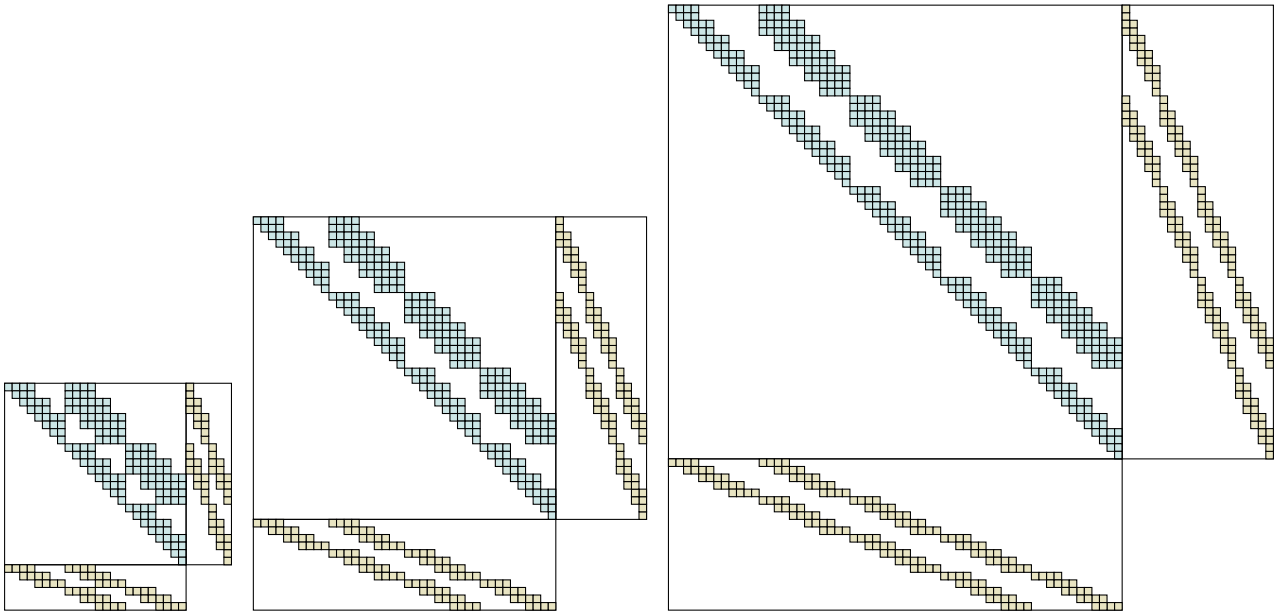
$$NZ = 4[4 \times 4 + 4 \times 6 + 2 \times 6 + 2 \times 9] = 280$$

In the case of the here-above matrix, the arrays COLIND and RWPTR will look like:

$$\begin{aligned} COLIND &= (0, 1, 2, 3, 8, 9, 10, 11, 0, 1, 2, 3, 8, 9, 10, 11, \dots) \\ RWPTR &= (0, 8, 16, 28, \dots) \end{aligned}$$

Assuming we are using $Q_1 \times P_0$ elements, the structure of the matrix \mathbb{G}_{el}^T is as follows (the 6 pressure dofs are connected to 24 velocity dofs):

$$\left(\begin{array}{cccccccccccccccccccccccccccc} & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 & 21 & 22 & 23 \end{array} \right) \quad (9.21)$$



From left to right: Nonzero structures of the assembled Stokes matrix for a 3×2 , 4×3 and 5×4 mesh of $Q_1 \times P_0$ elements.

Assuming we are now using $Q_1 \times Q_1$ elements (without bubble), the structure of the matrix \mathbb{G}_{el}^T is different: we now have 12 pressure dofs which are coupled to 24 velocity dofs:

$$\left(\begin{array}{cccccccccccccccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 & 19 & 20 & 21 & 22 & 23 & 24 \\ 0 & \square & \square & \square & \square & & & & & \square & \square & \square & \square & & & & & & & & & & & & \\ 1 & \square & \square & \square & \square & \square & \square & & & \square & \square & \square & \square & \square & \square & & & & & & & & & & \\ 2 & & & \square & \square & \square & \square & \square & \square & & & \square & \square & \square & \square & \square & \square & & & & & & & & \\ 3 & & & & & \square & \square & \square & \square & & & & \square & \square & \square & \square & \square & & & & & & & & \\ & & & & & & & & & & & & & & & & & & & & & & & & \\ \dots & & & & & & & & & & & & & & & & & & & & & & & & \\ 9 & & & & & & & & & \square & \square & \square & \square & \square & \square & & & & & \square & \square & \square & \square & & & \\ 10 & & & & & & & & & & \square & \square & \square & \square & \square & \square & \square & & & \square & \square & \square & \square & \square & \square & \\ 11 & & & & & & & & & & & \square & \square & \square & \square & \square & \square & & & & & \square & \square & \square & \square \end{array} \right) \quad (9.22)$$

If now the velocity dofs are organised as follows

$$(u_0, u_1, u_2, \dots, u_{11}, v_0, v_1, v_2, \dots, v_{11})$$

then the sparsity pattern of the assembled \mathbb{K} matrix looks like

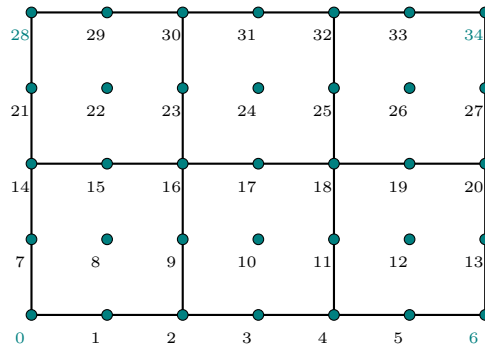
$$(9.23)$$

$$\mathbb{K} = \begin{pmatrix} \mathbb{K}_{xx} & \mathbb{K}_{xy} \\ \mathbb{K}_{yx} & \mathbb{K}_{yy} \end{pmatrix} = \begin{pmatrix} \mathbb{K}_1 & \mathbb{K}_1 \\ \mathbb{K}_1 & \mathbb{K}_1 \end{pmatrix}$$

When there are now two degrees of freedom per node, such as in the case of the Stokes equation in two-dimensions, the size of the \mathbb{K} matrix is given $NfemV = nnx * nny * ndofV$ where $NfemV$ is the total number of velocity degrees of freedom. What is different here is that for Q_2 elements we have $nnx = 2 * nelx + 1$ and $nny = 2 * nely + 1$.



In the case of the small grid above, we have then $nex = 3$, $nely = 2$, so that $nnx = 7$ and $nnx = 5$, and then $N_{femV} = 7 * 5 * 2 = 70$ and elemental matrices are now 18×18 in size.



(tikz_3x2_Q2.tex)

Concretely here:

- nodes 0,6,28,34 see 9 nodes (corners)
- nodes 1,3,5,7,8,10,12,13,21,22,24,26,27,29,31,33 see 9 nodes
- nodes 2,4,9,11,14,15,17,19,20,23,25,30,32, see 15 nodes
- nodes 16,18 see 25 nodes

If there was only one dof per node, we would find the number of non zeros as follow:

$$NZ = 4 * 9 + 16 * 9 + 13 * 15 + 2 * 25 = 36 + 144 + 195 + 50 = 425$$

But since there are two velocity dofs per node, we find that the total number of nonzeros is 4 times higher, i.e.

$$NZ = 1700$$

And if we choose for a symmetric CSR storage:

$$NZ_{symm} = \frac{NZ - n}{2} + n = \frac{1700 - 70}{2} + 70 = 885$$

Let us now turn to the real case of 2 dofs per node and establish who sees who:

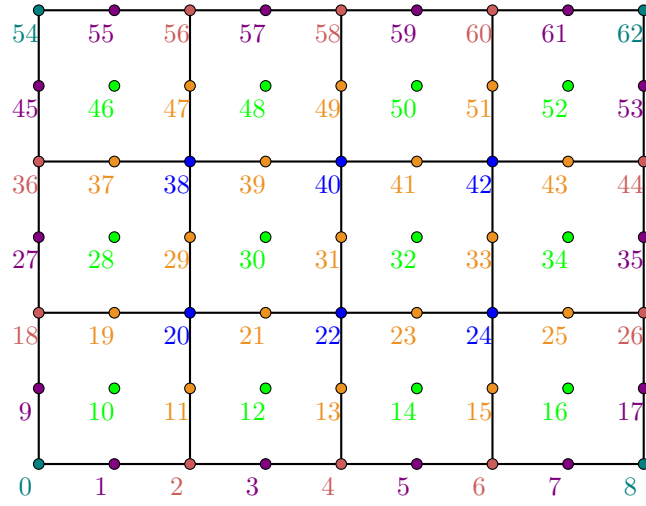
| dof | sees other dofs | total |
|-----|-------------------------------------------------------------------------------------|-------|
| 0 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 1 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 2 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 3 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 4 | 0,1,2,3,4,5,6,7,8,9,14,15,16,17,18,19,20,21,22,23,28,29,30,31,32,33,34,35,36,37 | 30 |
| 5 | 0,1,2,3,4,5,6,7,8,9,14,15,16,17,18,19,20,21,22,23,28,29,30,31,32,33,34,35,36,37 | 30 |
| 6 | 4,5,6,7,8,9,18,19,20,21,22,23,32,33,34,35,36,37 | 18 |
| 7 | 4,5,6,7,8,9,18,19,20,21,22,23,32,33,34,35,36,37 | 18 |
| 8 | 4,5,6,7,8,9,10,11,12,13,18,19,20,21,22,23,24,25,26,27,32,33,34,35,36,37,38,39,40,41 | 30 |
| 9 | 4,5,6,7,8,9,10,11,12,13,18,19,20,21,22,23,24,25,26,27,32,33,34,35,36,37,38,39,40,41 | 30 |
| 10 | 8,9,10,11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | 18 |
| 11 | 8,9,10,11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | 18 |
| 12 | 8,9,10,11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | 18 |
| 13 | 8,9,10,11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | 18 |
| 14 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 15 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 16 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 17 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| ... | ... | ... |
| 68 | 36,37,38,39,50,51,52,53,54,55,64,65,66,67,68,69 | 18 |
| 69 | 36,37,38,39,50,51,52,53,54,55,64,65,66,67,68,69 | 18 |

The second column of this array is the content of the ja array if a full storage is used. If we now use a symmetric storage then:

| dof | sees other dofs | total |
|-----|-----------------------------------------------------------------------------|-------|
| 0 | 0,1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 18 |
| 1 | 1,2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 17 |
| 2 | 2,3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 16 |
| 3 | 3,4,5,14,15,16,17,18,19,28,29,30,31,32,33 | 15 |
| 4 | 4,5,6,7,8,9,14,15,16,17,18,19,20,21,22,23,28,29,30,31,32,33,34,35,36,37 | |
| 5 | 5,6,7,8,9,14,15,16,17,18,19,20,21,22,23,28,29,30,31,32,33,34,35,36,37 | |
| 6 | 6,7,8,9,18,19,20,21,22,23,32,33,34,35,36,37 | |
| 7 | 7,8,9,18,19,20,21,22,23,32,33,34,35,36,37 | |
| 8 | 8,9,10,11,12,13,18,19,20,21,22,23,24,25,26,27,32,33,34,35,36,37,38,39,40,41 | |
| 9 | 9,10,11,12,13,18,19,20,21,22,23,24,25,26,27,32,33,34,35,36,37,38,39,40,41 | |
| 10 | 10,11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | |
| 11 | 11,12,13,22,23,24,25,26,27,36,37,38,39,40,41 | |
| 12 | 12,13,22,23,24,25,26,27,36,37,38,39,40,41 | |
| 13 | 13,22,23,24,25,26,27,36,37,38,39,40,41 | |
| 14 | 14,15,16,17,18,19,28,29,30,31,32,33 | |
| 15 | 15,16,17,18,19,28,29,30,31,32,33 | |
| 16 | 16,17,18,19,28,29,30,31,32,33 | |
| 17 | 17,18,19,28,29,30,31,32,33 | |
| ... | ... | ... |
| 68 | 68,69 | 2 |
| 69 | 69 | 1 |

In order establish a pattern we will need a bigger mesh:

(tikz_4x3_Q2.tex)



We have

- 4 corner nodes which have 9 neighbours
- nel mid-element nodes which have 9 neighbours
- $2 * nelx + 2 * nely$ mid-edge nodes on sides which have 9 neighbours
- $(nelx - 1) * nely + nelx * (nely - 1)$ internal mid-edges nodes which have 15 neighbours
- $2 * (nelx - 1) + 2 * (nely - 1)$ side nodes that have 15 neighbours
- $(nelx - 1) * (nely - 1)$ nodes which have 25 neighbours

In the end, the number of non-zeros (Q_1) is given by

$$\begin{aligned}
 NZ &= 4 * 9 \\
 &+ nel * 9 \\
 &+ (2 * nelx + 2 * nely) * 9 \\
 &+ [(nelx - 1) * nely + nelx * (nely - 1)] * 15 \\
 &+ [2 * (nelx - 1) + 2 * (nely - 1)] * 15 \\
 &+ (nelx - 1) * (nely - 1) * 25
 \end{aligned}$$

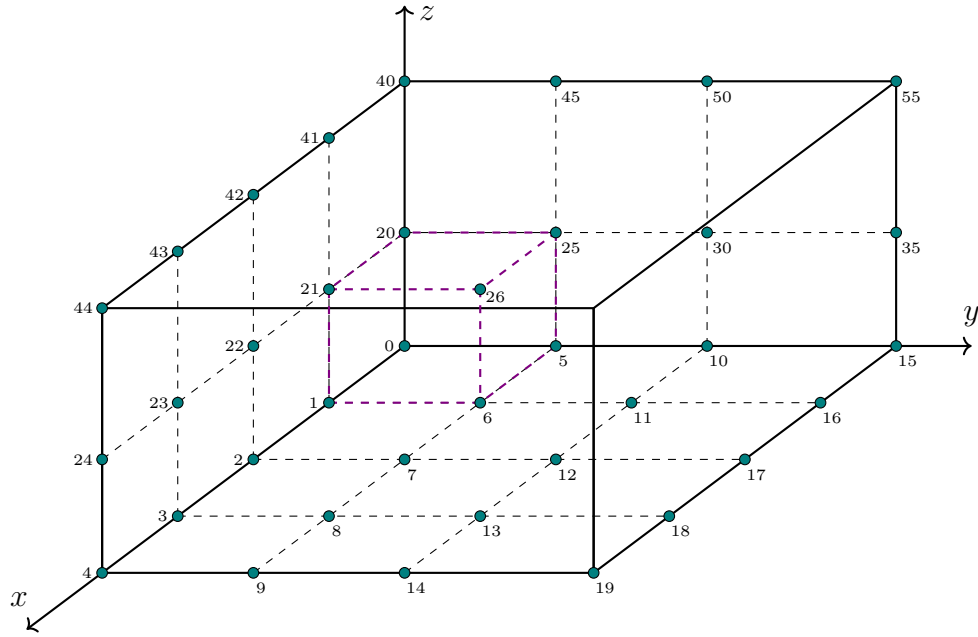
Verification: $nelx = 3$, $nely = 2$:

$$\begin{aligned}
 NZ &= 4 * 9 + 6 * 9 + (2 * 3 + 2 * 2) * 9 + [(3 - 1) * 2 + 3 * (2 - 1)] * 15 + [2 * (3 - 1) + 2 * (2 - 1)] * 15 + \\
 &= 36 + 54 + 10 * 9 + 7 * 15 + 6 * 15 + 2 * 25 \\
 &= 36 + 54 + 90 + 105 + 90 + 50 \\
 &= 425
 \end{aligned}$$

as expected.

9.5.5 3D domain - Q_1 - CSR storage - One degree of freedom

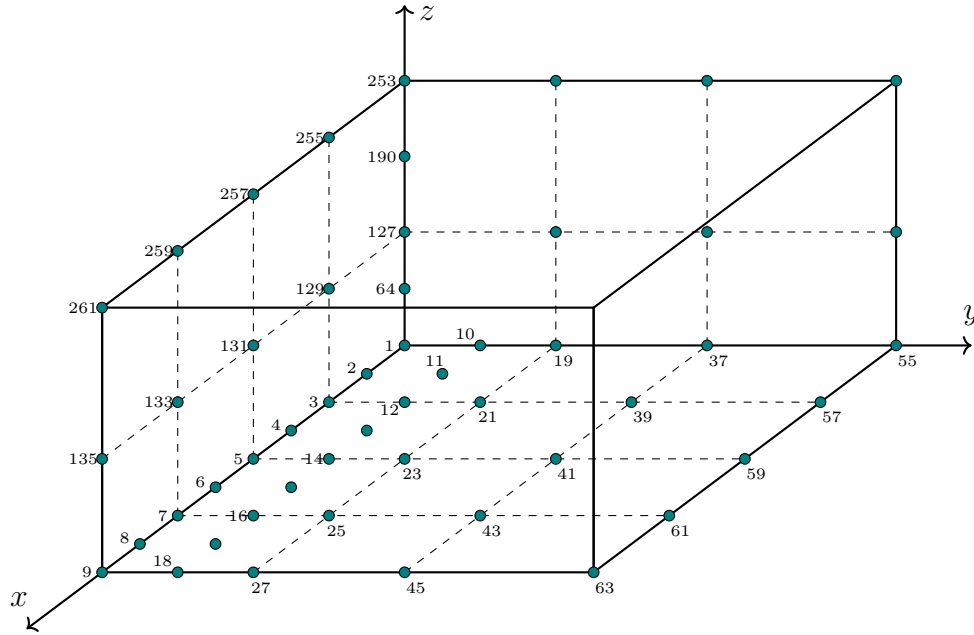
Let us consider a $3 \times 4 \times 2$ grid which counts $nnx \cdot nny \cdot nnz = 5 \cdot 4 \cdot 3 = 60$ nodes. The assembled FEM matrix \mathbb{K} size is then $N = nnx \times nny \times nnz \times ndof = 180$.



The total number of nonzeros in the case $ndof = 1$ would be decomposed as follows:

- 8 corners 'see' 8 neighbours
- 4 edges with $(nnx - 2)$ nodes in the x direction see 12 nodes
- 4 edges with $(nny - 2)$ nodes in the y direction see 12 nodes
- 4 edges with $(nnz - 2)$ nodes in the z direction see 12 nodes
- $2(nnx - 2)(nny - 2)$ nodes see 18 nodes
- $2(nnx - 2)(nnz - 2)$ nodes see 18 nodes
- $2(nny - 2)(nnz - 2)$ nodes see 18 nodes
- $(nnx - 2)(nny - 2)(nnz - 2)$ interior nodes see 27 nodes

9.5.6 3D domain - Q_2 - CSR storage - one degree of freedom



9.5.7 Matrix Storage in fieldstone

The majority of the early codes have the FE matrix being a full array

```
a_mat = np.zeros((Nfem,Nfem),dtype=np.float64)
```

and it is converted to CSR format on the fly in the solve phase:

```
sol = sps.linalg.spsolve(sps.csr_matrix(a_mat),rhs)
```

Note that linked list storages can be used (lil_matrix). Substantial memory savings but much longer compute times since it takes longer to write in such arrays. A conversion to CSR format is still necessary before calling the solver.

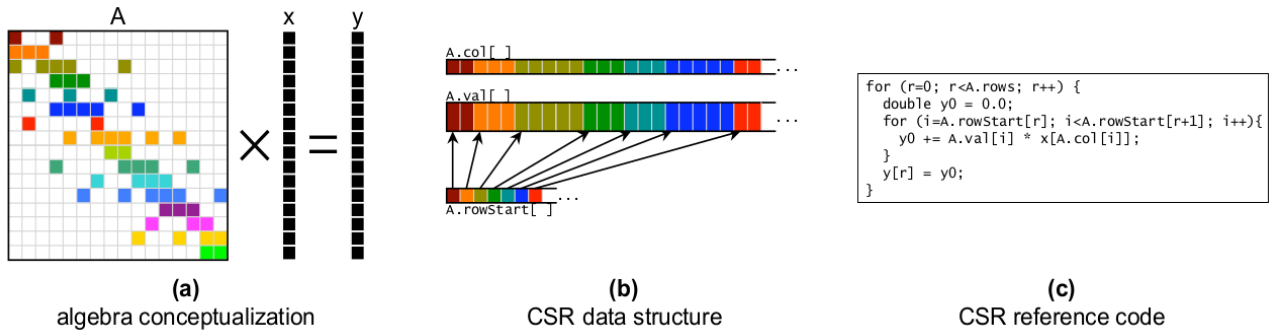
9.5.8 About Sparse Matrix-Vector multiplication

When/if the matrix M is stored in a two-dimensional array, its (left or right) multiplication by a vector is trivial. Either one resorts to writing a double for loop (not recommended), either one uses `numpy.dot`⁸ in python, or `matmul` in Fortran.

However, when the matrix is stored as a single continuous array, say CSR, how does this work? This question is *very important* since iterative solvers such as the Conjugate Gradient solver (see Section 9.33) rely extensively on multiplying the matrix by many different vectors.

The Sparse Matrix-Vector multiplication operation is often abbreviated SpMV. To quote Knepley [712]: "The Sparse Matrix-Vector Product (SpMV) is today a workhorse of scientific computing. It is a central kernel in iterative linear and nonlinear solvers for PDE, and now for many graph algorithms." As explained in Williams *et al.* (2007) [1361] (and in many other sources on the topic), the algorithm for a basic SpMV implementation is rather simple in its naive form.

⁸<https://numpy.org/doc/stable/reference/generated/numpy.dot.html>



Taken from Williams *et al.* (2008) [1362]. Sparse Matrix Vector Multiplication (SpMV). (a) visualization of the algebra: $\vec{y} \leftarrow \mathbf{A} \cdot \vec{x}$.
 (b) Standard compressed sparse row (CSR) representation of the matrix.

(c) The standard implementation of SpMV for a matrix stored in CSR. The outer loop is trivially parallelized without any data dependencies.

Let us assume that we wish to compute $\vec{y} = \mathbf{A} \cdot \vec{x}$ where \mathbf{A} is in CSR format. The pseudo code then goes as follows:

```
for i in range(0,m):
    y0=0
    for k in range(ROWPTR[i],ROWPTR[i+1]):
        y0 += VAL[k] * x[COLIND[k]]
    y[i]=y0
```

Although technically correct, this algorithm is problematic because the vector x array is accessed indirectly and this causes a non-optimal use of the processor, which in the end makes the calculation take longer than it should.


The following piece of code comes from ELEFANT . Note that here (ROWPTR= ia , COLIND= ja , VAL= mat)

```
subroutine spmv (nr , nc , nz , x , y , mat , ja , ia )
implicit none
integer , intent(in) :: nr , nc , nz
real(8) , intent(in) :: x(nc) , mat(nz)
real(8) , intent(out) :: y(nr)
integer , intent(in) :: ja(nz) , ia(nr+1)
real(8) t
integer i , k

do i = 1 , nr
    t = 0.0d0
    do k=ia(i) , ia(i+1)-1
        t = t + mat(k)*x(ja(k))
    end do
    y(i) = t
end do

end subroutine
```

How to make this calculation as efficiently as possible on CPUs and GPUs, on one thread or multiple threads has given rise to a lot of literature.

 **Relevant Literature** Krotkiewski & Dabrowski [733], Section 9.4 of Kepley [712], Williams *et al.* (2008) [1362]

9.5.9 SpMV and SpMV-T with the CSR format - a concrete example

(What follows was originally written for ELEFANT so that code excerpts and loop indexing are those of Fortran.)

Let us consider a sparse matrix \mathbb{G} which is not square (size is 3×5):

$$\mathbb{G}^T = \begin{pmatrix} 1 & 0 & 4 & 1 & 2 \\ 0 & 1 & 1 & 1 & 0 \\ 3 & 0 & 0 & 7 & 1 \end{pmatrix}$$

The number of rows is $nr = 3$, the number of columns is $nc = 5$ and the number of nonzeros is $nz = 10$.

Let us consider two vectors $\vec{\mathcal{V}}^T = (1, 1, 1, 1, 1)$ and $\vec{\mathcal{P}}^T = (1, 1, 1)$. Obviously, we have:

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} = \begin{pmatrix} 8 \\ 3 \\ 11 \end{pmatrix} \quad \text{and} \quad \mathbb{G} \cdot \vec{\mathcal{P}} = \begin{pmatrix} 4 \\ 1 \\ 5 \\ 9 \\ 3 \end{pmatrix}$$

The CSR storage of \mathbb{G}^T requires three arrays: ia (integer, size $nr + 1$), ja (integer, size nz) and mat (real, size nz). In the case of the small matrix above:

$$\begin{aligned} ia &= (1, 5, 8, 11) \\ ja &= (1, 3, 4, 5, 2, 3, 4, 1, 4, 5) \\ mat &= (1, 4, 1, 2, 1, 1, 1, 3, 7, 1) \end{aligned}$$

The sparse matrix vector multiplication kernel SpMV for $\vec{y} = \mathbf{A} \cdot \vec{x}$ has been explained above, and it is trivial to carry out this algorithm by hand and verify that the vector y is given by $y^T = (8, 3, 11)$.

Let us now turn to an interesting problem. Is it possible with the same arrays ia, ja, mat to compute the multiplication of the transpose of the matrix with a vector? The answer is of course positive and the code is given hereunder:

```
y=0.d0
do i = 1, nr
  do k=ia(i), ia(i+1)-1
    y(ja(k))=y(ja(k))+mat(k)*x(i)
  end do
end do
```

Let us take $i = 1$. The variable k then goes from 1 to 4. The inner loop does:

```
y(1)=y(1)+mat(1)*x(1)
y(3)=y(3)+mat(2)*x(1)
y(4)=y(4)+mat(3)*x(1)
y(5)=y(5)+mat(4)*x(1)
```

Let us take $i = 2$. The variable k then goes from 5 to 7. The inner loop does:

```
y(2)=y(2)+mat(5)*x(2)
y(3)=y(3)+mat(6)*x(2)
y(4)=y(4)+mat(7)*x(2)
```

Let us take $i = 3$. The variable k then goes from 8 to 10. The inner loop does:

```
y(1)=y(1)+mat(8)*x(3)
y(4)=y(4)+mat(9)*x(3)
y(5)=y(5)+mat(10)*x(3)
```

So in total, we have:

```

y(1)=mat(1)*x(1)+mat(8)*x(3)
y(2)=mat(5)*x(2)
y(3)=mat(2)*x(1)+mat(6)*x(2)
y(4)=mat(3)*x(1)+mat(7)*x(2)+mat(9)*x(3)
y(5)=mat(4)*x(1)+mat(10)*x(3)

```

which is indeed the result of the transposed of the matrix multiplied by a vector \vec{x} .

Let us consider a simple matrix \mathbb{K} which is square (size is 5×5):

$$\mathbb{K} = \begin{pmatrix} 1 & 0 & 4 & 1 & 2 \\ 0 & 1 & 0 & 1 & 0 \\ 4 & 0 & 0 & 7 & 1 \\ 1 & 1 & 7 & 4 & 0 \\ 2 & 0 & 1 & 0 & 5 \end{pmatrix}$$

In this case , NZ=16.

```

ia  = (1,5,7,10,14,17)
ja  = (1,3,4,5, 2,4, 1,4,5, 1,2,3,4, 1,3,5)
mat = (1,4,1,2,1,1,4,7,1,1,1,7,4,2,1,5)

```

The sparse matrix vector multiplication kernel SpMV for $\vec{y} = \mathbf{A} \cdot \vec{x}$ is given as follows in its simplest form. Since the matrix is symmetric, there is no use to store the whole matrix. Its upper half (for instance) will do. In this case, NZ= and then

```

ia  = (1,5,7,9,10,11)
ja  = (1,3,4,5, 2,4, 4,5, 4, 5)
mat = (1,4,1,2, 1,1, 7,1, 4, 5)

```

All is good and well until one wishes to multiply the real matrix by a vector. The SpMV routines described above will not work since it will return the upper half of the matrix multiplied by the vector.

One can then write a dedicated algorithm:

```

do i = 1,nr

! multiply the upper half by the vector

do k=ia(i), ia(i+1)-1
  y(i) = y(i) + mat(k)*x(ja(k))
end do

! multiply the transpose of matrix by vector
! but omit diagonal

do k=ia(i), ia(i+1)-1
  if (i/=ja(k)) then
    y(ja(k))=y(ja(k))+mat(k)*x(i)
  end if
end do

```


end do

end do

Example:

```
y(1)
=y(1) + mat(1)*x(ja(1)) + mat(2)*x(ja(2)) + mat(3)*x(ja(3)) + mat(4)*x(ja(4))
=y(1) + mat(1)*x(1) + mat(2)*x(3) + mat(3)*x(4) + mat(4)*x(5)
```

etc ...

Finish?

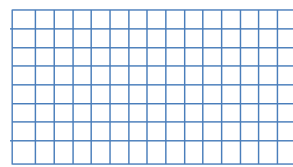
9.6 Mesh generation

meshes.tex

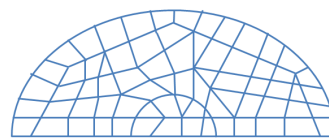
Before basis functions can be defined and PDEs can be discretised and solved we must first tessellate the domain with polygons, e.g. triangles and quadrilaterals in 2D, tetrahedra, prisms and hexahedra in 3D.

When the domain is itself simple (e.g. a rectangle, a sphere, ...) the mesh (or grid) can be (more or less) easily produced and the connectivity array filled with straightforward algorithms [1259]. However, real life applications can involve extremely complex geometries (e.g. a bridge, a human spine, a car chassis and body, etc ...) and dedicated algorithms/software must be used (see [1267, 419, 1374]).

We usually distinguish between two broad classes of grids: structured grids (with a regular connectivity) and unstructured grids (with an irregular connectivity).



Structured Grid



Unstructured Grid

Remark. Various families of so-called meshless methods exist and are commonly employed in Computational Fluid Dynamics [797, 783, 796, 798]. They are however very rarely used in Computational geodynamics, with a noticeable exception [530].

9.6.1 Quadrilateral-based meshes

Let us now focus on the case of a rectangular computational domain of size $L_x \times L_y$ with a regular mesh composed of $n_{elx} \times n_{ely} = n_{el}$ quadrilaterals. There are then $n_{nx} \times n_{ny} = n_{np}$ grid points. The elements are of size $h_x \times h_y$ with $h_x = L_x / n_{elx}$.

We have no reason to come up with an irregular/illogical node numbering so we can number nodes row by row or column by column as shown on the example hereunder of a 3×2 grid:

```

8=====9=====10=====11
|         |         |         |
|  (3)   |  (4)   |  (5)   |
|         |         |         |
4=====5=====6=====7
|         |         |         |
|  (0)   |  (1)   |  (2)   |
|         |         |         |
0=====1=====2=====3

```

"row by row"

```

2=====5=====8=====11
|         |         |         |
|  (1)   |  (3)   |  (5)   |
|         |         |         |
1=====4=====7=====10
|         |         |         |
|  (0)   |  (2)   |  (4)   |
|         |         |         |
0=====3=====6=====9

```

"column by column"

The numbering of the elements themselves could be done in a somewhat chaotic way but we follow the numbering of the nodes for simplicity. The row by row option is the adopted one in fieldstone and the coordinates of the points are computed as follows:

```

x = np.empty(nnp, dtype=np.float64)
y = np.empty(nnp, dtype=np.float64)
counter = 0
for j in range(0, nny):
    for i in range(0, nnx):
        x[counter] = i * hx
        y[counter] = j * hy
        counter += 1

```

The inner loop has i ranging from 0 to $nnx-1$ first for $j=0, 1, \dots$ up to $nny-1$ which indeed corresponds to the row by row numbering.

We now turn to the connectivity. As mentioned before, this is a structured mesh so that the so-called connectivity array, named **icon** in our case, can be filled easily. For each element we need to store the node identities of its vertices. Since there are **nel** elements and **m=4** corners, this is a $m \times nel$ array. The algorithm goes as follows:

```

icon = np.zeros((m, nel), dtype=np.int16)
counter = 0
for j in range(0, nely):
    for i in range(0, nelx):
        icon[0, counter] = i + j * nnx
        icon[1, counter] = i + 1 + j * nnx
        icon[2, counter] = i + 1 + (j + 1) * nnx
        icon[3, counter] = i + (j + 1) * nnx
        counter += 1

```

In the case of the 3×2 mesh, the **icon** is filled as follows:

| element id→ | 0 | 1 | 2 | 3 | 4 | 5 |
|-------------|---|---|---|---|----|----|
| node id↓ | | | | | | |
| 0 | 0 | 1 | 2 | 4 | 5 | 6 |
| 1 | 1 | 2 | 3 | 5 | 6 | 7 |
| 2 | 5 | 6 | 7 | 9 | 10 | 11 |
| 3 | 4 | 5 | 6 | 8 | 9 | 10 |

It is to be understood as follows: element #4 is composed of nodes 5, 6, 10 and 9. Note that nodes are always stored in a counter clockwise manner, starting at the bottom left. This is very important since the corresponding basis functions and their derivatives will be labelled accordingly.

In three dimensions things are very similar. The mesh now counts $nelx \times nely \times nelz = nel$ elements which represent a cuboid of size $Lx \times Ly \times Lz$. The position of the nodes is obtained as follows:

```

x = np.empty(nnp, dtype=np.float64)
y = np.empty(nnp, dtype=np.float64)
z = np.empty(nnp, dtype=np.float64)
counter = 0
for i in range(0, nnx):
    for j in range(0, nny):
        for k in range(0, nnz):
            x[counter] = i * hx
            y[counter] = j * hy
            z[counter] = k * hz
            counter += 1

```

The connectivity array is now of size $m \times nel$ with $m=8$:

```

icon = np.zeros((m, nel), dtype=np.int16)
counter = 0
for i in range(0, nelx):
    for j in range(0, nely):

```

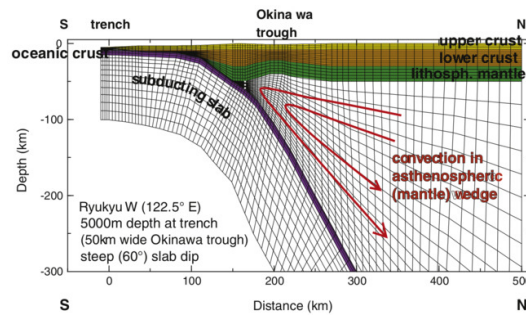
```

for k in range(0,nelz):
    icon[0,counter]=nny*nnz*(i )+nnz*(j )+k
    icon[1,counter]=nny*nnz*(i+1)+nnz*(j )+k
    icon[2,counter]=nny*nnz*(i+1)+nnz*(j+1)+k
    icon[3,counter]=nny*nnz*(i )+nnz*(j+1)+k
    icon[4,counter]=nny*nnz*(i )+nnz*(j )+k+1
    icon[5,counter]=nny*nnz*(i+1)+nnz*(j )+k+1
    icon[6,counter]=nny*nnz*(i+1)+nnz*(j+1)+k+1
    icon[7,counter]=nny*nnz*(i )+nnz*(j+1)+k+1
    counter += 1

```

produce drawing of node numbering

Although it is not very common in geosciences, quadrilateral meshes are sometimes employed in a boundary-fitted way, as shown hereunder:

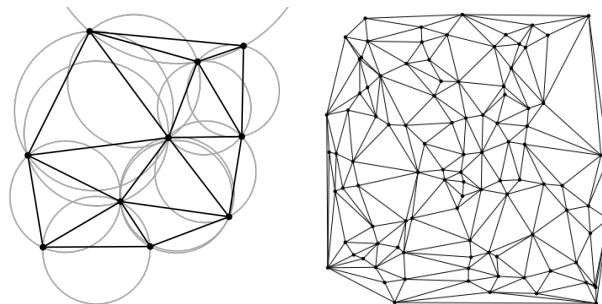


Taken from Gutscher et al. [517] (2016).

 Relevant Literature: Joun and Lee [657] (1997)

9.6.2 Delaunay triangulation and Voronoi cells, and triangle-based meshes

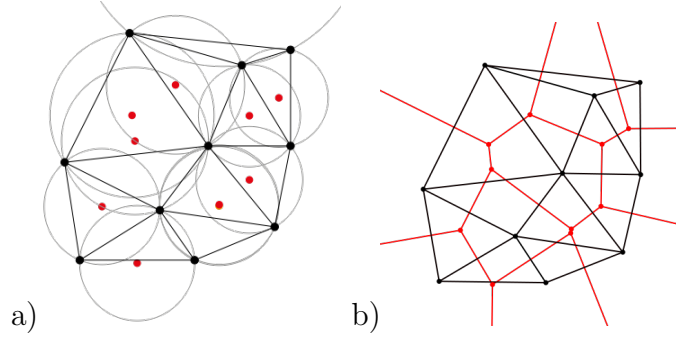
The topic of Delaunay⁹ triangulation is vast, but a simple definition can be written as follows: "a Delaunay triangulation for a set P of points in a plane is a triangulation $DT(P)$ such that no point in P is inside the circumcircle of any triangle in $DT(P)$." [wikipedia] Other properties of such triangulations are that they maximize the minimum angle of all the angles of the triangles in the triangulation. Note that for four or more points on the same circle (e.g., the vertices of a rectangle) the Delaunay triangulation is not unique and that points on a line also cannot yield a valid triangulation (for the simple reason that they do not form a triangle).



a) A Delaunay triangulation in the plane with circumcircles shown. b) The Delaunay triangulation of a random set of 100 points in a plane.


The Delaunay triangulation of a discrete point set P in general corresponds to the dual graph of the Voronoi diagram for P . A Voronoi diagram is composed of non-overlapping Voronoi cells which make a partition of the plane. For each point there is a corresponding region consisting of all points closer to that point than to any other: this region is the Voronoi cell of that point.

⁹The triangulation is named after Boris Delaunay for his work on this topic from 1934.

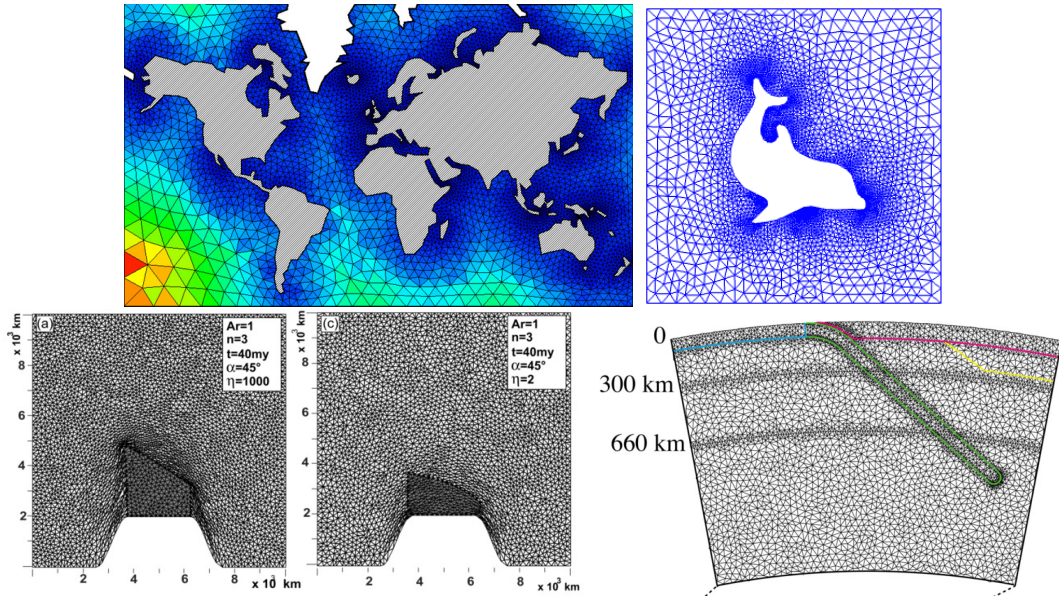


a) The Delaunay triangulation with all the circumcircles and their centers (in red). b) Connecting the centers of the circumcircles produces the Voronoi diagram (in red).

The Delaunay triangulation is used in the DOUAR code which is based on a particle levelset function to track materials. These particles are connected by means of a Delaunay triangulation (usually in a plane at startup, and then in a local Euclidean geometry once the surface is deformed) [136].

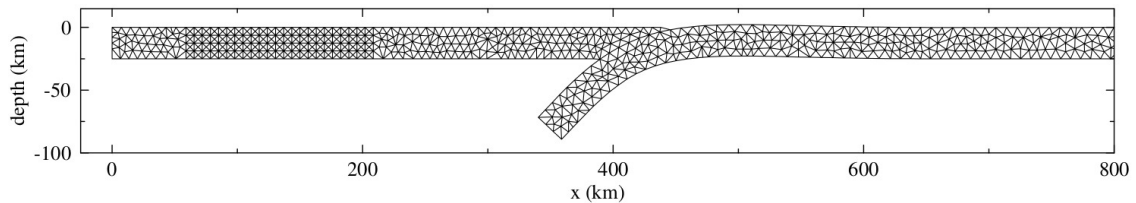
 Relevant Literature: [444].

Once a Delaunay triangulation has been obtained it can be used as a FEM mesh. Triangle-based meshes are obviously better suited for simulations of complex geometries:



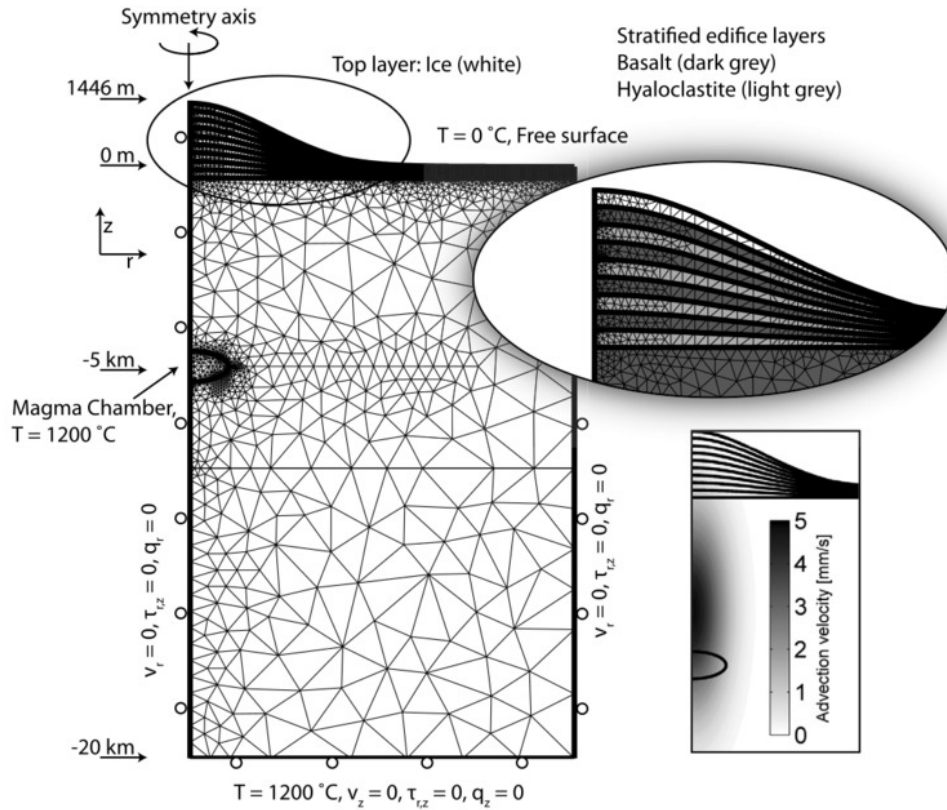
Bottom row. Left: Robl and Stüwe [1081] (2005), Right: Gerault, Becker, Kaus, Faccenna, Moresi, and Husson [446] (2012).

A very practical 2D triangle mesher is the code *Triangle*¹⁰ written by J.R. Shewchuk [1157, 1158, 1159]. Triangle is specialized for creating two-dimensional finite element meshes, but can also perform simpler related tasks such as forming Delaunay triangulations under various assumptions. Another very common mesher tool is Gmsh [457].

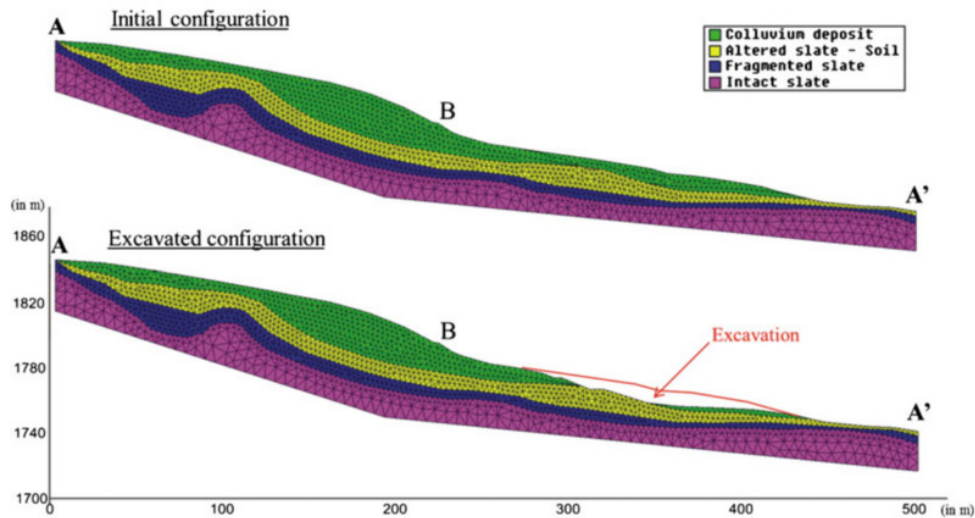


Taken from Buiter *et al.* [162]. Finite element grid. The subducting plate initially extends to 1226 km in the horizontal direction and is not completely shown here. Discretization in the subducting plate is slightly coarser towards the right edge.

¹⁰<https://www.cs.cmu.edu/~quake/triangle.html>

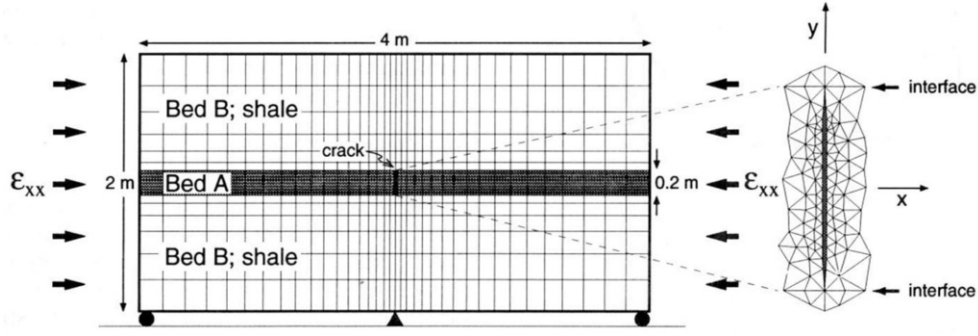


Taken from Bakker, Frehner, and Lupi [39] (2016). Numerical model setup of the 2D axisymmetric half-space with all applied boundary conditions to study the effects of ice-cap unloading on shallow volcanic systems.



Taken from Fernández-Merodo, García-Davalillo, Herrera, Mira, and Pastor [391] (2014). Modelling of slow landslides. Finite element mesh in the initial and excavated configuration.

Although it is rarely used in practice it is possible to produce meshes which contain both quadrilateral and triangular elements:




Taken from Fischer, Gross, Engelder, and Greenfield [395]. Mesh used to analyse the stress distribution around a pressurized crack in a layered elastic medium.

Mesh quality In Cioncolini and Boffi [258] (2019) the authors check the mesh quality of their triangulation by computing the following measures per element (they also refer to Field [394]):

$$q_1 = \frac{(b + c - a)(c + a - b)(a + b - c)}{abc}$$

$$q_2 = \frac{4\sqrt{3}A_T}{a^2 + b^2 + c^2}$$

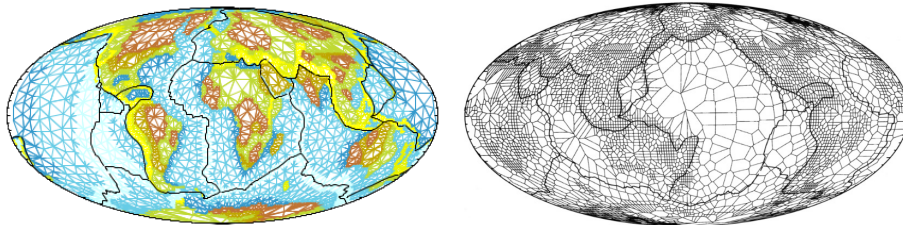
where a , b , c are the triangle side lengths and A_T is the triangle area. An equilateral triangle has $q_1 = q_2 = 1$ while a degenerate, zero area triangle has $q_1 = q_2 = 0$. As a rule-of-thumb, in a good quality mesh all triangles should have q_1 , q_2 above about 0.4-0.5.

 **Relevant Literature:** E. Mulyukova, B. Steinberger, M. Dabrowski, and S.V. Sobolev. “Survival of LLSVPs for billions of years in a vigorously convecting mantle: Replenishment and destruction of chemical anomaly”. In: *J. Geophys. Res.* 120 (2015), pp. 3824–3847. DOI: 10.1002/2014JB011688 Mirko Velić, Dave May, and Louis Moresi. “A fast robust algorithm for computing discrete voronoi diagrams”. In: *Journal of Mathematical Modelling and Algorithms* 8.3 (2009), pp. 343–355. DOI: 10.1007/s10852-008-9097-6

Remark. *The Natural Neighbour Interpolation method of Sambridge et al. [1100, 1099] is based on the Delaunay triangulation.*

Remark. *Moresi & Mather [904] have released Stripy, a Python module for (constrained) triangulation in Cartesian coordinates and on a sphere, which is based on Stripack [1060, 1061].*

write about gmesh



Taken from Gudmundsson & Sambridge (1998) [500]. Boundaries of Voronoi cells around 4100 of the original 16,200 2x2 degree cells selected to sample the details of the regionalization.

9.6.3 Tetrahedra

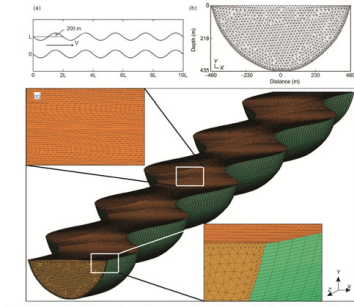
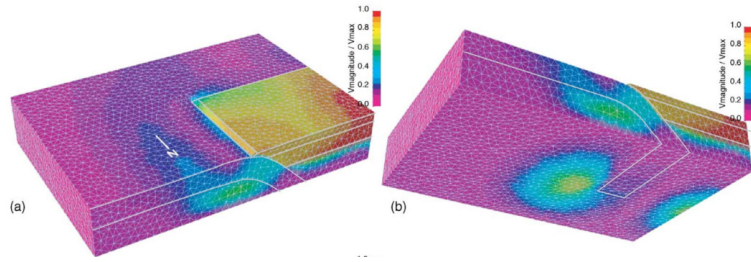
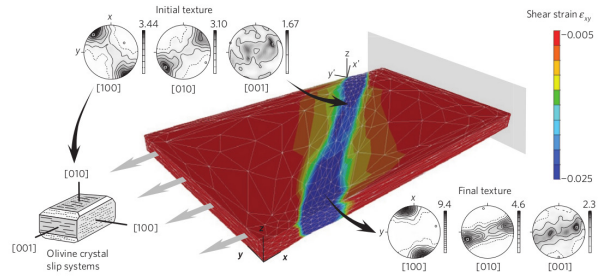
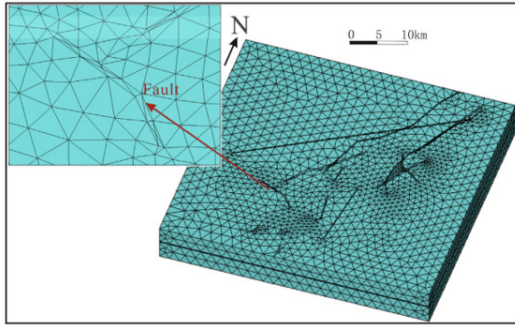


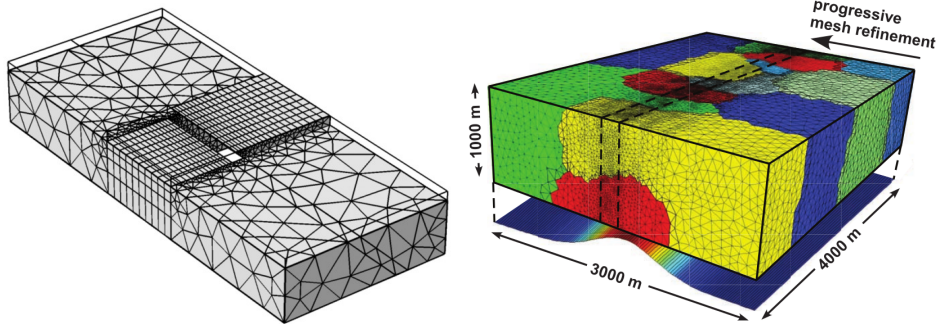
Figure 4 The top view of the initial geometric model (a), the cross section of the initial geometric model (b), the initial geometric model and its grid (c).



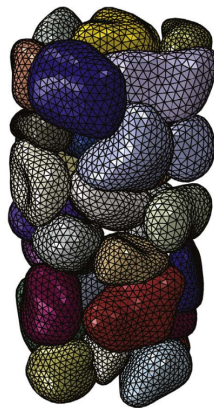
Left: Example of 3D mesh Yang and Shi [1379] (2015); Right: Normalized velocities of a STEP subduction model Govers and Wortel [480] (2005).



Left: 3D finite element grid in Damintun area, including prescribed faults, Guo, Yao, and Ren [509] (2016); Right: Structural reactivation in plate tectonics controlled by olivine crystal anisotropy, Tommasi, Knoll, Vauchez, Signorelli, Thoraval, and Logé [1271] (2009) - based on $\mathbb{P}_1 \times P_1$ elements.



Left: Mesh used for the three-dimensional model. A high resolution mesh is used in the wedge and subslab domains, while the mesh resolution decays to lower values toward the edge of the model. All elements are quadratic, allowing for twice the resolution visualized here, Paczkowski, Montési, Long, and Thissen [969] (2014); Right: Mid-Ocean Ridge Hydrothermal System: 3D mesh consisting of 2.5 m tetrahedron elements. Resolution is refined toward the axial center, with the finest resolution between the dashed lines, and colors indicate computational domains assigned to separate processors, Coumou, Driesner, and Heinrich [280] (2008).

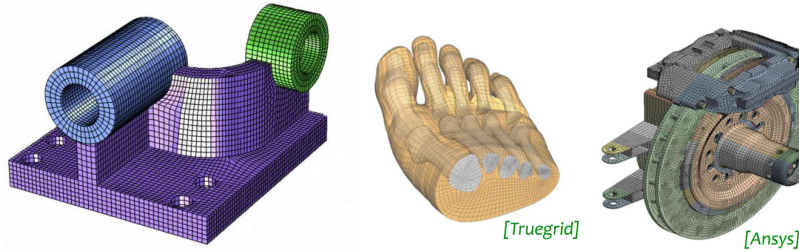



Grains of sand before a compression experiment with FEM Imseeh and Alshibli [621] (2018).

Check TetGen mesher Hang Si. “TetGen, a Delaunay-based quality tetrahedral mesh generator”. In: *ACM Transactions on Mathematical Software (TOMS)* 41.2 (2015), pp. 1–36. DOI: 10.1145/2629697.

9.6.4 Hexahedra

A hexahedron is a convex polytope isomorphic to the cube $[0, 1]^3$. Edges are line segments, facets are strictly **planar** convex polygons.



 Relevant Literature Efficient Volume computation for Three- Dimensional hexahedral Cells [349, 481]

9.6.5 Adaptive Mesh Refinement

Let us do a simple calculation and assume we wish to model mantle convection on Earth. The inner radius is $R_1 = 3485$ km and the bottom of the lithosphere is at $R_2 = 6250$ km. The volume of fluid is then

$$V = \frac{4}{3}\pi(R_2^3 - R_1^3) \simeq 8.5 \times 10^{11} \text{ km}^3$$

Let us further assume that we are satisfied with an average resolution of 10 km. Each element/cell is then 10^3 km^3 and the total number of elements/cell is then

$$N \simeq 8.5 \times 10^8 \sim \mathcal{O}(10^9)$$

This is a very large number. The resulting linear systems from the discretisation of the equations on such a mesh will be very even larger for the Stokes equations and solving these systems will require *very* large numbers of CPUs and long compute times.

Aside from these considerations it is quite obvious that a high resolution mesh is not needed in parts of the mantle where large scale upwellings and downwellings occur, but probably even higher resolution will be needed in the vicinity of thin plumes and boundary layers. This means that a uniform mesh is a sub-optimal way of discretising space for such problems.

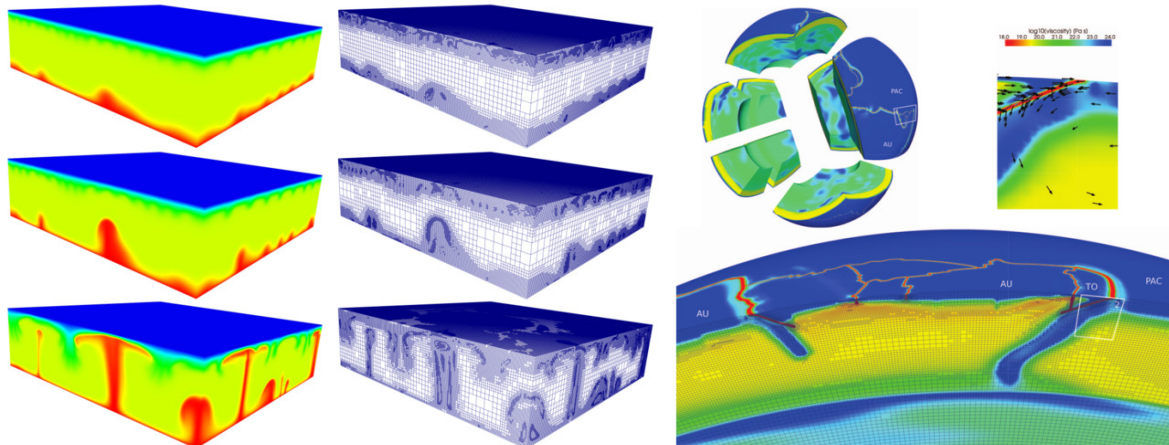
The same reasoning also holds in the lithosphere where for instance narrow plate boundaries need to be adequately resolved while the inside of rigid plates can be modelled with coarser meshes.

Finally, although one could employ meshing software to arrive at well balanced meshes in space, the dynamic character of the geodynamics modelling renders this approach cumbersome. A subduction zone, a mid-ocean rift or an ascending plume will evolve in time and the mesh will have to evolve in time too.

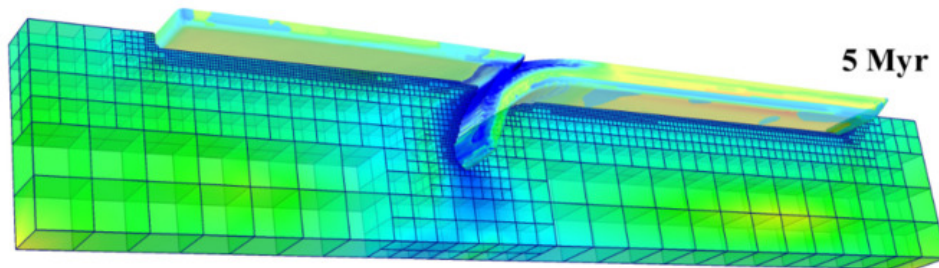
In light of all this, it was only a matter of time before Adaptive Mesh Refinement was adopted in computational geodynamics. However, since the use and update of such meshes is somewhat complex in terms of numerical algorithms, its introduction came somewhat late (00’s and later). The DOUAR code (see Section ??) developed originally by J. Braun and Ph. Fullsack is a prime example of an early multi-purpose code relying on a self-written Octree library [136]. More recently the ASPECT code was developed on top of the Octree library p4est [191]. Note the 2007 and 2008

papers by Davies et al [307, 310] which explore adaptive mesh refinement with the ConMan code (see Appendix ??).

For further reading I suggest you read the review by May, Schellart & Moresi on this topic [847].



Taken from [190] and [192]

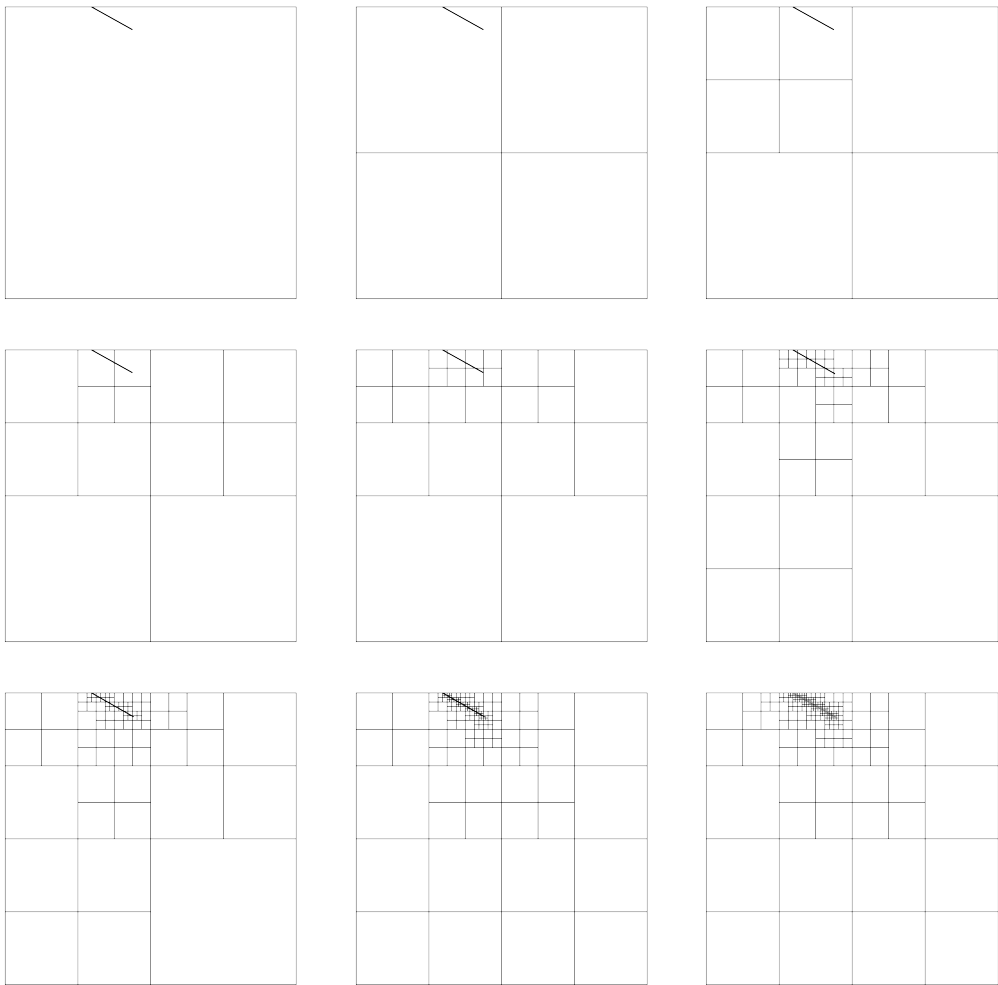


Taken from [467]

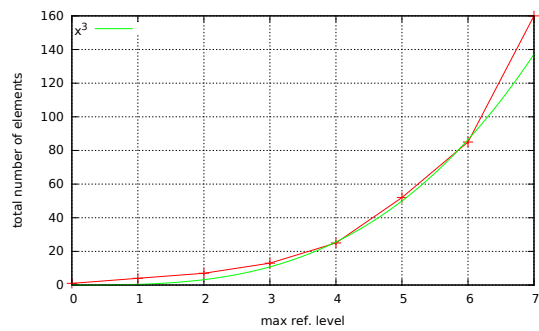
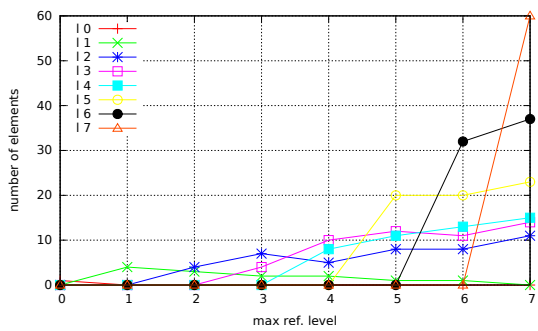
Relevant Literature:

- C. Burstedde et al. “Scalable Adaptive Mantle Convection Simulation on Petascale Supercomputers”. In: *ACM/IEEE SC Conference Series, 2008* (2008)
- C. Burstedde, O. Ghattas, G. Stadler, T. Tu, and L.C. Wilcox. “Parallel scalable adjoint-based adaptive solution of variable-viscosity Stokes flow problems”. In: *Computer Methods in Applied Mechanics and Engineering* 198 (2009), pp. 1691–1700. DOI: 10.1016/j.cma.2008.12.015
- Carsten Burstedde et al. “Extreme-scale AMR”. in: *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE Computer Society. 2010, pp. 1–12. DOI: 10.1109/SC.2010.25
- W. Leng and S. Zhong. “Implementation and application of adaptive mesh refinement for thermochemical mantle convection studies”. In: *Geochem. Geophys. Geosyst.* 12.4 (2011). DOI: 10.1029/2010GC003425
- Y. Mishin. “Adaptive multiresolution methods for problems of computational geodynamics”. PhD thesis. ETH Zurich, 2011
- K. Sverdrup, N. Nikiforakis, and A. Almgren. “Highly parallelisable simulations of time-dependent viscoplastic fluid flow simulations with structured adaptive mesh refinement”. In: *Physics of Fluids* 30 (2018), p. 093102. DOI: 10.1063/1.5049202
- Marc Fehling and Wolfgang Bangerth. “Algorithms for parallel generic hp-adaptive finite element software”. In: *ACM Transactions on Mathematical Software* 49.3 (2023), pp. 1–26. DOI: 10.1145/3603372

A short illustrative exercise .

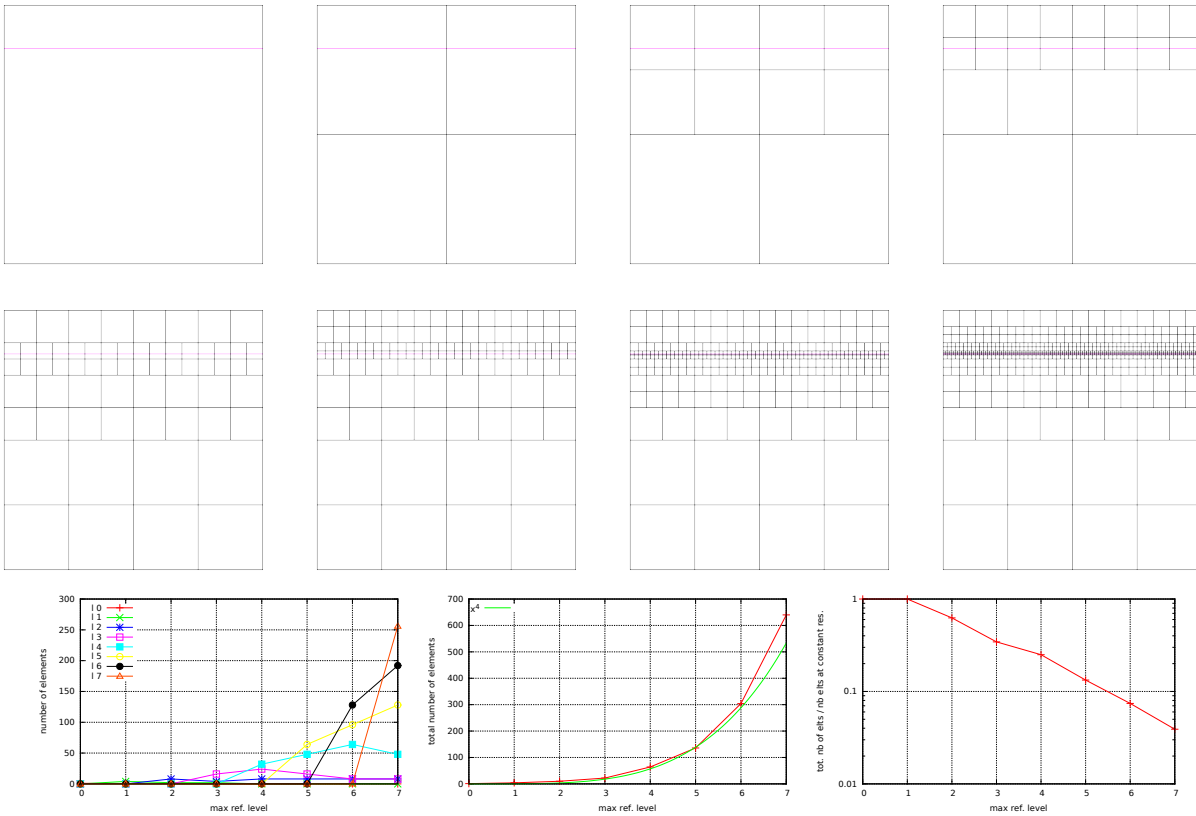


| | # 10 | # 11 | # 12 | # 13 | # 14 | # 15 | # 16 | # 17 | # 18 |
|--------------|------|------|------|------|------|------|------|------|------|
| max level= 0 | 1 | | | | | | | | |
| max level= 1 | 0 | 4 | | | | | | | |
| max level= 2 | 0 | 3 | 4 | | | | | | |
| max level= 3 | 0 | 2 | 7 | 4 | | | | | |
| max level= 4 | 0 | 2 | 5 | 10 | 8 | | | | |
| max level= 5 | 0 | 1 | 8 | 12 | 11 | 20 | | | |
| max level= 6 | 0 | 1 | 8 | 11 | 13 | 20 | 32 | | |
| max level= 7 | 0 | 0 | 11 | 14 | 15 | 23 | 37 | 60 | |
| max level= 8 | 0 | 0 | 11 | 13 | 17 | 27 | 43 | 72 | 116 |



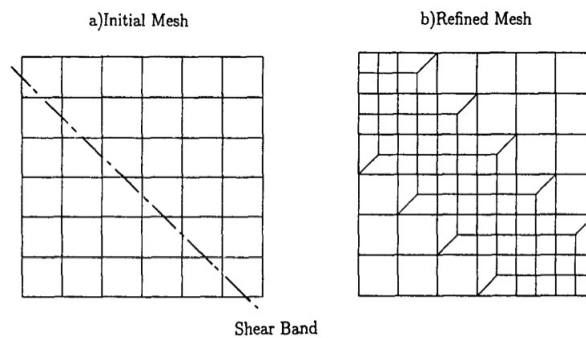
In the particular case presented here, even though the inclusion in a short two-dimensional line, the total number of elements grows faster than the third power of the refinement level. While of course the total number of elements remains much smaller than the constant resolution counterpart,

this observation tells us that authorising a unit increase of the maximum refinement level can have a substantial effect on the total number of elements.

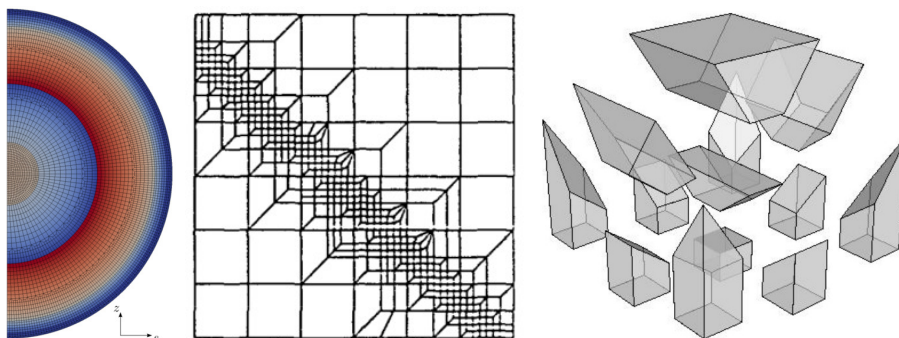


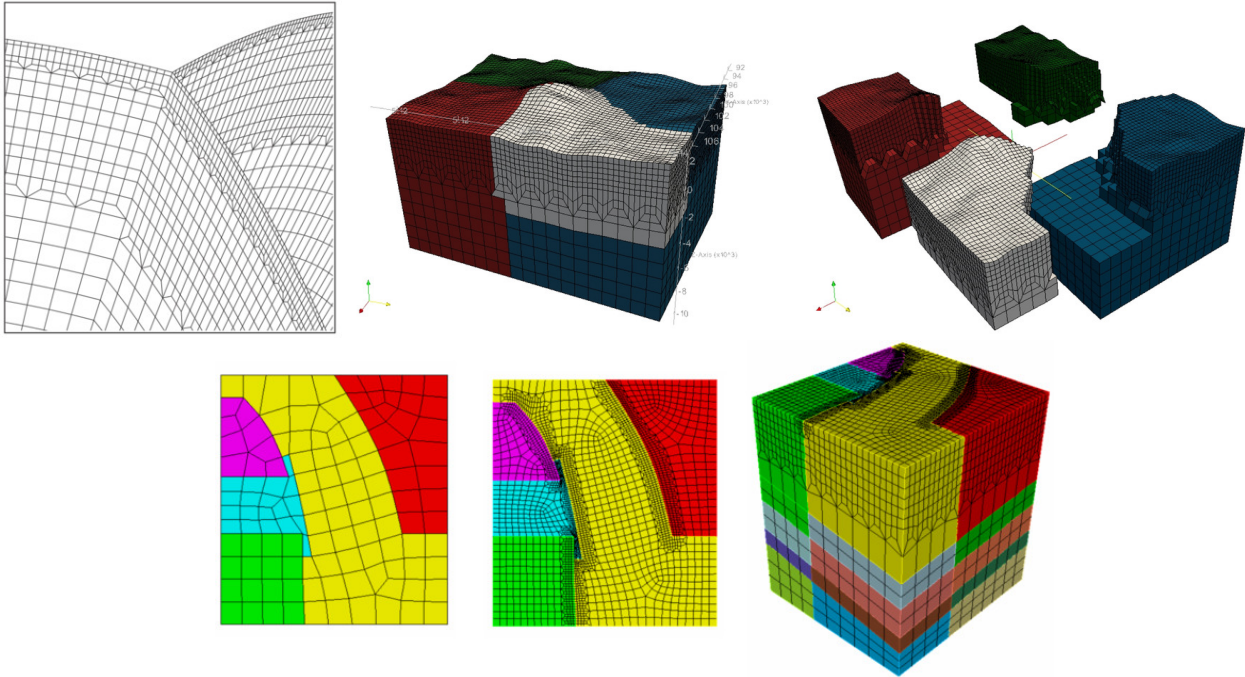
9.6.6 Conformal Mesh Refinement

The quadtree/octree mesh refinement presented above is one option when it comes to mesh refinement (or h -refinement). However their massive drawback is the presence of hanging nodes which require special attention. Another approach to mesh refinement is conformal mesh refinement as best exemplified on the following figures:

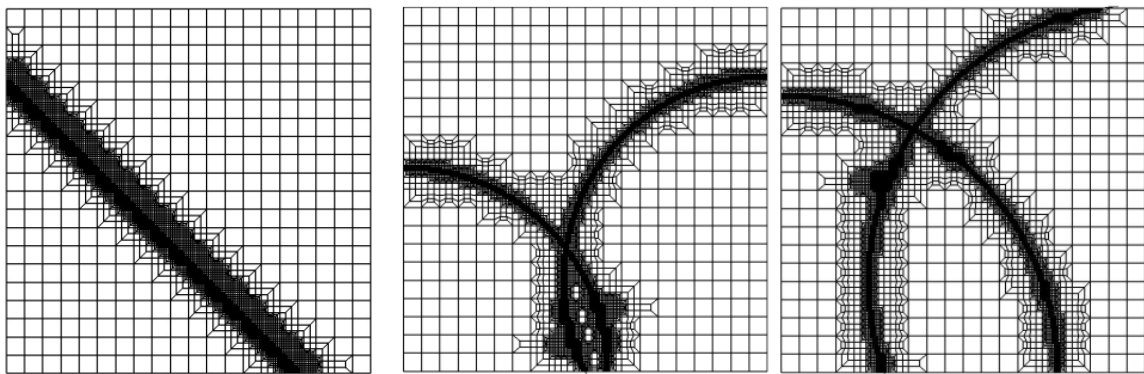


Taken from Deb *et al.* (1996) [324]. A typical instance of the outcome of the refinement procedure. Notice that the 'spill-over' is reduced to one row on each side of the 'localized' elements.

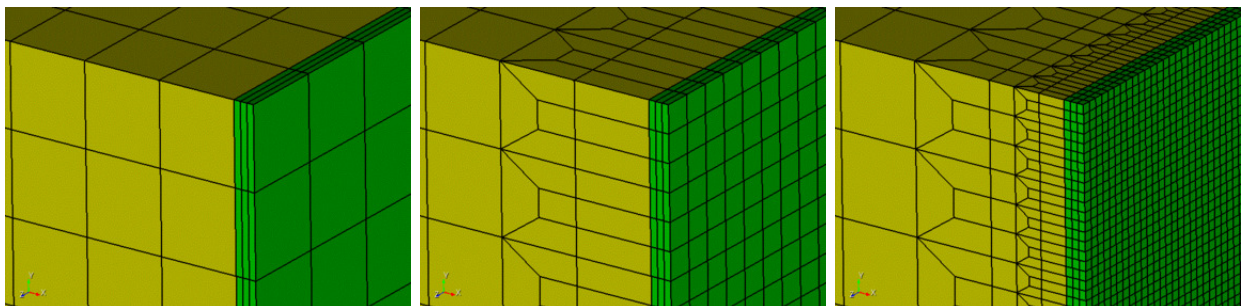
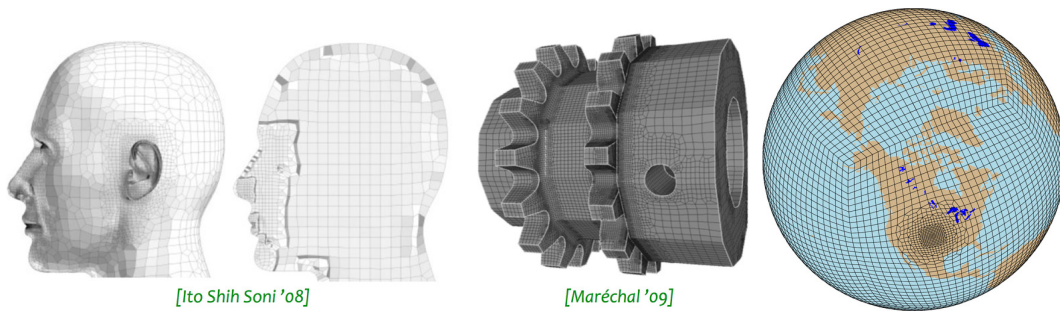




Top row, From left to right: van Driel *et al.* (2015) [346]; Deb *et al.* (1996) [324]; Harris *et al.* (2004) [549]; Komatitsch *et al.* (2005) [719]; Middle row: Specfem manual; Bottom row: I don't know anymore.



Taken from Garimella [436] (2009).



https://cubit.sandia.gov/public/14.0/help_manual/WebHelp/mesh_generation/mesh_modification/mesh_refinement.htm

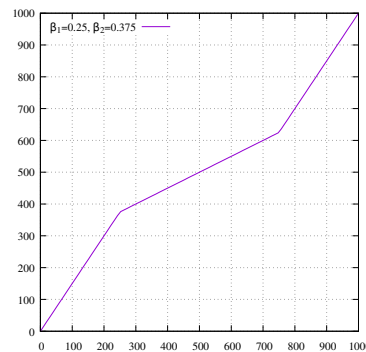
Relevant Literature:

- Düster & Rank [354],
- Harris *et al.* (2004) [549],
- Anderson *et al.* (2009) [20],
- Anderson [19],
- Garimella (2009) [436],
- Nicolas & Fouquet (2013) [940, 939].
- Parrish [979],
- Schneiders [1132, 1131, 1134, 1133],
- Schneiders *et al.* [1135],
- Staten & Canann [1194],
- book by Ramm *et al.* [1039].

9.6.7 Stretching the mesh

In some cases the topology of the mesh can be regular but one can for instance stretch the mesh such that (for instance) the vertical resolution is higher at the top than at the bottom, or higher in the middle than on the sides.

The idea behind the transformation is a piecewise-linear function which maps $[0,L]$ to $[0,L]$ where L is the length of the domain in the x -direction. For instance, this transformation can take the following form:



Parameters β_1 and β_2 control the shape of the lines.

The kinks in the line occur at $\beta_1 L$ and $(1 - \beta_1)L$ (see code here under).

The (minimal) code to transform the mesh is as follows:

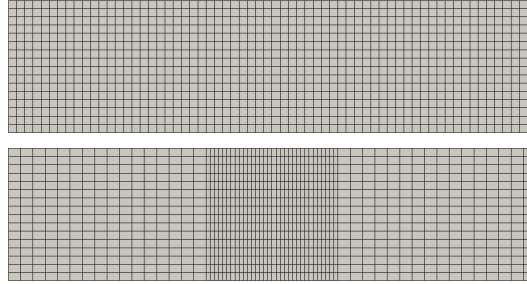
```
def stretch_towards_center(x,L,beta1,beta2):
    if x<beta1*L:
        val = beta2/beta1*x
    elif x<(1.-beta1)*L:
        val = (1-2*beta2)/(1-2*beta1)*(x-beta1*L)+beta2*L
    else:
        val=beta2/beta1*(x-(1-beta1)*L)+(1-beta2)*L
    return val
```

```
[...]

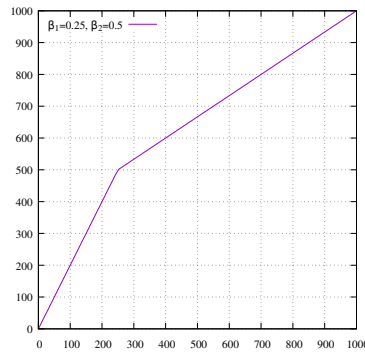
beta1=0.25
beta2=0.375

for i in range(0,NV):
    x[i]=stretch_towards_center(x[i],Lx,beta1,beta2)
```

The following meshes count 64x16 elements. The top one is a regular mesh, with square elements, while the second one has been stretched by means of the transformation above:



Concerning the stretching towards the top of the model domain, the transformation line is as follows:



Parameters β_1 and β_2 control the shape of the lines. The kinks in the line occur at $\beta_1 L$ and $(1 - \beta_1)L$.

The slope of the left line is $\beta_2/\beta_1 x$.

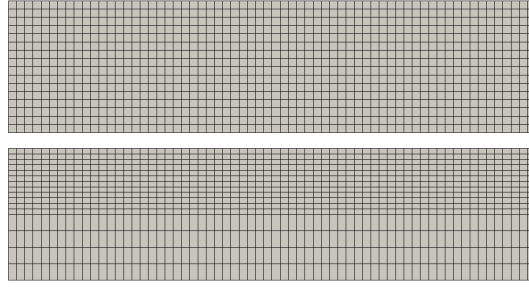
The (minimal) code to transform the mesh is as follows:

```
def stretch_towards_top(x,L,beta1,beta2):
    if x<beta1*L:
        val=beta2/beta1*x
    else:
        val=(1-beta2)/(1-beta1)*(x-beta1*L)+beta2*L
    return val

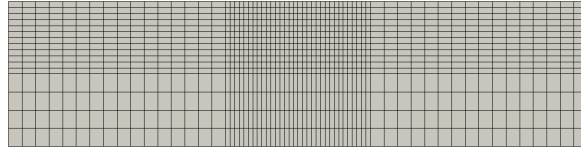
[...]

beta1=0.25
beta2=0.5
for i in range(0,NV):
    y[i]=stretch_towards_top(y[i],Ly,beta1,beta2)
```

The following meshes count 64x16 elements. The top one is a regular mesh, with square elements, while the second one has been stretched by means of the transformation above.

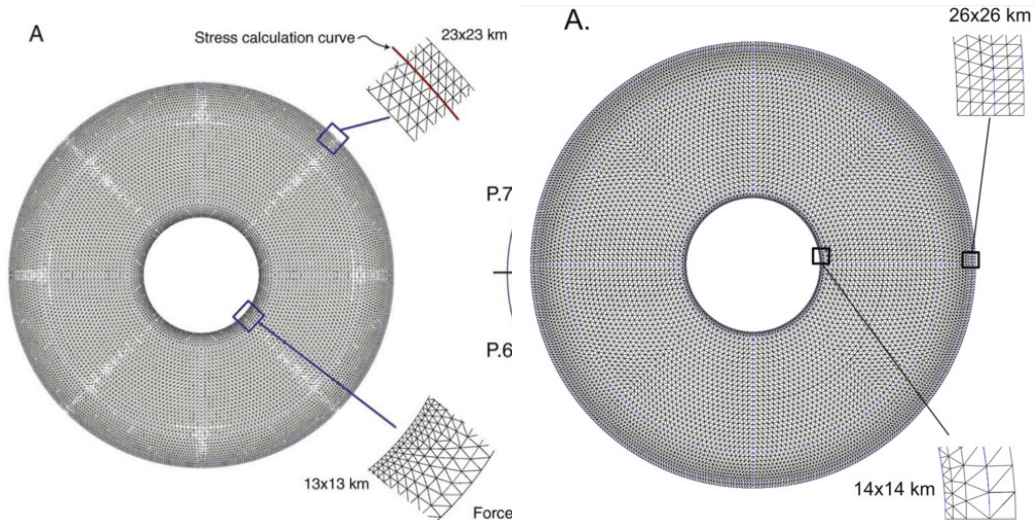


Finally both transformations can be applied to the same mesh:



This approach is used in [STONE 67](#).

9.6.8 Meshes in an annulus



The quadratic finite element mesh as used in Brandenburg *et al.* [133, 134]

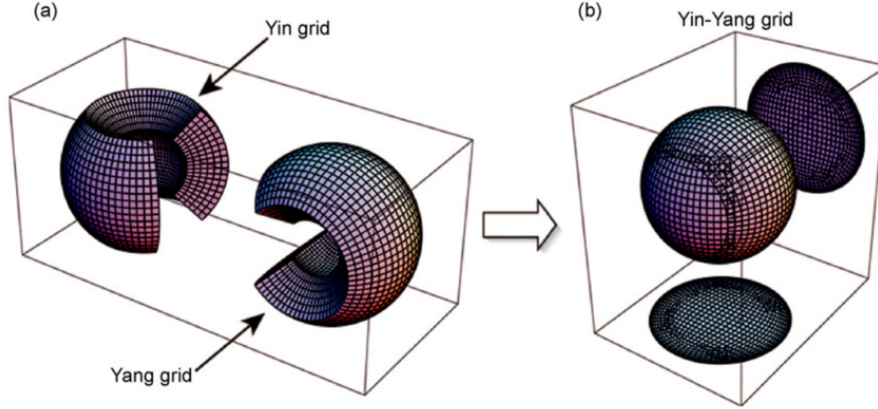
9.6.9 Meshes in/on a hollow sphere

The following is for the most part published in Thieulot (2018) Thieulot [1259] (2018).

To a first approximation the Earth is a sphere: the Earth’s polar diameter is about 43 kilometers shorter than its equatorial diameter, a negligible difference of about 0.3%. As a consequence, modelling physical processes which take place in the planet require the discretisation of a sphere. Furthermore, because core dynamics occur on vastly different time scales than mantle dynamics, mantle modelling usually leaves the core out, thereby requiring simulations to be run on a hollow sphere mesh (with the noticeable exception of Gerya and Yuen [450]).

Although so-called latitude-longitude grids would seem appealing, they suffer from the convergence of meridians at the poles (resulting in over sampling at poles) and the juxtaposition of triangles near the poles and quadrilaterals elsewhere. As a consequence more regular, but more complex, grids have been designed over the years which tessellate the surface of the sphere into triangles or quadrilaterals (sometimes overlapping). There is the ‘cubed sphere’ [1084, 564, 235, 1205, 236, 147, 1385], the Yin-Yang grid [663, 1387, 1389, 665, 1228, 283, 284], the Yin-Yang-zhong grid [554], the Yin-yang

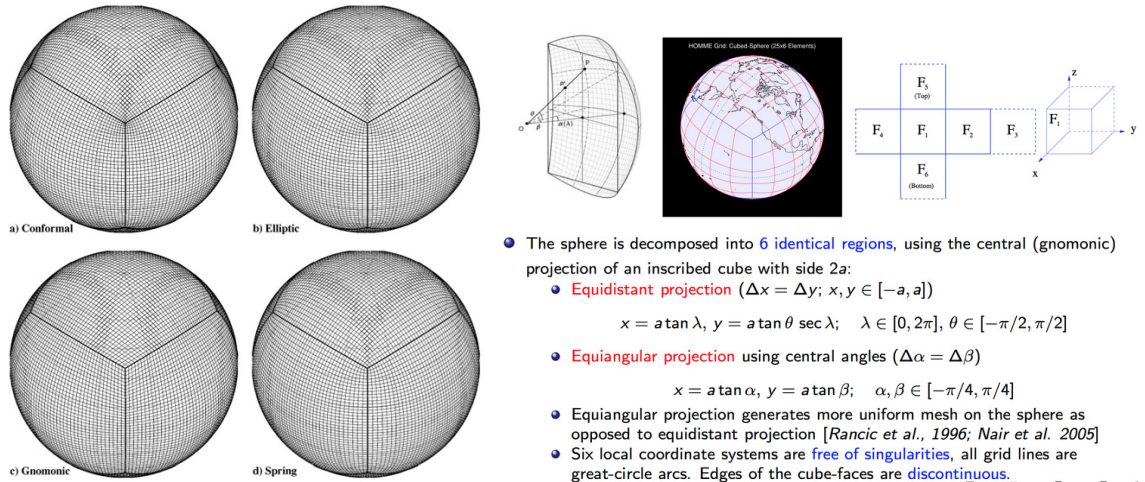
grid of Shahnas & Peltier [1151], the spiral grid [617], an icosahedron-based grid [58], or a grid composed of 12 blocks further subdivided into quadrilaterals [1414] as used in the CitcomS code. Note that [956] have also presented a method for generating a numerical grid on a spherical surface which allows the grid to be based on several different regular polyhedrons (including octahedron, cube, icosahedron, and rhombic dodecahedron). Ideally, one wishes to generate a mesh that is regular, i.e. angles between edges/faces as close to 90° as possible, of approximately similar volumes.



Example of Yin-Yang grid. Taken from Kameyama, Kageyamab, and Sato [665] (2008).

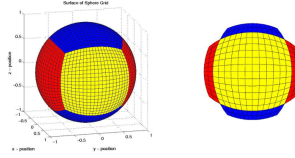
How such meshes are built is often not discussed in the literature. It is a tedious exercise of three-dimensional geometry and it can be time-consuming, especially the connectivity array generation. In Thieulot (2018) [1259] I present an open source mesh generator for three hollow sphere meshes: the 'cubed sphere' mesh, the CitcomS mesh and the icosahedral mesh:

- The cubed sphere ('HS06'), composed of 6 blocks which are themselves subdivided into $N_b \times N_b$ quadrilateral shaped cells [1095, 1084, 564, 189]. Four types of cubed spheres meshes have been proposed: the conformal, elliptic, gnomonic and spring types [1022]:

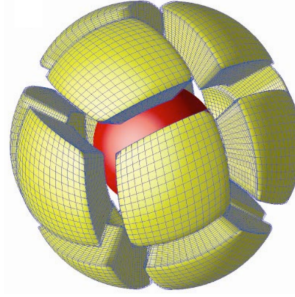


Left: The cubed-sphere grids at 2° resolution displaying cells on the sphere, the image focuses on the distribution of grid cells near one corner of the grid; (a) conformal mapping [1044, 850], (b) the gnomonic grid modified by elliptic solver, (c) equiangular gnomonic mapping and (d) the gnomonic grid modified by spring dynamics. [1022]. Right: Taken from presentation by R. Nair, see Nair2008.pdf

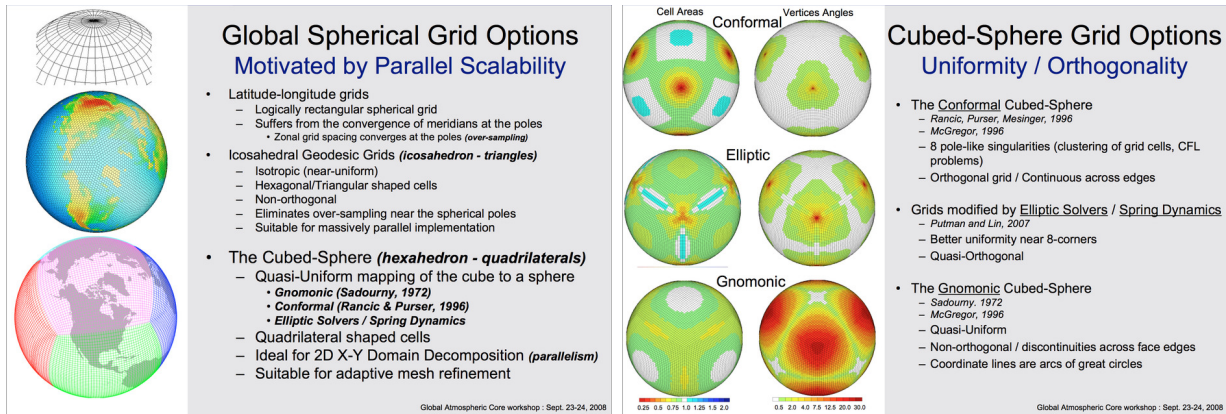
However only gnomonic meshes are considered in Thieulot (2018): these are obtained by inscribing a cube within a sphere and expanding to the surface of the sphere. The cubed sphere has been used in large-scale mantle convection simulation in conjunction with Adaptive Mesh Refinement [7, 189].



- The CitcomS mesh ('HS12') composed of 12 blocks also subdivided into $N_b \times N_b$ quadrilateral shaped cells [1414, 1205, 1412, 29]. Note that ASPECT [732, 560], a relatively new code aimed at superseding CitcomS can generate and use this type of mesh [1256] but is not limited to it.



- The icosahedral mesh ('HS20') composed of 20 triangular blocks [58, 57] subdivided into triangles, which is used in the TERRA code [174, 173, 172, 308].

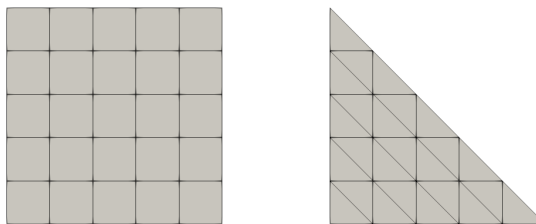


source?

Given the regularity and symmetry of these meshes determining the location of the mesh nodes in space is a relatively straightforward task. Building the mesh connectivity in an efficient manner is where the difficulty lies.

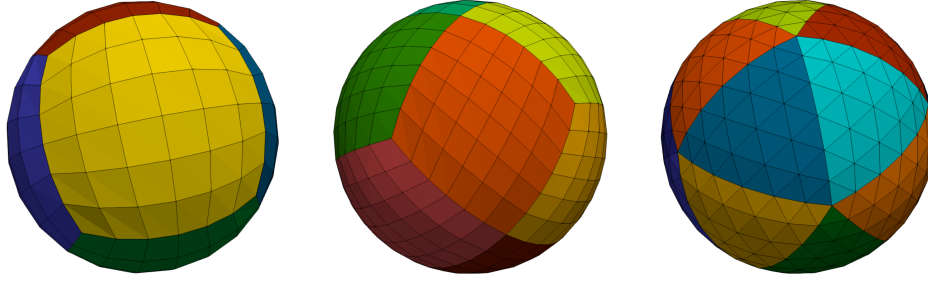
The approach to building all three meshes is identical:

1. A reference square or triangle is populated with cells, parametrised by a level l : the square is subdivided into $l \times l$ quadrilaterals while the triangle is subdivided into l^2 triangles.



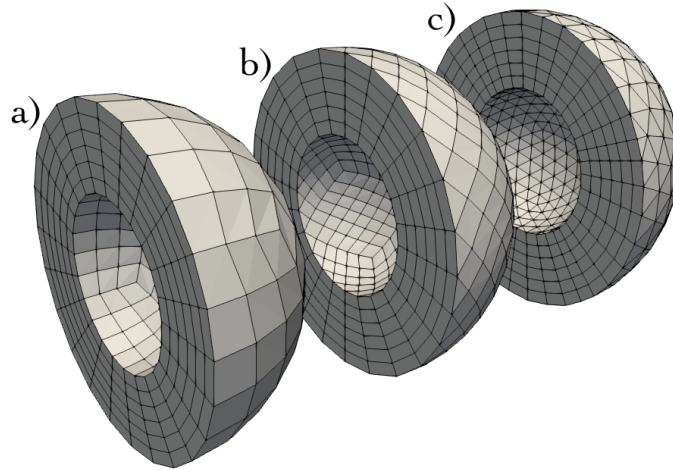
Reference square and triangles meshes at level 5.

2. This reference square or triangle is then replicated n_{block} times (6, 12 or 20) and mapped onto a portion of a unit sphere. The blocks are such that their union covers a full sphere but they cannot overlap except at the edges:



From left to right: HS06, HS12 and HS20 shells coloured by block number.

3. All block meshes are then merged together to generate a shell mesh. This task is rather complex as duplicate nodes must be removed and all connectivity arrays of the blocks must then be mended accordingly.
4. Shell meshes are replicated $n_{layer}+1$ times outwards with increasing radii. The n_{layer} shells are then merged together to form a hollow sphere mesh:



a) HS06 mesh composed of 6 blocks containing each 6^3 cells; b) HS12 mesh composed of 12 blocks containing each 6^3 cells; c) HS20 mesh composed of 20 blocks containing each 6^3 cells.

More information on these steps is available in the manual of the code. In the following table the number of nodes and cells for a variety of resolutions for all three mesh types is reported. Looking at the CitcomS literature of the past 20 years, we find that the mesh data presented in this table cover the various resolutions used, e.g. 12×48^3 [859, 29], 12×64^3 [165] 12×96^3 [166], 12×128^3 [65, 1348, 1349]. Note that in the case of the HS06 and HS12 meshes the mesh nodes are mapped out to the 6 or 12 blocks following either an equidistant or equiangle approach (see [1022] for details on both approaches).

| type | level | N | N_{el} | structure |
|------|-------|-------------|-------------|-------------------|
| HS06 | 2 | 78 | 48 | 6×2^3 |
| HS06 | 4 | 490 | 384 | 6×4^3 |
| HS06 | 8 | 3,474 | 3,072 | 6×8^3 |
| HS06 | 16 | 26,146 | 24,576 | 6×16^3 |
| HS06 | 32 | 202,818 | 196,608 | 6×32^3 |
| HS06 | 64 | 1,597,570 | 1,572,864 | 6×64^3 |
| HS06 | 128 | 12,681,474 | 12,582,912 | 6×128^3 |
| HS06 | 256 | 101,057,026 | 100,663,296 | 6×256^3 |
| HS12 | 2 | 150 | 96 | 12×2^3 |
| HS12 | 4 | 970 | 768 | 12×4^3 |
| HS12 | 8 | 6,930 | 6,144 | 12×8^3 |
| HS12 | 16 | 52,258 | 49,152 | 12×16^3 |
| HS12 | 32 | 405,570 | 393,216 | 12×32^3 |
| HS12 | 48 | 1,354,850 | 1,327,104 | 12×48^3 |
| HS12 | 64 | 3,195,010 | 3,145,728 | 12×64^3 |
| HS12 | 128 | 25,362,690 | 25,165,824 | 12×128^3 |
| HS12 | 256 | 202,113,538 | 201,326,592 | 12×256^3 |
| HS20 | 2 | 126 | 160 | 20×2^3 |
| HS20 | 4 | 810 | 1,280 | 20×4^3 |
| HS20 | 8 | 5,778 | 10,240 | 20×8^3 |
| HS20 | 16 | 43,554 | 81,920 | 20×16^3 |
| HS20 | 32 | 337,986 | 655,360 | 20×32^3 |
| HS20 | 64 | 2,662,530 | 5,242,880 | 20×64^3 |
| HS20 | 128 | 21,135,618 | 41,943,040 | 20×128^3 |
| HS20 | 256 | 168,428,034 | 335,544,320 | 20×256^3 |

Number of nodes N and elements/cells N_{el} for the three types of meshes and for various levels.

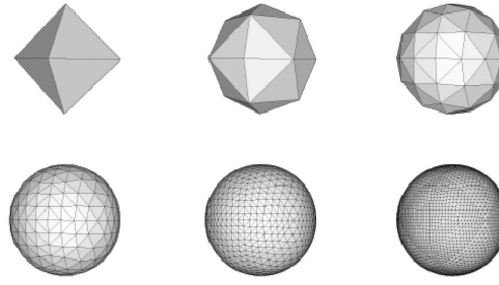
HS06: cubed sphere; HS12: CitcomS mesh; HS20: icosahedral mesh.

There are also many possibilities offered by the use of tetrahedral cells/elements:

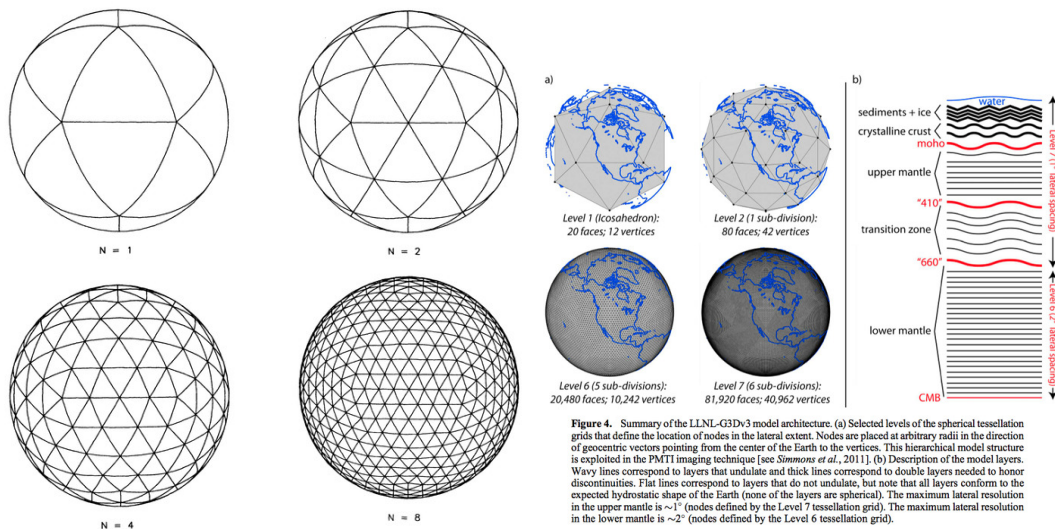


Left: Grid of a global neo-tectonic SHELLS model coupled to a global mantle circulation model; colours represent temperatures (red=hot, blue=cold) at a depth of 200km below the surface. Taken from Oeser *et al.* (2009) [952]. Right: Taken from the GeoTess software ¹¹ manual.

¹¹<https://www.sandia.gov/geotess/>



Example of a Hierarchical Triangular Mesh



Baumgardner and Frederickson [58] (1985), Simmons, Myers, Johannesson, and Matzel [1171] (2012)

Relevant Literature:

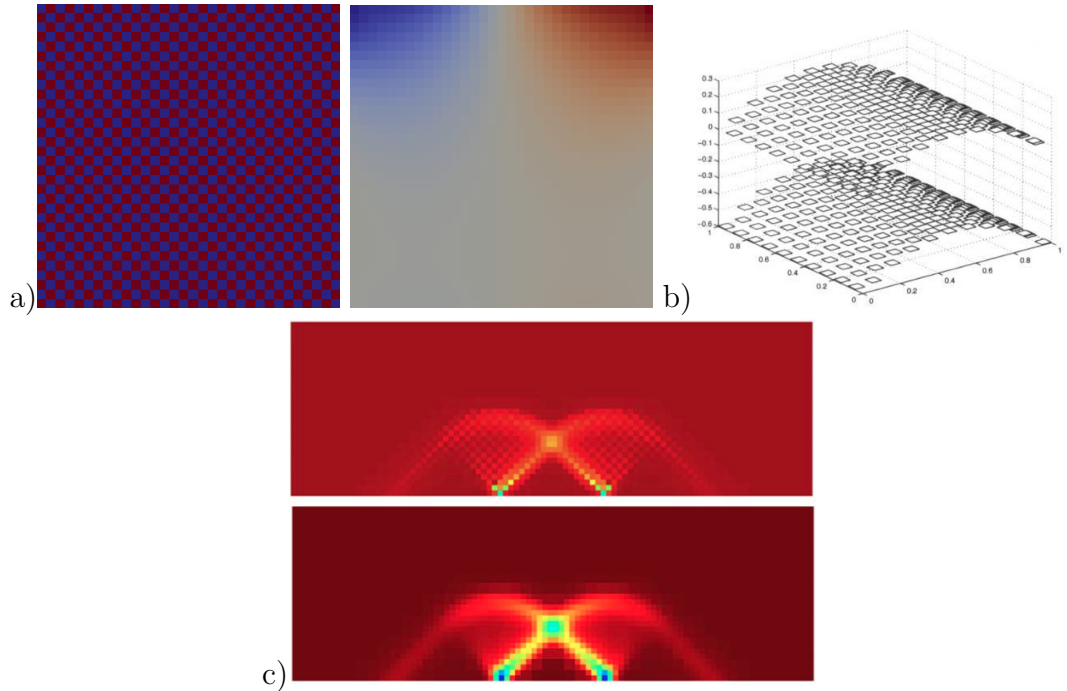
- Phillips, Davies, and Oldham [996] (2019) present an algorithm which builds polyhedral-based grids.
- Upadhyaya, Sharma, Mittal, and Fatima [1299] (2011) on icosahedral-hexagonal grids on a sphere for CFD applications.
- Saff and Kuijlaars [1096] (1997) on distributing many points on a sphere.
- Swinbank and Purser [1222] (2006) on Fibonacci grids which possess virtually uniform and isotropic resolution, with an equal area for each grid point.
- Hardin, Saff, et al. [547] (2004) on discretizing manifolds via Minimum Energy Points. Element Software

9.7 Pressure smoothing/filtering/recovery for $Q_1 \times P_0$ elements

It has been widely documented that the use of the $Q_1 \times P_0$ element is not without problems. Aside from the consequences it has on the FE matrix properties, we will here focus on another unavoidable side effect: the spurious pressure checkerboard modes.

These modes have been thoroughly analysed decades ago, see for instance Hughes, Liu, and Brooks [608] (1979), Sani, Gresho, Lee, and Griffiths [1108] and Sani, Gresho, Lee, Griffiths, and Engelman [1109] (1981), Griffiths and Silvester [495] (1994). They can be filtered out (Chen, Pan, and Chang [221] (1995)) or simply smoothed (Lee, Gresho, and Sani [759] (1979)), as we will see later. Nodes on edges and corners may need special treatment as documented in Sani, Gresho, Lee, and Griffiths [1108] (1981) or Lee, Gresho, and Sani [759] (1979). The list of 8 schemes is not exhaustive with regards to the above mentioned publications. There has been considerable amount of work on the topic and this section is unfortunately not representing the literature appropriately.

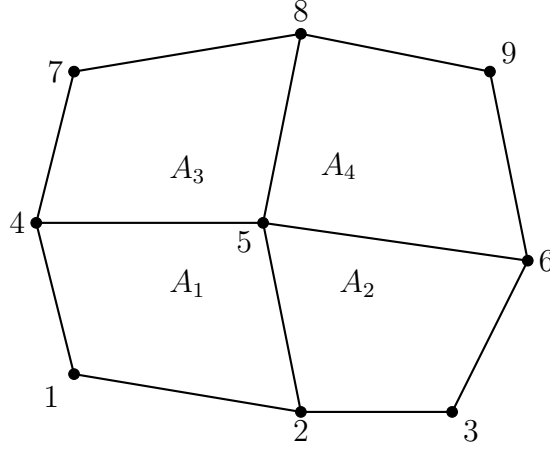
On the following figure (a,b), pressure fields for the lid driven cavity experiment are presented for both an even and un-even number of elements. We see that the amplitude of the modes can sometimes be so large that the 'real' pressure signal is not visible under the checkerboard and that something as simple as the number of elements in the domain can trigger those or not at all.



a) element pressure for a 32x32 grid and for a 33x33 grid;
b) image from [341, p307] for a manufactured solution; c) elemental pressure and smoothed pressure for the punch experiment [1261]

9.7.1 Scheme 1

The easiest post-processing step that can be used (especially when a regular grid is used) is explained in Thieulot *et al.* (2008) [1261]: "The element-to-node interpolation is performed by averaging the elemental values from elements common to each node; the node-to-element interpolation is performed by averaging the nodal values element-by-element. This method is not only very efficient but produces a smoothing of the pressure that is adapted to the local density of the octree. Note that these two steps can be repeated until a satisfying level of smoothness (and diffusion) of the pressure field is attained."



$$q_5^{(1)} = \frac{1}{4} \sum_{e=1}^4 p_e$$

In the codes which rely on the $Q_1 \times P_0$ element, the (elemental) pressure is simply defined as

```
p=np.zeros(nel , dtype=np.float64 )
```

while the nodal pressure is then defined as¹²

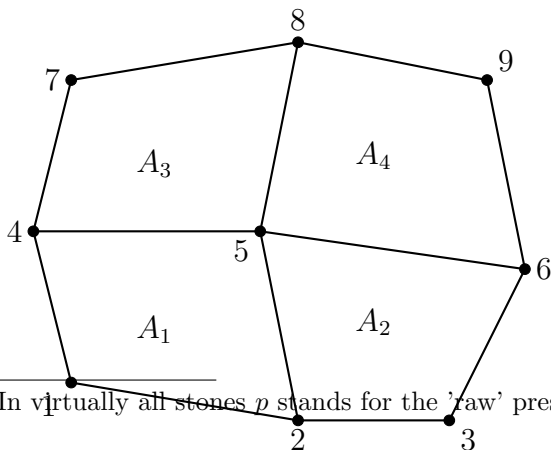
```
q=np.zeros(nnp , dtype=np.float64 )
```

The element-to-node algorithm is then simply (in 2D):

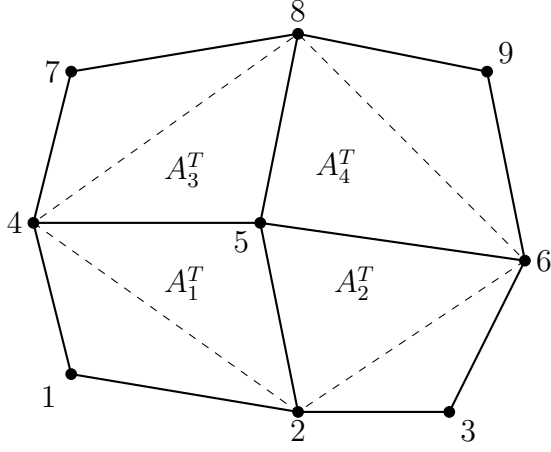
```
count=np.zeros(nnp, dtype=np.int32)
for iel in range(0,nel):
    q[icon[0,iel]]+=p[iel]
    q[icon[1,iel]]+=p[iel]
    q[icon[2,iel]]+=p[iel]
    q[icon[3,iel]]+=p[iel]
    count[icon[0,iel]]+=1
    count[icon[1,iel]]+=1
    count[icon[2,iel]]+=1
    count[icon[3,iel]]+=1
q=q/count
```

9.7.2 Schemes 2,3

Schemes 2,3 are very similar and are presented in Sani *et al.* (1981) [1108, 1109]. Scheme 2 uses the areas of the surrounding elements as weights for the arithmetic averaging while scheme 3 uses the area of the triangles:



¹²In virtually all stones p stands for the 'raw' pressure and q stands for its projection onto the velocity mesh.



$$q_5^{(2)} = \frac{\sum_{e=1}^4 A_e p_e}{\sum_{e=1}^4 A_e} \quad q_5^{(3)} = \frac{\sum_{e=1}^4 A_e^T p_e}{\sum_{e=1}^4 A_e^T}$$

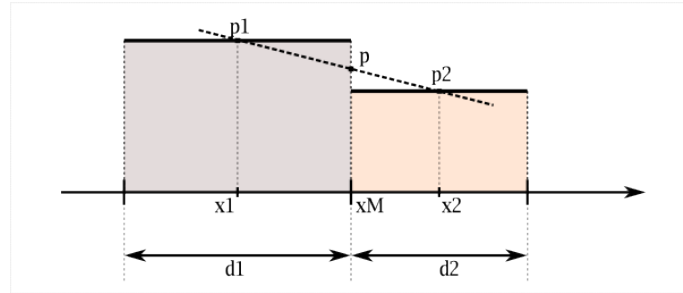
Remark. Although Schemes 1,2,3 are similar, scheme 1 is the simplest and fastest to implement since the areas of neighbouring elements/triangles are not needed.

Remark. Schemes 1,2,3 are identical if all elements are rectangles of identical dimensions.

9.7.3 Scheme 4

This scheme has been designed by me. It resembles the last three ones, but the weighing is in this case different.

Let us consider a 1D problem:



Elemental pressures p_1 and p_2 corresponding to elements 1 and 2 respectively are known at locations x_1 and x_2 . The two elements have a different size, characterised in this case by the distances d_1 and d_2 to their common edge.

The equation of the line passing through points (x_1, p_1) and (x_2, p_2) is

$$p(x) = \frac{p_2 - p_1}{x_2 - x_1}(x - x_1) + p_1$$

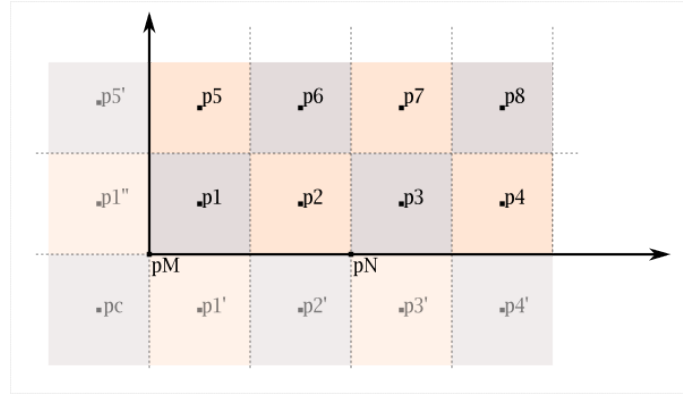
The x coordinate of the common edge is given by $x = x_1 + d_1/2$, and since $x_2 - x_1 = (d_1 + d_2)/2$, the pressure at this location writes:

$$p(x_M) = \frac{p_2 - p_1}{d_1 + d_2}d_1 + p_1 = \frac{\frac{p_1}{d_1} + \frac{p_2}{d_2}}{\frac{1}{d_1} + \frac{1}{d_2}}$$

Extrapolating this formula to 2D, d_1 and d_2 are in fact the element volumes, so that

$$q_5^{(4)} = \frac{\sum_{j=1}^4 \frac{p_j^e}{A_j^e}}{\sum_{j=1}^4 \frac{1}{A_j^e}} = \frac{\frac{p_1^e}{A_1^e} + \frac{p_2^e}{A_2^e} + \frac{p_3^e}{A_3^e} + \frac{p_4^e}{A_4^e}}{\frac{1}{A_1^e} + \frac{1}{A_2^e} + \frac{1}{A_3^e} + \frac{1}{A_4^e}}$$

There remains a problem, due to the presence of the boundary nodes for which the sums present in the above equation do not run up to 4. A boundary node only has three neighbours and a corner node only two. Additional measures are required for these nodes.



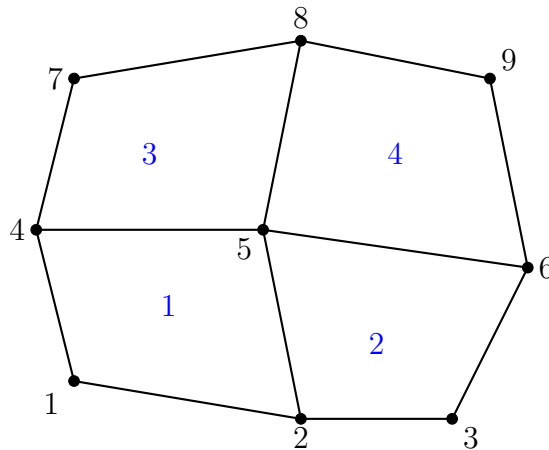
The pressure value p_N is obtained as follows:

$$q_N = \frac{\frac{p_2^e}{A_2^e} + \frac{p_3^e}{A_3^e} + \frac{p_{2'}^e}{A_{2'}^e} + \frac{p_{3'}^e}{A_{3'}^e}}{\frac{1}{A_2^e} + \frac{1}{A_3^e} + \frac{1}{A_{2'}^e} + \frac{1}{A_{3'}^e}}$$

The areas and pressures of the mirrored elements 2' and 3' are extrapolated from the areas of elements 2 and 6, and 3 and 7 respectively. Likewise the pressure p_M at the corner node is obtained through the pressures of its surrounding elements.

9.7.4 Scheme 5 - Least squares

This scheme is presented (among other places) in Lee *et al.* (1979) [759]. Let us start from the patch of 4 Q_1 elements counting 9 nodes:



We are looking for a field q living on the nodes. We build the quantity

$$J = \iint_{\Omega} (q - p)^2 dV$$

where p is the elemental field. To make things clearer we split the integral into the sum of elemental integrals:

$$J = \iint_{\Omega_1} (q(x, y) - p_1)^2 dV + \iint_{\Omega_2} (q(x, y) - p_2)^2 dV + \iint_{\Omega_3} (q(x, y) - p_3)^2 dV + \iint_{\Omega_4} (q(x, y) - p_4)^2 dV$$

Inside each element the field $q(x, y)$ is given by a bilinear interpolation so that:

$$\begin{aligned} J &= \iint_{\Omega_1} (\mathcal{N}_1(x, y)q_1 + \mathcal{N}_2(x, y)q_2 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_4(x, y)q_4 - p_1)^2 dV \\ &+ \iint_{\Omega_2} (\mathcal{N}_2(x, y)q_2 + \mathcal{N}_3(x, y)q_3 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_5(x, y)q_5 - p_2)^2 dV \\ &+ \iint_{\Omega_3} (\mathcal{N}_4(x, y)q_4 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_8(x, y)q_8 + \mathcal{N}_7(x, y)q_7 - p_3)^2 dV \\ &+ \iint_{\Omega_4} (\mathcal{N}_5(x, y)q_5 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_9(x, y)q_9 + \mathcal{N}_8(x, y)q_8 - p_4)^2 dV \end{aligned} \quad (9.24)$$

where the N_i functions are the basis functions (unusually expressed in x, y coordinates). The least square procedure looks for the set of q_i such that

$$\frac{\partial J}{\partial q_i} = 0 \quad \forall i = 1, \dots, 9$$

and this yields 9 equations/constraints for 9 unknowns.

$$\begin{aligned}
\frac{\partial J}{\partial q_1} &= \iint_{\Omega_1} 2(\mathcal{N}_1(x, y)q_1 + \mathcal{N}_2(x, y)q_2 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_4(x, y)q_4 - p_1)\mathcal{N}_1(x, y)dV \\
\frac{\partial J}{\partial q_2} &= \iint_{\Omega_1} 2(\mathcal{N}_1(x, y)q_1 + \mathcal{N}_2(x, y)q_2 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_4(x, y)q_4 - p_1)\mathcal{N}_2(x, y)dV \\
&+ \iint_{\Omega_2} 2(\mathcal{N}_2(x, y)q_2 + \mathcal{N}_3(x, y)q_3 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_5(x, y)q_5 - p_2)\mathcal{N}_2(x, y)dV \\
\frac{\partial J}{\partial q_3} &= \iint_{\Omega_2} 2(\mathcal{N}_2(x, y)q_2 + \mathcal{N}_3(x, y)q_3 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_5(x, y)q_5 - p_2)\mathcal{N}_3(x, y)dV \\
\frac{\partial J}{\partial q_4} &= \iint_{\Omega_1} 2(\mathcal{N}_1(x, y)q_1 + \mathcal{N}_2(x, y)q_2 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_4(x, y)q_4 - p_1)\mathcal{N}_4(x, y)dV \\
&+ \iint_{\Omega_3} 2(\mathcal{N}_4(x, y)q_4 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_8(x, y)q_8 + \mathcal{N}_7(x, y)q_7 - p_3)\mathcal{N}_4(x, y)dV \\
\frac{\partial J}{\partial q_5} &= \iint_{\Omega_1} 2(\mathcal{N}_1(x, y)q_1 + \mathcal{N}_2(x, y)q_2 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_4(x, y)q_4 - p_1)\mathcal{N}_5(x, y)dV \\
&+ \iint_{\Omega_2} 2(\mathcal{N}_2(x, y)q_2 + \mathcal{N}_3(x, y)q_3 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_5(x, y)q_5 - p_2)\mathcal{N}_5(x, y)dV \\
&+ \iint_{\Omega_3} 2(\mathcal{N}_4(x, y)q_4 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_8(x, y)q_8 + \mathcal{N}_7(x, y)q_7 - p_3)\mathcal{N}_5(x, y)dV \\
&+ \iint_{\Omega_4} 2(\mathcal{N}_5(x, y)q_5 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_9(x, y)q_9 + \mathcal{N}_8(x, y)q_8 - p_4)\mathcal{N}_5(x, y)dV \\
\frac{\partial J}{\partial q_6} &= \iint_{\Omega_2} 2(\mathcal{N}_2(x, y)q_2 + \mathcal{N}_3(x, y)q_3 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_5(x, y)q_5 - p_2)\mathcal{N}_6(x, y)dV \\
&+ \iint_{\Omega_4} 2(\mathcal{N}_5(x, y)q_5 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_9(x, y)q_9 + \mathcal{N}_8(x, y)q_8 - p_4)\mathcal{N}_6(x, y)dV \\
\frac{\partial J}{\partial q_7} &= \iint_{\Omega_3} 2(\mathcal{N}_4(x, y)q_4 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_8(x, y)q_8 + \mathcal{N}_7(x, y)q_7 - p_3)\mathcal{N}_7(x, y)dV \\
\frac{\partial J}{\partial q_8} &= \iint_{\Omega_3} 2(\mathcal{N}_4(x, y)q_4 + \mathcal{N}_5(x, y)q_5 + \mathcal{N}_8(x, y)q_8 + \mathcal{N}_7(x, y)q_7 - p_3)\mathcal{N}_8(x, y)dV \\
&+ \iint_{\Omega_4} 2(\mathcal{N}_5(x, y)q_5 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_9(x, y)q_9 + \mathcal{N}_8(x, y)q_8 - p_4)\mathcal{N}_8(x, y)dV \\
\frac{\partial J}{\partial q_9} &= \iint_{\Omega_4} 2(\mathcal{N}_5(x, y)q_5 + \mathcal{N}_6(x, y)q_6 + \mathcal{N}_9(x, y)q_9 + \mathcal{N}_8(x, y)q_8 - p_4)\mathcal{N}_9(x, y)dV \quad (9.25)
\end{aligned}$$

The factor 2 are removed and the terms $\int p_i N_j$ are known so they end up in the right hand side.

$$\begin{aligned}
&\iint_{\Omega_1} (\mathcal{N}_1 \mathcal{N}_1 q_1 + \mathcal{N}_1 \mathcal{N}_2 q_2 + \mathcal{N}_1 \mathcal{N}_5 q_5 + \mathcal{N}_1 \mathcal{N}_4 q_4) dV = \iint_{\Omega_1} p_1 \mathcal{N}_1 dV \\
&\iint_{\Omega_1} (\mathcal{N}_2 \mathcal{N}_1 q_1 + \mathcal{N}_2 \mathcal{N}_2 q_2 + \mathcal{N}_2 \mathcal{N}_5 q_5 + \mathcal{N}_2 \mathcal{N}_4 q_4) dV \\
&+ \iint_{\Omega_2} (\mathcal{N}_2 \mathcal{N}_2 q_2 + \mathcal{N}_3 \mathcal{N}_2 q_3 + \mathcal{N}_6 \mathcal{N}_2 q_6 + \mathcal{N}_5 \mathcal{N}_2 q_5) dV = \iint_{\Omega_1} p_1 \mathcal{N}_2 dV + \iint_{\Omega_2} p_2 \mathcal{N}_2 dV \\
&\dots = \dots \\
&\iint_{\Omega_4} (\mathcal{N}_9 \mathcal{N}_5 q_5 + \mathcal{N}_9 \mathcal{N}_6 q_6 + \mathcal{N}_9 \mathcal{N}_9 q_9 + \mathcal{N}_9 \mathcal{N}_8 q_8) dV = \iint_{\Omega_4} p_4 \mathcal{N}_9 dV \quad (9.26)
\end{aligned}$$

The mass matrices corresponding to the four elements are

$$\begin{aligned} \mathbf{M}_1 &= \int_{\Omega_1} \begin{pmatrix} \mathcal{N}_1\mathcal{N}_1 & \mathcal{N}_1\mathcal{N}_2 & \mathcal{N}_1\mathcal{N}_5 & \mathcal{N}_1\mathcal{N}_4 \\ \mathcal{N}_2\mathcal{N}_1 & \mathcal{N}_2\mathcal{N}_2 & \mathcal{N}_2\mathcal{N}_5 & \mathcal{N}_2\mathcal{N}_4 \\ \mathcal{N}_5\mathcal{N}_1 & \mathcal{N}_5\mathcal{N}_2 & \mathcal{N}_5\mathcal{N}_5 & \mathcal{N}_5\mathcal{N}_4 \\ \mathcal{N}_4\mathcal{N}_1 & \mathcal{N}_4\mathcal{N}_2 & \mathcal{N}_4\mathcal{N}_5 & \mathcal{N}_4\mathcal{N}_4 \end{pmatrix} dV & \mathbf{M}_2 &= \int_{\Omega_2} \begin{pmatrix} \mathcal{N}_2\mathcal{N}_2 & \mathcal{N}_2\mathcal{N}_3 & \mathcal{N}_2\mathcal{N}_6 & \mathcal{N}_2\mathcal{N}_5 \\ \mathcal{N}_3\mathcal{N}_2 & \mathcal{N}_3\mathcal{N}_3 & \mathcal{N}_3\mathcal{N}_6 & \mathcal{N}_3\mathcal{N}_5 \\ \mathcal{N}_6\mathcal{N}_2 & \mathcal{N}_6\mathcal{N}_3 & \mathcal{N}_6\mathcal{N}_6 & \mathcal{N}_6\mathcal{N}_5 \\ \mathcal{N}_5\mathcal{N}_2 & \mathcal{N}_5\mathcal{N}_3 & \mathcal{N}_5\mathcal{N}_6 & \mathcal{N}_5\mathcal{N}_5 \end{pmatrix} dV \\ \mathbf{M}_3 &= \int_{\Omega_3} \begin{pmatrix} \mathcal{N}_4\mathcal{N}_4 & \mathcal{N}_4\mathcal{N}_5 & \mathcal{N}_4\mathcal{N}_8 & \mathcal{N}_4\mathcal{N}_7 \\ \mathcal{N}_5\mathcal{N}_4 & \mathcal{N}_5\mathcal{N}_5 & \mathcal{N}_5\mathcal{N}_8 & \mathcal{N}_5\mathcal{N}_7 \\ \mathcal{N}_8\mathcal{N}_4 & \mathcal{N}_8\mathcal{N}_5 & \mathcal{N}_8\mathcal{N}_8 & \mathcal{N}_8\mathcal{N}_7 \\ \mathcal{N}_7\mathcal{N}_4 & \mathcal{N}_7\mathcal{N}_5 & \mathcal{N}_7\mathcal{N}_8 & \mathcal{N}_7\mathcal{N}_7 \end{pmatrix} dV & \mathbf{M}_4 &= \int_{\Omega_4} \begin{pmatrix} \mathcal{N}_5\mathcal{N}_5 & \mathcal{N}_5\mathcal{N}_6 & \mathcal{N}_5\mathcal{N}_9 & \mathcal{N}_5\mathcal{N}_8 \\ \mathcal{N}_6\mathcal{N}_5 & \mathcal{N}_6\mathcal{N}_6 & \mathcal{N}_6\mathcal{N}_9 & \mathcal{N}_6\mathcal{N}_8 \\ \mathcal{N}_9\mathcal{N}_5 & \mathcal{N}_9\mathcal{N}_6 & \mathcal{N}_9\mathcal{N}_9 & \mathcal{N}_9\mathcal{N}_8 \\ \mathcal{N}_8\mathcal{N}_5 & \mathcal{N}_8\mathcal{N}_6 & \mathcal{N}_8\mathcal{N}_9 & \mathcal{N}_8\mathcal{N}_8 \end{pmatrix} dV \end{aligned}$$

so that the 9 equations above are actually the result of the assembly process of these four elemental systems:

$$\left(\iint_{\Omega_e} \vec{\mathcal{N}}^T \vec{\mathcal{N}} dV \right) \cdot \vec{q}_e = \iint_{\Omega_e} \vec{\mathcal{N}}^T p_e dV \quad e = 1, 2, 3, 4$$

Also check section 4.5.4 of Glaisner and Tezduyar [464] (1987), in which the authors present a two-step algorithm: 1) pressure is averaged over each element. 2) the nodal values of the pressure are recovered through a least-squares approach.

9.7.5 Scheme 6 - Consistent pressure recovery

This is the method presented in Zienkiewicz and Nakazawa [1421] (1982). In the second part of this publication the authors wish to establish a simple and effective numerical method to calculate variables eliminated by the penalisation process. The method involves an additional finite element solution for the nodal pressures using the same finite element basis and numerical quadrature as used for the velocity.

Let us start with¹³:

$$q = -\lambda \vec{\nabla} \cdot \vec{v}$$

We are going to treat this equation like any other PDE in the context of the FE method, i.e. we are going to establish its weak form. We assume that the pressure is given inside an element by

$$q(x, y) = \sum_{i=1}^{m_v} \mathcal{N}_i(x, y) q_i = \vec{\mathcal{N}} \cdot \vec{q}$$

and the velocity:

$$\vec{v} = (u, v) \quad u(x, y) = \sum_{i=1}^{m_v} \mathcal{N}_i(x, y) u_i \quad v(x, y) = \sum_{i=1}^{m_v} \mathcal{N}_i(x, y) v_i$$

where the \mathcal{N}_i are the Q_1 basis functions and q_i are the sought after nodal values. We multiply the equation above by a Q_1 basis function \mathcal{N}_i and integrate over the whole domain:

$$\iint_{\Omega} \mathcal{N}_i(x, y) q(x, y) dx dy = -\lambda \iint_{\Omega} \mathcal{N}_i \vec{\nabla} \cdot \vec{v} dx dy$$

As before we now focus on the above expression inside a single element e :

$$\iint_{\Omega_e} \mathcal{N}_i(x, y) q(x, y) dx dy = -\lambda \iint_{\Omega_e} \mathcal{N}_i \vec{\nabla} \cdot \vec{v} dx dy$$

¹³I here voluntarily use q instead of p

After $\mathcal{N}_i \rightarrow \vec{\mathcal{N}} = (\mathcal{N}_1, \mathcal{N}_2, \mathcal{N}_3, \mathcal{N}_4)^T$, the left hand side term becomes:

$$\iint_{\Omega_e} \vec{\mathcal{N}}^T q(x, y) dx dy = \iint_{\Omega_e} \vec{\mathcal{N}}^T \vec{\mathcal{N}} \cdot \vec{q} dx dy = \left(\underbrace{\iint_{\Omega_e} \vec{\mathcal{N}}^T \vec{\mathcal{N}} dx dy}_{\mathbf{M}_e} \right) \cdot \vec{q}$$

where \mathbf{M}_e is the elemental mass matrix. We now turn to the right hand side. We have

$$\vec{\nabla} \cdot \vec{v} = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = \sum_i \frac{\partial \mathcal{N}_i}{\partial x} u_i + \sum_i \frac{\partial \mathcal{N}_i}{\partial y} v_i$$

We here too define $\vec{V}_e = (u_1, v_1, u_2, v_2, u_3, v_3, u_4, v_4)^T$ so that

$$\begin{aligned} & \iint_{\Omega_e} \vec{\mathcal{N}} \vec{\nabla} \cdot \vec{v} dV \\ &= \iint_{\Omega_e} \vec{\mathcal{N}}^T \sum_{i=1}^{m_v} \left(\frac{\partial \mathcal{N}_i}{\partial x} u_i + \frac{\partial \mathcal{N}_i}{\partial y} v_i \right) dV \\ &= \iint_{\Omega_e} \begin{pmatrix} \mathcal{N}_1 \left(\sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} u_i + \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} v_i \right) \\ \mathcal{N}_2 \left(\sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} u_i + \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} v_i \right) \\ \mathcal{N}_3 \left(\sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} u_i + \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} v_i \right) \\ \mathcal{N}_4 \left(\sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial x} u_i + \sum_{i=1}^4 \frac{\partial \mathcal{N}_i}{\partial y} v_i \right) \end{pmatrix} dV \\ &= \int_{\Omega_e} \begin{pmatrix} \mathcal{N}_1 & \mathcal{N}_1 & 0 \\ \mathcal{N}_2 & \mathcal{N}_2 & 0 \\ \mathcal{N}_3 & \mathcal{N}_3 & 0 \\ \mathcal{N}_4 & \mathcal{N}_4 & 0 \end{pmatrix} \cdot \begin{pmatrix} \sum_i \frac{\partial \mathcal{N}_i}{\partial x} u_i \\ \sum_i \frac{\partial \mathcal{N}_i}{\partial y} v_i \\ \sum_i \left(\frac{\partial \mathcal{N}_i}{\partial y} u_i + \frac{\partial \mathcal{N}_i}{\partial x} v_i \right) \end{pmatrix} dV \\ &= \int_{\Omega_e} \underbrace{\begin{pmatrix} \mathcal{N}_1 & \mathcal{N}_1 & 0 \\ \mathcal{N}_2 & \mathcal{N}_2 & 0 \\ \mathcal{N}_3 & \mathcal{N}_3 & 0 \\ \mathcal{N}_4 & \mathcal{N}_4 & 0 \end{pmatrix}}_{\mathbf{N}} \cdot \underbrace{\begin{pmatrix} \partial_x \mathcal{N}_1 & 0 & \partial_x \mathcal{N}_2 & 0 & \partial_x \mathcal{N}_3 & 0 & \partial_x \mathcal{N}_4 & 0 \\ 0 & \partial_y \mathcal{N}_1 & 0 & \partial_y \mathcal{N}_2 & 0 & \partial_y \mathcal{N}_3 & 0 & \partial_y \mathcal{N}_4 \\ \partial_y \mathcal{N}_1 & \partial_x \mathcal{N}_1 & \partial_y \mathcal{N}_2 & \partial_x \mathcal{N}_2 & \partial_y \mathcal{N}_3 & \partial_x \mathcal{N}_3 & \partial_y \mathcal{N}_4 & \partial_x \mathcal{N}_4 \end{pmatrix}}_{\mathbf{B}} \cdot \vec{V}_e dV \\ &= \left(\int_{\Omega_e} \mathbf{N} \cdot \mathbf{B} dV \right) \cdot \vec{V}_e \\ &= -\mathbb{G}_e^T \cdot \vec{V}_e \end{aligned} \tag{9.27}$$

This makes sense since \mathbb{G}^T is the discrete divergence operator. However, it is not very efficient to build \mathbb{G}_e only to multiply it with a vector of already known quantities. In practice we implement Eq. (9.27) which implementation resembles the buoyancy term of the Stokes equation.

After assembly we arrive at

$$\mathbf{M} \cdot \vec{q} = \lambda \mathbb{G}^T \cdot \vec{V} \quad \text{with} \quad \mathbb{G}_e = - \int_{\Omega_e} \mathbf{N} \cdot \mathbf{B} dV$$

where \mathbf{M} is the global mass matrix, \vec{q} the vector of all nodal pressures, \mathbb{G} the discrete gradient matrix and \vec{V} the (velocity) solution vector. The system can be easily solved since the mass matrix is a friendly matrix. The vector \vec{q} contains the nodal pressure values directly, with no need for a smoothing scheme!

Remark. *Very importantly, the mass matrix \mathbf{M} is to be evaluated at the full integration points, while the constraint part (the right hand side of the equation) is to be evaluated at the reduced integration point, i.e. in the middle of the element.*

Remark. *As noted in [1421], it is interesting to note that when linear elements are used and the lumped matrices are used for the \mathbf{M} the resulting algebraic equation is identical to the smoothing scheme 1 only if a uniform square finite element mesh is used. In this respect this method is expected to yield different results when elements are not square or even rectangular.*

Remark. *The third column of the matrix \mathbf{N} and the last line of the \mathbf{B} matrix could be removed altogether. If your code is based on the mixed formulation, then you already have built matrix \mathbb{G} so you can easily re-use this piece of code to compute \mathbb{G} again, this time with a reduced integration quadrature. If you are using the penalty formulation then you need to program all from scratch and then simply do away with these unnecessary terms, or you can directly build the rhs as $\int_{\Omega_e} \vec{N}^T p_e$ (assuming you have previously computed the pressure in the middle of each element by means of $p = -\lambda \vec{\nabla} \cdot \vec{v}$).*

Remark. *This scheme is identical to the least square scheme!*

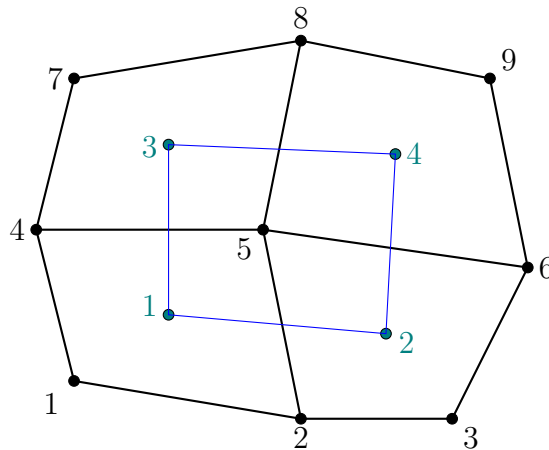
9.7.6 Scheme 7

Same as scheme 6, but with lumped mass matrix.

9.7.7 Scheme 8 - bilinear interpolation

Let us assume that the centers of the four elements make a Q_1 quadrilateral element, as shown on this figure:

(tikz_pscheme8.tex)



The values at the corners are p_1, p_2, p_3 and p_4 . Assuming that the pressure inside this element can be represented by a bilinear field, we have

$$p(x, y) = a + bx + cy + dxy$$

where the coefficients will be determined by ensuring that $p(x_i, y_i) = p_i$ for $i = 1, 2, 3, 4$, or:

$$a + bx_1 + cy_1 + dx_1y_1 = p_1 \quad (9.29)$$

$$a + bx_2 + cy_2 + dx_2y_2 = p_2 \quad (9.30)$$

$$a + bx_3 + cy_3 + dx_3y_3 = p_3 \quad (9.31)$$

$$a + bx_4 + cy_4 + dx_4y_4 = p_4 \quad (9.32)$$

i.e.

$$\begin{pmatrix} 1 & x_1 & y_1 & x_1y_1 \\ 1 & x_2 & y_2 & x_2y_2 \\ 1 & x_3 & y_3 & x_3y_3 \\ 1 & x_4 & y_4 & x_4y_4 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$$

There remains an issue with nodes which are on the boundaries of the domain. These are of course not 'surrounded' by four pressure values so the above algorithm does not apply directly. However, looking at the above figure, and assuming that node 1 is a lower left corner of a 2D domain, we can use the bilinear interpolation based on elements 1,2,3,4 to extrapolate a nodal pressure value at node 1. The same would apply for nodes 2 and 4 for example.

Remark. *This scheme is not applicable to quadtree-based meshed.*

9.8 The value of the timestep

The chosen time step δt used for time integration is chosen to comply with the Courant-Friedrichs-Lewy condition [22].

$$\delta t = C \min \left(\frac{\min h}{\max |\vec{v}|^p}, \frac{h^2}{\kappa} \right) \quad (9.33)$$

where h is a measure of the element diameter, p is the polynomial order of the element, $\kappa = k/\rho C_p$ is the thermal diffusivity and C is the so-called CFL number chosen in $[0, 1[$. $\min h$ is the smallest element diameter in the domain, while $\max |\vec{v}|$ is the maximum velocity (norm) in the domain.

In essence the CFL condition arises when solving hyperbolic PDEs . It limits the time step in many explicit time-marching computer simulations so that the simulation does not produce (too) incorrect results.

This condition is not needed when solving the Stokes equation but it is mandatory when solving the heat transport equation or any kind of advection-diffusion equation. Note that any increase of grid resolution (i.e. h becomes smaller) yields an automatic decrease of the time step value.

9.9 Exporting data to vtk/vtu format

This format seems to be the universally accepted format for 2D and 3D visualisation in Computational Geodynamics (and even CFD ?). Such files can be opened with open source softwares such as Paraview ¹⁴, MayaVi ¹⁵ or Visit ¹⁶.

Unfortunately it is my experience that no simple tutorial exists about how to build such files. There is an official document which describes the vtk format¹⁷ but it delivers the information in a convoluted way¹⁸. I therefore describe hereafter how fieldstone builds the vtk/vtu files¹⁹.

I hereunder show vtk file corresponding to a 3×2 grid made of linear elements. In this particular example there are:

- 12 nodes and 6 elements
- 1 elemental field (the pressure p)
- 2 nodal fields: 1 scalar (the smoothed pressure q), 1 vector (the velocity field u,v,0)

Note that vtk files are inherently 3D so that even in the case of a 2D simulation the z -coordinate of the points and for instance their z -velocity component must be provided. The file, usually called *solution.vtk* starts with a header:

```
<VTKFile type='UnstructuredGrid' version='0.1' byte_order='BigEndian'>
<UnstructuredGrid>
<Piece NumberOfPoints='12' NumberOfCells='6'>
```

We then proceed to write the node coordinates as follows:

```
<Points>
<DataArray type='Float32' NumberOfComponents='3' Format='ascii'>
0.000000e+00 0.000000e+00 0.000000e+00
3.333333e-01 0.000000e+00 0.000000e+00
6.666667e-01 0.000000e+00 0.000000e+00
1.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 5.000000e-01 0.000000e+00
3.333333e-01 5.000000e-01 0.000000e+00
6.666667e-01 5.000000e-01 0.000000e+00
1.000000e+00 5.000000e-01 0.000000e+00
0.000000e+00 1.000000e+00 0.000000e+00
3.333333e-01 1.000000e+00 0.000000e+00
6.666667e-01 1.000000e+00 0.000000e+00
1.000000e+00 1.000000e+00 0.000000e+00
</DataArray>
</Points>
```

These are followed by the elemental field(s):

```
<CellData Scalars='scalars'>
<DataArray type='Float32' Name='p' Format='ascii'>
-1.333333e+00
-3.104414e-10
1.333333e+00
-1.333333e+00
```

¹⁴<https://www.paraview.org/>

¹⁵<https://docs.enthought.com/mayavi/mayavi/>

¹⁶<https://wci.llnl.gov/simulation/computer-codes/visit/>

¹⁷<https://www.vtk.org/wp-content/uploads/2015/04/file-formats.pdf>

¹⁸I only later realised it was also limited!

¹⁹I have found the following information about vtk vs. vtu: VTK denotes the simple legacy format file. VTU denotes a Serial Unstructured Grid format information for the XML-based syntax.

```

8.278417e-17
1.333333e+00
</DataArray>
</CellData>

```

Nodal quantities are written next:

```

<PointData Scalars='scalars'>
<DataArray type='Float32' NumberOfComponents='3' Name='velocity' Format='ascii'>
0.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
8.888885e-08 -8.278405e-24 0.000000e+00
8.888885e-08 1.655682e-23 0.000000e+00
0.000000e+00 0.000000e+00 0.000000e+00
1.000000e+00 0.000000e+00 0.000000e+00
1.000000e+00 0.000000e+00 0.000000e+00
1.000000e+00 0.000000e+00 0.000000e+00
1.000000e+00 0.000000e+00 0.000000e+00
</DataArray>
<DataArray type='Float32' NumberOfComponents='1' Name='q' Format='ascii'>
-1.333333e+00
-6.666664e-01
6.666664e-01
1.333333e+00
-1.333333e+00
-6.666664e-01
6.666664e-01
1.333333e+00
-1.333333e+00
-6.666664e-01
6.666664e-01
1.333333e+00
</DataArray>
</PointData>

```

To these informations we must append 3 more datasets. The first one is the connectivity, the second one is the offsets and the third one is the type. The first one is trivial since the required connectivity array is the same as the one needed for the Finite Elements. The second must be understood as follows: when reading the connectivity information in a linear manner the offset values indicate the beginning of each element (omitting the zero value). The third is simply the type of element as given in the vtk format document (9 corresponds to a generic quadrilateral with an internal numbering consistent with ours - more on this later).

```

<Cells>
<DataArray type='Int32' Name='connectivity' Format='ascii'>
0 1 5 4
1 2 6 5
2 3 7 6
4 5 9 8
5 6 10 9
6 7 11 10
</DataArray>
<DataArray type='Int32' Name='offsets' Format='ascii'>
4
8
12
16

```

```

20
24
</DataArray>
<DataArray type='Int32' Name='types' Format='ascii'>
9
9
9
9
9
9
</DataArray>
</Cells>

```

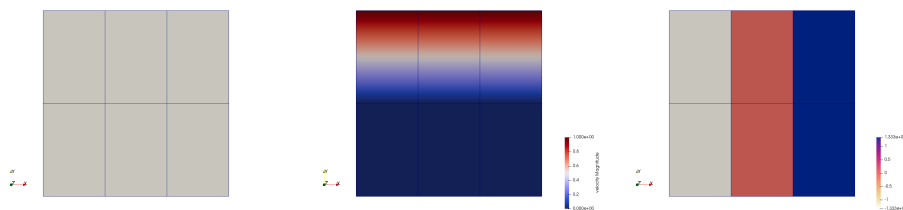
The file is then closed with

```

</Piece>
</UnstructuredGrid>
</VTKFile>

```

The *solution.vtu*²⁰ can then be opened with ParaView, MayaVi or Visit and the reader is advised to find tutorials online on how to install and use these softwares. Also check Appendix O.0.1 on how to use ParaView.



In the same folder `images/vtk` there is the python script `makevtu.py`²¹ which produces 3 different vtu files. The first one *solution1.vtu* is similar to the one above: an `nelx*nely` quadrilateral-based mesh in a unit square. The second one (*solution2.vtu*) looks identical when opened in Paraview but it is rather different: each element is exported as its own sub-mesh, so that if the mesh counts `nel` elements the number of vertices is `4*nel`, and not `(nel+1)*(nely+1)`. As such this file is larger. The icon array is needed to write down the positions of the four vertices of each element but not to write down the connectivity since the first 4 points are making the 1st element, the next four points are making the second element, etc ...

```

vtufile.write("<Points>\n")
vtufile.write("<DataArray type='Float32' NumberOfComponents='3' Format='ascii'>\n")
for iel in range(0,nel):
    if not flag[iel]:
        for k in range(0,m):
            vtufile.write("%10e_%10e_%10e\n" % (x[icon[k,iel]],y[icon[k,iel]],0.))
vtufile.write("</DataArray>\n")
vtufile.write("</Points>\n")
vtufile.write("<Cells>\n")
vtufile.write("<DataArray type='Int32' Name='connectivity' Format='ascii'>\n")
for iel in range(0,nel_left):
    vtufile.write("%d_%d_%d_%d\n" % (iel*4,iel*4+1,iel*4+2,iel*4+3))
vtufile.write("</DataArray>\n")
...
vtufile.write("</Cells>\n")

```

²⁰<https://raw.githubusercontent.com/cedrict/fieldstone/master/images/vtk/solution.vtu>

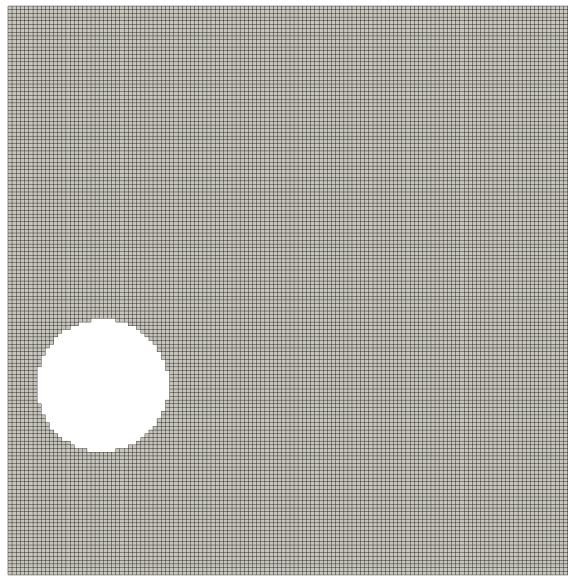
²¹<https://raw.githubusercontent.com/cedrict/fieldstone/master/images/vtk/makevtu.py>

This format is rather practical in the case of linear or higher order discontinuous fields. For example, in the case of the $Q_2 \times P_{-1}$ element pair, the pressure is linear inside each element and discontinuous across element edges. One can then assign pressure values at the four vertices of each element.

Finally a third mesh *solution3.vtu* is produced. It is based on the 2nd one, but since elements are now somewhat de-coupled, then one can export only a subset of the mesh. For instance one could not show elements which are too distorted, or below a certain line, or outside a certain volume, etc ... In [makevtu.py](#) all elements whose center is inside a circle are flagged and will not be exported into the vtu file:

```
for iel in range(0,nel):
    flag[iel]= (xc[iel]-0.333*Lx)**2+(yc[iel]-0.666*Ly)**2<0.234**2
nel_flagged=np.sum(flag)
nel_left=nel-nel_flagged
```

Once opened in ParaView this is how it looks like:



We can find at this address²² a list of supported cell types (although their internal numbering is not mentioned so enjoy the guess work). Here is a subset of these that can be relevant in our case:

```
// Linear cells
VTK_EMPTY_CELL = 0,
VTK_VERTEX = 1,
VTK_POLY_VERTEX = 2,
VTK_LINE = 3,
VTK_POLY_LINE = 4,
VTK_TRIANGLE = 5,
VTK_TRIANGLE_STRIP = 6,
VTK_POLYGON = 7,
VTK_PIXEL = 8,
VTK_QUAD = 9,
VTK_TETRA = 10,
VTK_VOXEL = 11,
VTK_HEXAHEDRON = 12,
VTK_WEDGE = 13,
```

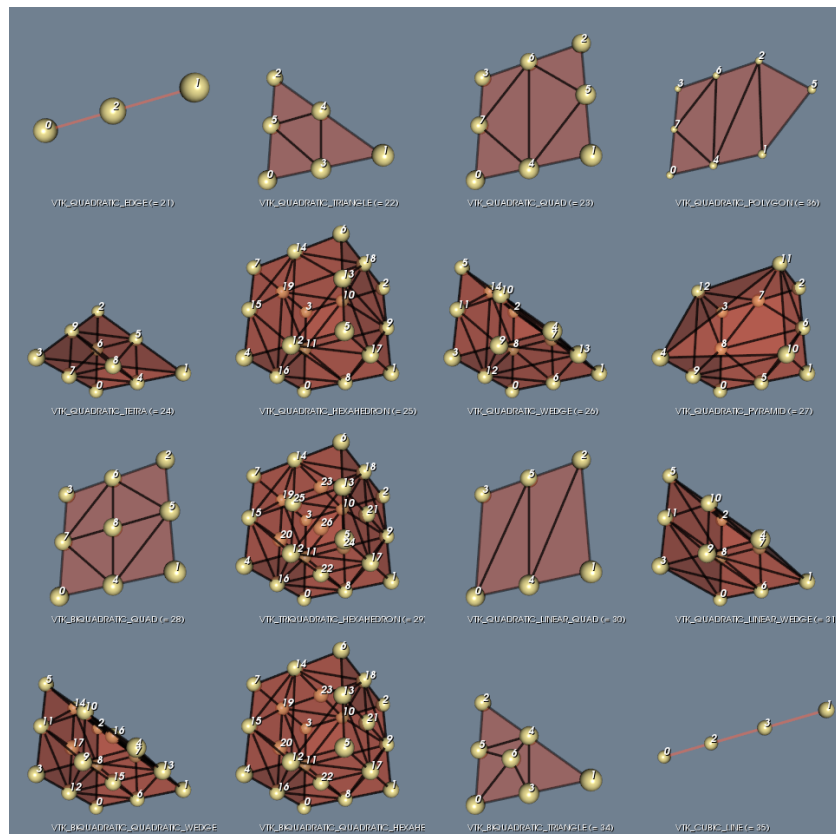
²²https://vtk.org/doc/nightly/html/vtkCellType_8h_source.html

```

VTK_PYRAMID = 14,
VTK_PENTAGONAL_PRISM = 15,
VTK_HEXAGONAL_PRISM = 16,
// Quadratic, isoparametric cells
VTK_QUADRATIC_EDGE = 21,
VTK_QUADRATIC_TRIANGLE = 22,
VTK_QUADRATIC_QUAD = 23,
VTK_QUADRATIC_POLYGON = 36,
VTK_QUADRATIC_TETRA = 24,
VTK_QUADRATIC_HEXAHEDRON = 25,
VTK_QUADRATIC_WEDGE = 26,
VTK_QUADRATIC_PYRAMID = 27,
VTK_BIQUADRATIC_QUAD = 28,
VTK_TRIQUADRATIC_HEXAHEDRON = 29,
VTK_TRIQUADRATIC_PYRAMID = 37,
VTK_QUADRATIC_LINEAR_QUAD = 30,
VTK_QUADRATIC_LINEAR_WEDGE = 31,
VTK_BIQUADRATIC_QUADRATIC_WEDGE = 32,
VTK_BIQUADRATIC_QUADRATIC_HEXAHEDRON = 33,
VTK_BIQUADRATIC_TRIANGLE = 34,
// Cubic, isoparametric cell
VTK_CUBIC_LINE = 35,

```

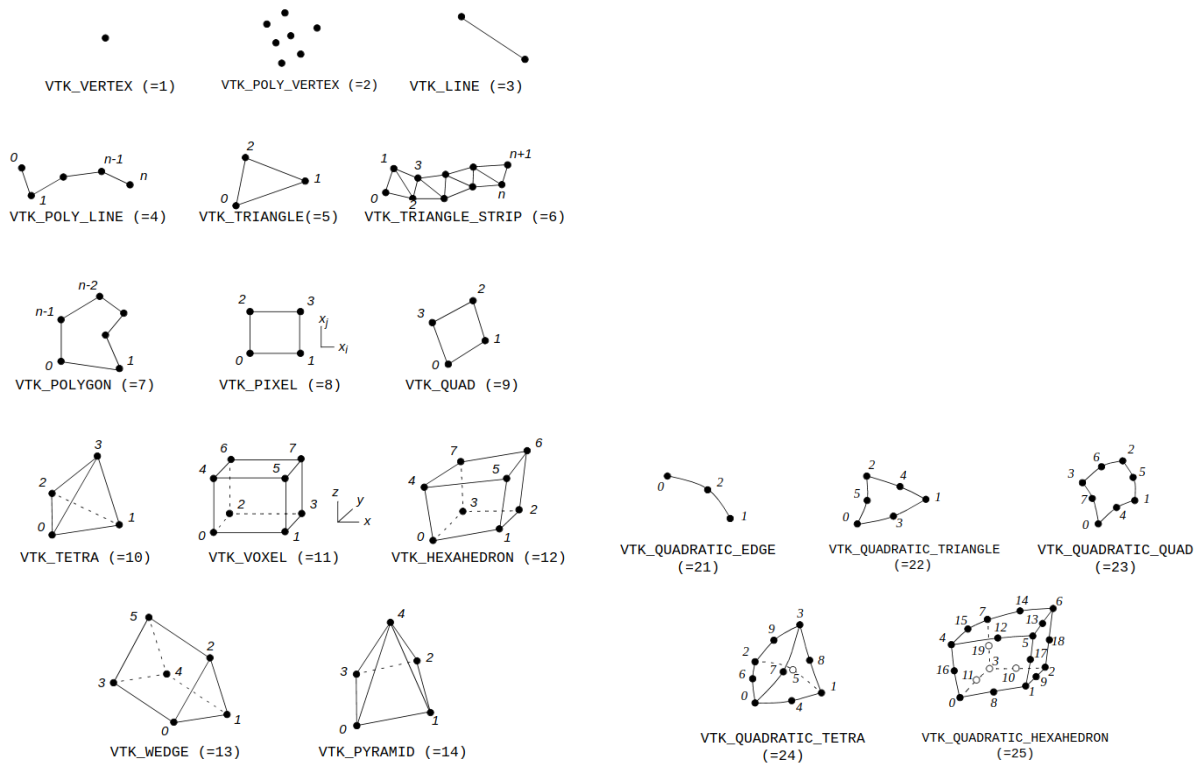
After some online searching I found this figure²³:



Fortunately these cell types seem to coincide with those in this document²⁴ mentioned earlier where we find these schematics:

²³<https://examples.vtk.org/site/Cxx/GeometricObjects/IsoparametricCellsDemo/>

²⁴<https://www.vtk.org/wp-content/uploads/2015/04/file-formats.pdf>



9.10 Runge-Kutta methods

These methods were developed around 1900 by the German mathematicians Carl Runge and Martin Kutta. The RK methods are methods for the numerical integration of ODEs²⁵. These methods are well documented in any numerical analysis textbook and the reader is referred to [455, 626]. Any Runge-Kutta method is uniquely identified by its Butcher tableau (REF?) which contains all necessary coefficients to build the algorithm.

missing
refs for
Butcher
tableau

The simplest Runge-Kutta method is the (forward) Euler method. Its tableau is:

| | |
|---|---|
| 0 | |
| 1 | 1 |

The standard second-order RK method (also called midpoint method) is:

| | | |
|-----|-----|---|
| 0 | | |
| 1/2 | 1/2 | |
| | 0 | 1 |

Another second-order RK method, called Heun's method²⁶ is follows:

| | | |
|---|-----|-----|
| 0 | | |
| 1 | 1 | |
| | 1/2 | 1/2 |

A third-order RK method is as follows:

| | | | |
|-----|-----|-----|-----|
| 0 | | | |
| 1/2 | 1/2 | | |
| 1 | -1 | 2 | |
| | 1/6 | 4/6 | 1/6 |

The RK4 method falls in this framework. Its tableau is:

| | | | | |
|-----|-----|-----|-----|-----|
| 0 | | | | |
| 1/2 | 1/2 | | | |
| 1/2 | 0 | 1/2 | | |
| 1 | 0 | 0 | 1 | |
| | 1/6 | 1/3 | 1/3 | 1/6 |

A slight variation of the standard RK4 method is also due to Kutta in 1901 and is called the 3/8-rule. Almost all of the error coefficients are smaller than in the standard method but it requires slightly more FLOPs per time step. Its Butcher tableau is

| | | | | |
|-----|------|-----|-----|-----|
| 0 | | | | |
| 1/3 | 1/3 | | | |
| 2/3 | -1/3 | 1 | | |
| 1 | 1 | -1 | 1 | |
| | 1/8 | 3/8 | 3/8 | 1/8 |

²⁵https://en.wikipedia.org/wiki/Runge-Kutta_methods

²⁶https://en.wikipedia.org/wiki/Heun's_method

- The RK2 method is also simple but requires a bit more work.

| | |
|---|---------|
| 0 | |
| 1 | 1 |
| | 1/2 1/2 |

Carry out a loop over markers and

1. interpolate velocity \vec{v}_m onto each marker m at position \vec{r}_m
2. compute new intermediate position as follows: $\vec{r}_m^{(1)}(t + \delta t) = \vec{r}_m(t) + \vec{v}_m \delta t / 2$
3. compute velocity $\vec{v}_m^{(1)}$ at position $\vec{r}_m^{(1)}$
4. compute new position: $\vec{r}_m(t + \delta t) = \vec{r}_m(t) + \vec{v}_m^{(1)} \delta t$

Note that the intermediate positions could be in a different element of the mesh so extra care must be taken when computing intermediate velocities.

- The RK3 method introduces two intermediate steps.

| | | | |
|-----|---------------|---------------|---------------|
| 0 | | | |
| 1/2 | $\frac{1}{2}$ | | |
| 1 | -1 | 2 | |
| | $\frac{1}{6}$ | $\frac{4}{6}$ | $\frac{1}{6}$ |

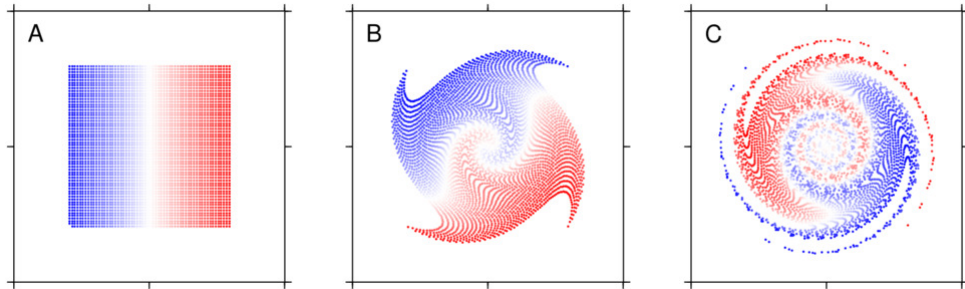
Carry out a loop over markers and

1. interpolate velocity \vec{v}_m onto each marker m at position \vec{r}_m
2. compute new intermediate position as follows: $\vec{r}_m^{(1)}(t + \delta t) = \vec{r}_m(t) + \frac{1}{2} \vec{v}_m \delta t$
3. compute velocity $\vec{v}_m^{(1)}$ at position $\vec{r}_m^{(1)}$
4. compute new intermediate position as follows: $\vec{r}_m^{(2)}(t + \delta t) = \vec{r}_m(t) + (-1 \vec{v}_m + 2 \vec{v}_m^{(1)}) \delta t$
5. compute velocity $\vec{v}_m^{(2)}$ at position $\vec{r}_m^{(2)}$
6. compute new position: $\vec{r}_m(t + \delta t) = \vec{r}_m(t) + (\frac{1}{6} \vec{v}_m + \frac{4}{6} \vec{v}_m^{(1)} + \frac{1}{6} \vec{v}_m^{(2)}) \delta t$

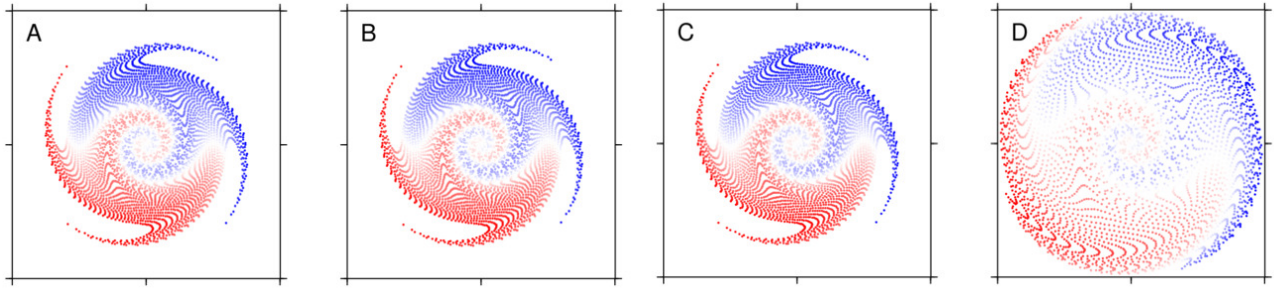
The following example is borrowed from [825], itself borrowed from Fullsack [426, Section 5.4]. It is a whirl flow [964], a flow with rotational symmetry in which concentric layers of material rotate around a centre with an angular velocity:

$$\omega(r) = \omega_0 \frac{r}{r_0} \exp\left(-\frac{r}{r_0}\right)$$

The box is $[-0.5, 0.5] \times [-0.5, 0.5]$, $r_0 = 0.25$, $\omega_0 = 0.3$ and $\delta t = 1$. 60×60 particles are regularly positioned inside the $[-0.3, 0.3] \times [-0.3, 0.3]$ square. Maierova [825] has carried out this experiment for the above Runge-Kutta methods.



Model domain with particles colored at three different time-steps: (A) $t = 0$ (initial position of particles), (B) $t = 50$, and (C) $t = 200$. The advection is computed using the fourth-order Runge-Kutta scheme. Taken from [825]



The same plot as above, but for different advection schemes at $t = 100$. Advection was computed using (A) the fourth-order Runge-Kutta scheme, (B) the mid- point method, (C) Heun's method and (D) the explicit Euler method. Taken from [825]

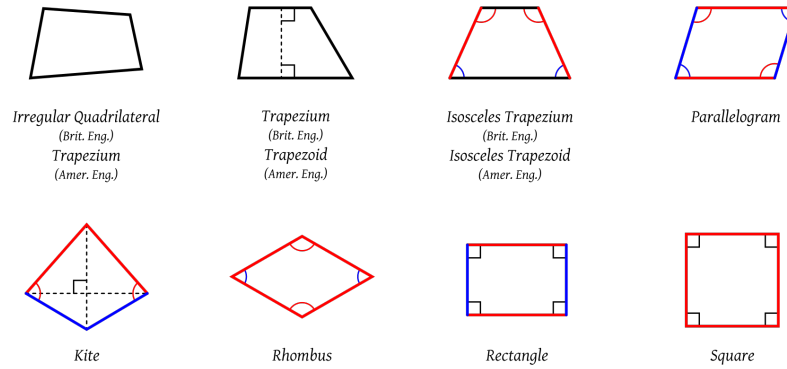
9.11 Am I in or not? - finding reduced coordinates

It is quite common that at some point one must answer the question: "Given a mesh and its connectivity on the one hand, and the coordinates of a point on the other, how do I accurately and quickly determine in which element the point resides?"

One typical occurrence of such a problem is linked to the use of the Particle-In-Cell technique: particles are advected and move through the mesh, and need to be localised at every time step. This question could arise in the context of a benchmark where certain quantities need to be measured at specific locations inside the domain.

Two-dimensional space

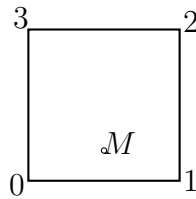
We shall first focus on quadrilaterals. There are many kinds of quadrilaterals as shown hereunder:



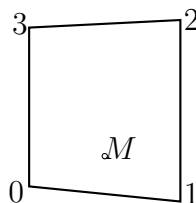
Taken from Wikipedia <https://en.wikipedia.org/wiki/Quadrilateral#/media/File:Quadrilaterals.svg>

The trivial case of rectangular elements Testing whether the point M is inside the element is trivial. For $x_0 \leq x_M \leq x_2$ and $y_0 \leq y_M \leq y_2$, its reduced coordinates are given by

$$\begin{aligned}
 r_M &= \frac{2}{x_2 - x_0}(x_M - x_0) - 1 = \frac{2}{h_x}(x_M - x_0) - 1 \\
 s_M &= \frac{2}{y_2 - y_0}(y_M - y_0) - 1 = \frac{2}{h_y}(y_M - y_0) - 1
 \end{aligned} \tag{9.34}$$



An intermediate case We make the following assumption that the lateral sides of the element are vertical while the bottom and top are not necessarily horizontal:



Because the sides are verical then if $x_0 \leq x_M \leq x_2$ then

$$r_M = \frac{2}{x_2 - x_0}(x_M - x_0) - 1$$

Then, if M is inside the element then its y coordinate is given by

$$y_M = \sum_i \mathcal{N}_i(r_M, s_M) y_i$$

where \mathcal{N}_i are the four Q_1 basis functions associated to the vertices. Assuming we know r_M then we can solve for s_M :

$$\begin{aligned} y_M &= \frac{1}{4}(1 - r_M)(1 - s_M)y_0 + \frac{1}{4}(1 + r_M)(1 - s_M)y_1 + \frac{1}{4}(1 + r_M)(1 + s_M)y_2 + \frac{1}{4}(1 - r_M)(1 + s_M)y_3 \\ &= \frac{1}{4}[(1 - r)y_0 + (1 + r)y_1 + (1 + r)y_2 + (1 - r)y_3 + s_M[-(1 - r)y_0 - (1 + r)y_1 + (1 + r)y_2 + (1 - r)y_3]] \end{aligned}$$

or,

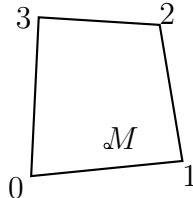
$$s_M = \frac{4y_M - [(1 - r_M)y_0 + (1 + r_M)y_1 + (1 + r_M)y_2 + (1 - r_M)y_3]}{-(1 - r_M)y_0 - (1 + r_M)y_1 + (1 + r_M)y_2 + (1 - r_M)y_3}$$

If the obtained value is in $[-1, 1]$ then the point M is in the element. Verification: when $y_1 = y_0$ and $y_2 = y_3$ then

$$\begin{aligned} s_M &= \frac{4y_M - [(1 - r_M)y_0 + (1 + r_M)y_0 + (1 + r_M)y_3 + (1 - r_M)y_3]}{-(1 - r_M)y_0 - (1 + r_M)y_0 + (1 + r_M)y_3 + (1 - r_M)y_3} \\ &= \frac{4y_M - [2y_0 + 2y_3]}{-2y_0 + 2y_3} \\ &= \frac{1}{y_3 - y_0}[2y_M - (y_0 + y_3)] \\ &= \frac{1}{y_3 - y_0}[2y_M - 2y_0 + y_0 - y_3] \\ &= \frac{2}{y_3 - y_0}(y_M - y_0) - 1 \end{aligned} \tag{9.35}$$

which is the expression that corresponds to a rectangular element as seen previously.

A generic quadrilateral We wish to arrive at a single algorithm which is applicable to all quadrilaterals and we now focus on an irregular quadrilateral (no face is parallel to the axis of the coordinate system).



Several rather simple options exist:

- we could subdivide the quadrilateral into two triangles and check whether point M is inside any of them (as it turns out, this problem is rather straightforward for triangles. Simply google it.)

- We could check that point M is always on the left side of segments $0 \rightarrow 1$, $1 \rightarrow 2$, $2 \rightarrow 3$, $3 \rightarrow 0$.
- ...

Any of these approaches will work although some might be faster than others. In three-dimensions all will however become cumbersome to implement and might not even work at all. Fortunately, there is an elegant way to answer the question, as detailed in the following subsection, which works both in 2D and 3D.

Three-dimensional space

If point M is inside the quadrilateral, there exist a set of reduced coordinates $r, s, t \in [-1 : 1]^3$ such that

$$\sum_{i=1}^4 \mathcal{N}_i(r_M, s_M, t_M) x_i = x_M \quad \sum_{i=1}^4 \mathcal{N}_i(r_M, s_M, t_M) y_i = y_M \quad \sum_{i=1}^4 \mathcal{N}_i(r_M, s_M, t_M) z_i = z_M$$

This can be cast as a system of three equations and three unknowns. Unfortunately, each basis function \mathcal{N}_i contains a term rst (as well as rs , rt , and st) so that it is not a linear system. We must then use an iterative technique: the algorithm starts with a guess for values r_M, s_M, t_M and improves on their value iteration after iteration. In what follows the subscript M is dropped from r, s, t .

The classical way of solving nonlinear systems of equations is Newton's method. We can rewrite the equations above as $\mathbf{F}(r, s, t) = 0$:

$$\begin{aligned} \sum_{i=1}^8 \mathcal{N}_i(r, s, t) x_i - x_M &= 0 \\ \sum_{i=1}^8 \mathcal{N}_i(r, s, t) y_i - y_M &= 0 \\ \sum_{i=1}^8 \mathcal{N}_i(r, s, t) z_i - z_M &= 0 \end{aligned} \tag{9.36}$$

or,

$$\begin{aligned} F_r(r, s, t) &= 0 \\ F_s(r, s, t) &= 0 \\ F_t(r, s, t) &= 0 \end{aligned}$$

so that we now have to find the zeroes of continuously differentiable functions $\mathbf{F} : \mathbb{R} \rightarrow \mathbb{R}$. The recursion is simply:

$$\begin{pmatrix} r_{k+1} \\ s_{k+1} \\ t_{k+1} \end{pmatrix} = \begin{pmatrix} r_k \\ s_k \\ t_k \end{pmatrix} - J_F(r_k, s_k, t_k)^{-1} \begin{pmatrix} F_r(r_k, s_k, t_k) \\ F_s(r_k, s_k, t_k) \\ F_t(r_k, s_k, t_k) \end{pmatrix}$$

where J the Jacobian matrix:

$$J_F(r_k, s_k, t_k) = \begin{pmatrix} \frac{\partial F_r}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_r}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_r}{\partial t}(r_k, s_k, t_k) \\ \frac{\partial F_s}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_s}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_s}{\partial t}(r_k, s_k, t_k) \\ \frac{\partial F_t}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_t}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_t}{\partial t}(r_k, s_k, t_k) \end{pmatrix}$$

$$= \begin{pmatrix} \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial r}(r_k, s_k, t_k) x_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial s}(r_k, s_k, t_k) x_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial t}(r_k, s_k, t_k) x_i \\ \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial r}(r_k, s_k, t_k) y_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial s}(r_k, s_k, t_k) y_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial t}(r_k, s_k, t_k) y_i \\ \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial r}(r_k, s_k, t_k) z_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial s}(r_k, s_k, t_k) z_i & \sum_{i=1}^8 \frac{\partial \mathcal{N}_i}{\partial t}(r_k, s_k, t_k) z_i \end{pmatrix}$$

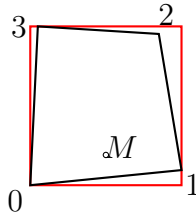
In practice, we solve the following system:

$$J_F(r_k, s_k, t_k) \left[\begin{pmatrix} r_{k+1} \\ s_{k+1} \\ t_{k+1} \end{pmatrix} - \begin{pmatrix} r_k \\ s_k \\ t_k \end{pmatrix} \right] = - \begin{pmatrix} F_r(r_k, s_k, t_k) \\ F_s(r_k, s_k, t_k) \\ F_t(r_k, s_k, t_k) \end{pmatrix}$$

Finally, the algorithm goes as follows:

- set guess values for r, s, t (typically 0)
- loop over $k=0, \dots$
- Compute $\text{rhs} = -\mathbf{F}(r_k, s_k, t_k)$
- Compute matrix $J_F(r_k, s_k, t_k)$
- solve system for (dr_k, ds_k, dt_k)
- update $r_{k+1} = r_k + dr_k, s_{k+1} = s_k + ds_k, t_{k+1} = t_k + dt_k$
- stop iterations when (dr_k, ds_k, dt_k) is small
- if $r_k, s_k, t_k \in [-1, 1]^3$ then M is inside.

This method converges quickly but involves iterations, and multiple solves of 3×3 systems which, when carried out for each marker and at each time step can prove to be expensive. A simple modification can be added to the above algorithm: iterations should be carried out *only* when the point M is inside of a cuboid of size $[\min_i x_i : \max_i x_i] \times [\min_i y_i : \max_i y_i] \times [\min_i z_i : \max_i z_i]$ where the sums run over the vertices of the element. In 2D this translates as follows: only carry out Newton iterations when M is inside the red rectangle!



Note that the algorithm above extends to high degree elements such as Q_2 and higher, even with curved sides. As shown in the 2D case if the element is a cuboid or if all its lateral faces are vertical then one can compute the reduced coordinates without using an iterative method.

Three-dimensional space - special case

We assume that the mesh is such that the cross section of all Q_1 elements is a rectangle in the xy -plane.

Let (x, y, z) be a point inside the element. The global coordinates x, y, z are obtained from the reduced coordinates r, s, t via the basis the basis functions:

$$x = \sum_{i=1}^8 \mathcal{N}_i(r, s, t) x_i \quad y = \sum_{i=1}^8 \mathcal{N}_i(r, s, t) y_i \quad z = \sum_{i=1}^8 \mathcal{N}_i(r, s, t) z_i \quad (9.37)$$

Let

$$\begin{aligned} \vec{v}_1 &= (+1, +1, +1, +1, +1, +1, +1, +1) \\ \vec{v}_2 &= (-1, +1, +1, -1, -1, +1, +1, -1) \\ \vec{v}_3 &= (-1, -1, +1, +1, -1, -1, +1, +1) \\ \vec{v}_4 &= (-1, -1, -1, -1, +1, +1, +1, +1) \\ \vec{v}_5 &= (+1, -1, +1, -1, +1, -1, +1, -1) \\ \vec{v}_6 &= (+1, -1, -1, +1, -1, +1, +1, -1) \\ \vec{v}_7 &= (+1, +1, -1, -1, -1, -1, +1, +1) \end{aligned}$$

and

$$\begin{aligned} \vec{x} &= (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8) \\ \vec{y} &= (y_1, y_2, y_3, y_4, y_5, y_6, y_7, y_8) \\ \vec{z} &= (z_1, z_2, z_3, z_4, z_5, z_6, z_7, z_8) \end{aligned}$$

then Eqs. (9.37) can also be written

$$\begin{aligned} x &= \frac{1}{8} (\vec{v}_1 + r\vec{v}_2 + s\vec{v}_3 + t\vec{v}_4 + rs\vec{v}_5 + rt\vec{v}_6 + st\vec{v}_7) \cdot \vec{x} \\ y &= \frac{1}{8} (\vec{v}_1 + r\vec{v}_2 + s\vec{v}_3 + t\vec{v}_4 + rs\vec{v}_5 + rt\vec{v}_6 + st\vec{v}_7) \cdot \vec{y} \\ z &= \frac{1}{8} (\vec{v}_1 + r\vec{v}_2 + s\vec{v}_3 + t\vec{v}_4 + rs\vec{v}_5 + rt\vec{v}_6 + st\vec{v}_7) \cdot \vec{z} \end{aligned} \quad (9.38)$$

If the element has a rectangular cross-section $s_x \times s_y$ then

$$\begin{aligned} \vec{x} &= (x_0, x_0 + s_x, x_0 + s_x, x_0, x_0, x_0 + s_x, x_0 + s_x, x_0) \\ \vec{y} &= (y_0, y_0, y_0 + s_y, y_0 + s_y, y_0, y_0, y_0 + s_y, y_0 + s_y) \end{aligned}$$

which yields

$$\begin{aligned} r &= 2 \frac{x - x_0}{s_x} - 1 \\ s &= 2 \frac{y - y_0}{s_y} - 1 \end{aligned}$$

Since the local coordinates r and s can be easily computed, one can use Eq. (9.38) to obtain t :

$$t = \frac{8z - (\vec{v}_1 + r\vec{v}_2 + s\vec{v}_3 + rs\vec{v}_5) \cdot \vec{z}}{(\vec{v}_4 + r\vec{v}_6 + s\vec{v}_7) \cdot \vec{z}}$$

9.12 Error measurements and convergence rates

errors.tex

What follows is written in the case of a two-dimensional model. Generalisation to 3D is trivial. What follows is mostly borrowed from [1254].

When measuring the order of accuracy of the primitive variables \vec{v} and p , it is standard to report errors in both the L_1 and the L_2 norm. For a scalar quantity Ψ , the L_1 and L_2 norms are computed as

$$\|\Psi\|_1 = \int_V |\Psi| dV \quad \|\Psi\|_2 = \sqrt{\int_V \Psi^2 dV} \quad (9.39)$$

For a vector quantity $\vec{k} = (k_x, k_y)$ in a two-dimensional space, the L_1 and L_2 norms are defined as:

$$\|\vec{k}\|_1 = \int_V (|k_x| + |k_y|) dV \quad \|\vec{k}\|_2 = \sqrt{\int_V (k_x^2 + k_y^2) dV} \quad (9.40)$$

To compute the respective norms the integrals in the above norms can be approximated by splitting them into their element-wise contributions. The element volume integral can then be easily computed by numerical integration using Gauss-Legendre quadrature.

The respective L_1 and L_2 norms for the pressure error can be evaluated via

$$e_p^h|_1 = \sum_{i=1}^{n_e} \sum_{q=1}^{n_q} |e_p^h(\vec{r}_q)| w_q |J_q| \quad e_p^h|_2 = \sqrt{\sum_{i=1}^{n_e} \sum_{q=1}^{n_q} |e_p^h(\vec{r}_q)|^2 w_q |J_q|} \quad (9.41)$$

where $e_p^h(\vec{r}_q) = p^h(\vec{r}_q) - p(\vec{r}_q)$ is the pressure error evaluated at the q -th quadrature associated with the i th element. n_e and n_q refer to the number of elements and the number of quadrature points per element. w_q and J_q are the quadrature weight and the Jacobian associated with point q .

The velocity error $e_{\vec{v}}^h$ is evaluated using the following two norms

$$e_{\vec{v}}^h|_1 = \sum_{i=1}^{n_e} \sum_{q=1}^{n_q} [|e_u^h(\vec{r}_q)| + |e_v^h(\vec{r}_q)|] w_q |J_q| \quad e_{\vec{v}}^h|_2 = \sqrt{\sum_{i=1}^{n_e} \sum_{q=1}^{n_q} [|e_u^h(\vec{r}_q)|^2 + |e_v^h(\vec{r}_q)|^2] w_q |J_q|} \quad (9.42)$$

where $e_u^h(\vec{r}_q) = u^h(\vec{r}_q) - u(\vec{r}_q)$ and $e_v^h(\vec{r}_q) = v^h(\vec{r}_q) - v(\vec{r}_q)$.

Another norm is very rarely used in the geodynamics literature but is preferred in the Finite Element literature: the H^1 norm. The mathematical basis for this norm and the nature of the $H^1(\Omega)$ Hilbert space is to be found in many FE books [341, 650, 604]. This norm is expressed as follows for a function f such that $f, |\nabla f| \in L^2(\Omega)$ ²⁸

$$\|f\|_{H^1} = \left(\int_{\Omega} (|f|^2 + |\nabla f|^2) d\Omega \right)^{1/2} \quad (9.43)$$

We then have

$$e_{\vec{v}}^h|_{H^1} = \|\vec{v}^h - \vec{v}\|_{H^1} = \sqrt{\sum_{i=1}^d \int_{\Omega} [(v_i^h - v_i)^2 + \vec{\nabla}(v_i^h - v_i) \cdot \vec{\nabla}(v_i^h - v_i)] d\Omega} \quad (9.44)$$

²⁸https://en.wikipedia.org/wiki/Sobolev_space

where d is the number of dimensions. Note that sometimes the following semi-norm is used [336, 101]:

$$e_{\vec{v}}^h|_{H^1} = \|\vec{v}^h - \vec{v}\|_{H^1} = \sqrt{\sum_{i=1}^d \int_{\Omega} [\vec{\nabla}(v_i^h - v_i) \cdot \vec{\nabla}(v_i^h - v_i)] d\Omega} \quad (9.45)$$

When computing the different error norms for e_p and $e_{\vec{v}}$ for a set of numerical experiments with varying resolution h we expect the error norms to follow the following relationships:

$$e_{\vec{v}}^h|_1 = Ch^{rvL_1} \quad e_{\vec{v}}^h|_2 = Ch^{rvL_2} \quad e_{\vec{v}}^h|_{H^1} = Ch^{rvH^1} \quad (9.46)$$

$$e_p^h|_1 = Ch^{rpL_1} \quad e_p^h|_2 = Ch^{rpL_2} \quad (9.47)$$

where C is a resolution-independent constant and $rpXX$ and $rvXX$ are the convergence rates for pressure and velocity in various norms, respectively. Using linear regression on the logarithm of the respective error norm and the resolution h , one can compute the convergence rates of the numerical solutions.

As mentioned in [336], when finite element solutions converge at the same rates as the interpolants we say that the method is optimal, i.e.:

$$e_{\vec{v}}^h|_{L_2} = \mathcal{O}(h^3) \quad e_{\vec{v}}^h|_{H^1} = \mathcal{O}(h^2) \quad e_p^h|_{L_2} = \mathcal{O}(h^2) \quad (9.48)$$

We note that when using discontinuous pressure space (e.g., P_0 , P_{-1}), these bounds remain valid even when the viscosity is discontinuous provided that the element boundaries conform to the discontinuity.

About extrapolation

Section contributed by W. Bangerth and part of Thieulot and Bangerth [1260] (2022) but it was ultimately not used.

In a number of numerical benchmarks we want to estimate the error $X_h - X^*$ between a quantity X_h computed from the numerical solution \vec{v}_h, p_h and the corresponding value X computed from the exact solution \vec{v}, p . Examples of such quantities X are the root mean square velocity \mathbf{v}_{rms} , but it could also be a mass flux across a boundary, an average horizontal velocity at the top boundary, or any other scalar quantity.

If the exact solution is known, then one can of course compute X from it. On the other hand, we would of course like to assess convergence also in cases where the exact solution is not known. In that case, one can compute an *estimate* X^* for X by way of *extrapolation*. To this end, we make the assumption that asymptotically, X_h converges to X at a fixed (but unknown) rate r , so that

$$e_h = |X_h - X| \approx Ch^r. \quad (9.49)$$

Here, X , C and r are all unknown constants to be determined, although we are not really interested in C . We can evaluate X_h from the numerical solution on successively refined meshes with mesh sizes h , $h/2$, and $h/4$. Then, in addition to (9.49) we also have

$$e_{h/2} = |X_{h/2} - X| \approx C \left(\frac{h}{2}\right)^r, \quad (9.50)$$

$$e_{h/4} = |X_{h/4} - X| \approx C \left(\frac{h}{4}\right)^r. \quad (9.51)$$

Taking ratios of equations (9.49)–(9.51), and replacing the unknown X by an *estimate* X^* , we then arrive at the following equation:

$$\frac{|X_h - X^*|}{|X_{h/2} - X^*|} = \frac{|X_{h/2} - X^*|}{|X_{h/4} - X^*|} = 2^r.$$

If one assumes that X_h converges to X uniformly either from above or below (rather than oscillate around X), then this equation allows us to solve for X^* and r :

$$(X_h - X^*)(X_{h/4} - X^*) = (X_{h/2} - X^*)(X_{h/2} - X^*)$$

$$X_h X_{h/4} - X^* X_{h/4} - X_h X^* + (X^*)^2 = X_{h/2}^2 - 2X^* X_{h/2} + (X^*)^2$$

$$X_h X_{h/4} - X^* X_{h/4} - X_h X^* = X_{h/2}^2 - 2X^* X_{h/2}$$

$$X_h X_{h/4} - X_{h/2}^2 = -2X^* X_{h/2} + X^* X_{h/4} + X_h X^*$$

$$X_h X_{h/4} - X_{h/2}^2 = X^*(-2X_{h/2} + X_{h/4} + X_h)$$

and finally:

$$X^* = \frac{X_h X_{h/4} - X_{h/2}^2}{X_h - 2X_{h/2} + X_{h/4}}, \quad r = \log_2 \frac{X_{h/2} - X^*}{X_{h/4} - X^*}.$$

In the determination of r , we could also have used X_h and $X_{h/2}$, but using $X_{h/2}$ and $X_{h/4}$ is generally more reliable because the higher order terms we have omitted in (9.49) are less visible on finer meshes.

In some cases, however, halving the mesh size multiple times is not really tractable (memory problem, or cpu time). Let us now start again from

$$e_h = |X_h - X| \approx Ch^r. \quad (9.52)$$

and assume that we run two other models at a resolution αh and βh , such that $1 > \alpha > \beta > 0$. In the example above we of course had $\alpha = 1/2$ and $\beta = 1/4$. Then we have

$$e_{\alpha h} = |X_{\alpha h} - X| \approx C(\alpha h)^r, \quad (9.53)$$

$$e_{\beta h} = |X_{\beta h} - X| \approx C(\beta h)^r. \quad (9.54)$$

which leads to

$$\frac{|X_h - X^*|}{|X_{\alpha h} - X^*|} = \frac{Ch^r}{C(\alpha h)^r} = (1/\alpha)^r \quad \text{and} \quad \frac{|X_{\alpha h} - X^*|}{|X_{\beta h} - X^*|} = \frac{C(\alpha h)^r}{C(\beta h)^r} = (\alpha/\beta)^r$$


In order for both to be equal we must have

$$(1/\alpha)^r = (\alpha/\beta)^r \quad \Rightarrow \quad 1/\alpha = \alpha/\beta \quad \Rightarrow \quad \beta = \alpha^2$$

So of course if $\alpha = 1/2$ then $\beta = 1/4$, but now we can also take $\alpha = 3/4$ and then $\beta = 9/16$. Etc ...

In the end, this approach might not be that useful since the mesh sizes would then be $h, 3h/4, 9h/16, 27h/64$ which may be hard to achieve in practice.

9.13 The initial temperature field

 Relevant Literature:

- Thermal gradients in the continental crust [219]
- Simple analytical approximation to the temperature structure in subduction zones [378]
- Thermal structure of subduction zone back arcs [296]
- Thermal Structure of Oceanic Lithosphere [1067]
- thermal structure of a subducting plate with finite length [598]

Single layer with imposed temperature b.c.

Let us take a single layer of material characterised by a heat capacity C_p , a heat conductivity k and a heat production term H .



The Heat transport equation writes

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H \quad (9.55)$$

At steady state and in the absence of a velocity field, assuming that the material properties to be independent of time and space, and assuming that there is no heat production ($H = 0$), this equation simplifies to

$$\Delta T = 0 \quad (9.56)$$

Assuming the layer to be parallel to the x -axis, the temperature is $T(x, y) = T(y) = \alpha y + \beta$. In order to specify the constants α and β , we need two constraints.

At the bottom of the layer $y = y_b$ a temperature T_b is prescribed while a temperature T_t is prescribed at the top with $y = y_t$. This ultimately yields a temperature field in the layer given by

$$T(y) = \frac{T_t - T_b}{y_t - y_b}(y - y_b) + T_b$$

If now the heat production coefficient is not zero, the differential equation reads

$$k \Delta T + H = 0 \quad (9.57)$$

The temperature field is then expected to be of the form

$$T(y) = -\frac{H}{2k}y^2 + \alpha y + \beta \quad (9.58)$$

Supplied again with the same boundary conditions, this leads to

$$\beta = T_b + \frac{H}{2k}y_b^2 - \alpha y_b$$

ie,

$$T(y) = -\frac{H}{2k}(y^2 - y_b^2) + \alpha(y - y_b) + T_b$$

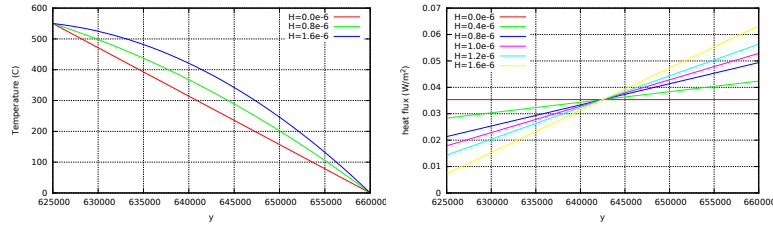
and finally

$$\alpha = \frac{T_t - T_b}{y_t - y_b} + \frac{H}{2k}(y_b + y_t)$$

or,

$$T(y) = -\frac{H}{2k}(y^2 - y_b^2) + \left(\frac{T_t - T_b}{y_t - y_b} + \frac{H}{2k}(y_b + y_t) \right) (y - y_b) + T_b$$

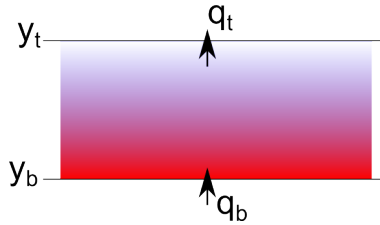
Taking $H = 0$ in this equation obviously yields the temperature field obtained previously. Taking $k = 2.25$, $T_t = 0C$, $T_b = 550C$, $y_t = 660km$, $y_b = 630km$ yields the following temperature profiles and heat fluxes when the heat production H varies:



Looking at the values at the top, which are somewhat estimated to be about $55 - 65 mW/m^2$ [639, table 8.6], one sees that value $H = 0.8e - 6$ yields a very acceptable heat flux. Looking at the bottom, the heat flux is then about $0.03 W/m^2$ which is somewhat problematic since the heat flux at the Moho is reported to be somewhere between 10 and 20 mW/m^2 in [639, table 7.1].

Single layer with imposed heat flux b.c.

Let us now assume that heat fluxes are imposed at the top and bottom of the layer:



We start again from the ODE

$$k\Delta T + H = 0$$

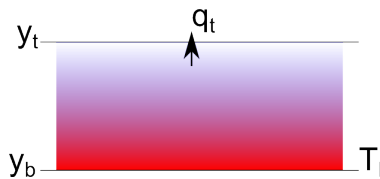
but only integrate it once:

$$k \frac{dT}{dy} + Hy + \alpha = 0$$

At the bottom $q = k(dT/dy)|_{y=y_b} = q_b$ and at the top $q = k(dT/dy)|_{y=y_t} = q_t$ so that

to finish

Single layer with imposed heat flux and temperature b.c.



to finish

Half cooling space

TODO.


 Relevant Literature[384]

Plate model

[856]

McKenzie slab

When doing thermo-mechanical modelling, the initial temperature field in the domain is of prime importance. This is especially true for the temperature in the slab for subduction modelling as its rheological behaviour is strongly temperature-dependent. One could easily design a simple geometrical initial field but it is unlikely to be close to the field of a slowly subducting slab at an angle in a hot mantle.

McKenzie [852] derived such approximate initial field from the steady-state energy equation in two dimensions:

$$\rho C_p \vec{v} \cdot \vec{\nabla} T = k \vec{\nabla}^2 T \quad (9.59)$$

We denote by T_l the temperature at the base of the lithosphere and l its thickness (i.e. the thickness of the slab).

Assuming $\vec{v} = (v_x, 0)$ yields

$$\rho C_p v_x \frac{\partial T}{\partial x} = k \frac{\partial^2 T}{\partial x^2}$$

and substitution of $T' = T/T_l$, $x' = x/l$ and $z' = z/l \in [0, 1]$ in this equation leads to

$$\rho C_p v_x \frac{T_l}{l} \frac{\partial T'}{\partial x'} = k \frac{T_l}{l^2} \left(\frac{\partial^2 T'}{\partial x'^2} + \frac{\partial^2 T'}{\partial z'^2} \right)$$

or

$$\frac{\rho C_p v_x l}{k} \frac{\partial T'}{\partial x'} = \frac{\partial^2 T'}{\partial x'^2} + \frac{\partial^2 T'}{\partial z'^2}$$

and finally (see Eq. 2.3 of [852]):

$$\frac{\partial^2 T'}{\partial x'^2} - 2R \frac{\partial T'}{\partial x'} + \frac{\partial^2 T'}{\partial z'^2} = 0$$

where R is the thermal Reynolds number

$$R = \frac{\rho C_p v_x l}{2k}$$

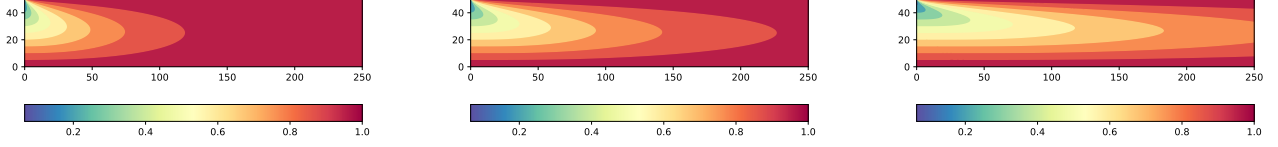
The general solution to this PDE with $T' = 1$ on the top, left and right boundary is

$$T'(x', z') = 1 + \sum_n C_n \exp \left[(R - (R^2 + n^2 \pi^2)^{1/2}) x' \right] \sin(n \pi z')$$

We now must make an assumption about the temperature on the left boundary ($x' = 0$), which is the temperature of the lithosphere. For simplicity McKenzie assumes that $T'(x' = 0, z') = 1 - z'$ so that $C_n = 2(-1)^n / n\pi$ and finally

$$T'(x', z') = 1 + 2 \sum_n \frac{(-1)^n}{n\pi} \exp \left[(R - (R^2 + n^2 \pi^2)^{1/2}) x' \right] \sin(n \pi z') \quad (9.60)$$

Let us build a simple temperature model for a $250\text{km} \times 50\text{km}$ slab, with $\rho = 3000$, $C_p = 1250$, $k = 3$. The python code is available in `images/mckenzie/mckenzie1.py`.

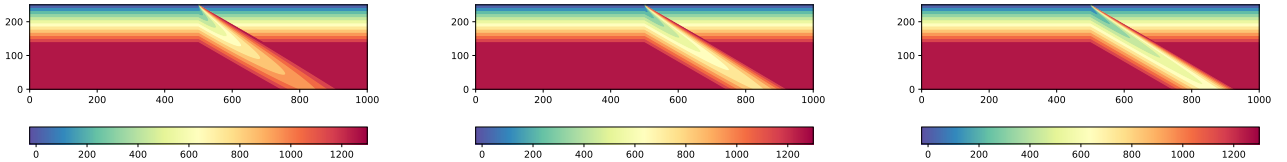


Left to right: Dimensionless temperature T' in a $250\text{km} \times 50\text{km}$ slab for $v_x = 0.5, 1, 2\text{cm/year}$

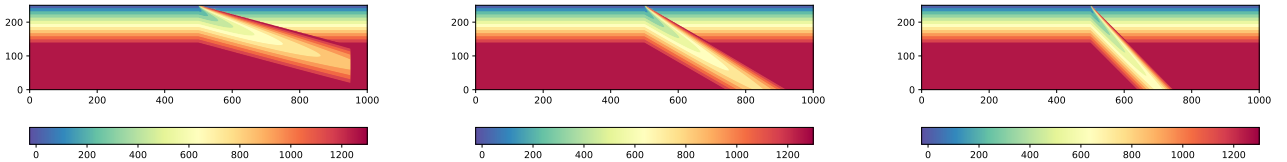
We logically recover the fact that the slower the slab penetrates the mantle the more temperature diffusion dominates over temperature advection. For $v = 0.5\text{cm/year}$ we see that the slab assumes a constant temperature $T' = 1$ at all depths $0 \leq z' \leq 1$ for $x' \geq 125\text{km}$.

Note that this field is a steady-state field, valid for a constant density, heat conductivity and heat capacity, zero heat production, that it implies that the velocity is constant and that the lithosphere temperature is linear.

One can also embed the slab in a more realistic context, a subduction zone, involving a subducting lithosphere, an over-riding plate and a mantle. The domain is $1000\text{km} \times 250\text{km}$. The mantle temperature is set to 1300°C . The slab dip can be varied and so can the velocity. The python code is available in `images/mckenzie/mckenzie2.py`.



Left to right: temperature T for $v_x = 0.5, 1, 2\text{cm/year}$ and $\phi = 30^\circ$.



Left to right: temperature T for $v_x = 1\text{cm/year}$ and $\phi = 15, 30, 45^\circ$.

Initial temperature for global mantle convection models

This is a difficult topic, and Gottschaldt *et al.* [476] list a few issues or facts to take into account:

- Frequent impacts may have determined the heat structure of the outer layers (Arrhenius and Lepland 2000), leading to an early thermally stable stratification.
- A global magma ocean (Solomatov 2000) or several large scale melting events (Kleine *et al.* . 2004) are also conceivable.
- Fractional crystallisation and subsequent overturn has the potential to result in compositionally or thermally stable layering, too (Elkins-Tanton *et al.* 2003; Zaranek and Parmentier 2004)

9.14 The consistent boundary flux (CBF)

cbf.tex

The Consistent Boundary Flux technique was devised to alleviate the problem of the accuracy of primary variables derivatives (mainly velocity and temperature) on boundaries. These derivatives are important since they are needed to compute the heat flux (and therefore the Nusselt number) or dynamic topography and geoid.

The idea was first introduced in Mizukami (1986) [886] and later used in geodynamics in Zhong *et al.* (1993) [1409]. It was finally implemented in the CITCOMS code [1412, 897] and more recently in the ASPECT code (dynamic topography postprocessor). Note that the CBF should be seen as a post-processor step as it does not alter the primary variables values.

The CBF method is implemented and used in [STONE](#) 27. It is also discussed but not explicitly named in Reddy's book [1051, p309]. Also see Larock & Herrmann (1976) [747], Gresho *et al.* (1987) [487], Marshall *et al.* [837].

The CBF applied to the Stokes equation

We start from the strong form:

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{b} = \vec{0} \quad (9.61)$$

and then write the weak form on an element e :

$$\int_{\Omega_e} N_i^\vee \vec{\nabla} \cdot \boldsymbol{\sigma} dV + \int_{\Omega_e} N_i^\vee \vec{b} dV = \vec{0} \quad (9.62)$$

We then use the two equations:

$$\vec{\nabla} \cdot (N\boldsymbol{\sigma}) = N\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{\nabla} N \cdot \boldsymbol{\sigma} \quad (\text{chain rule})$$

$$\int_{\Omega} (\vec{\nabla} \cdot \boldsymbol{\sigma}) dV = \int_{\Gamma} \boldsymbol{\sigma} \cdot \vec{n} dS \quad (\text{divergence theorem})$$

and integrate by parts in order to obtain:

$$\int_{\Gamma} N_i^\vee \boldsymbol{\sigma} \cdot \vec{n} dS - \int_{\Omega_e} \vec{\nabla} N_i^\vee \cdot \boldsymbol{\sigma} dV + \int_{\Omega_e} N_i^\vee \vec{b} dV = \vec{0} \quad (9.63)$$

and since the traction vector \vec{t} is given by $\vec{t} = \boldsymbol{\sigma} \cdot \vec{n}$ we have:

$$\int_{\Gamma_e} N_i^\vee \vec{t} dS = \int_{\Omega_e} \vec{\nabla} N_i^\vee \cdot \boldsymbol{\sigma} dV - \int_{\Omega_e} N_i^\vee \vec{b} dV \quad (9.64)$$

The core idea of the method lies in considering the traction vector as an unknown living on the nodes on the boundary, and assuming we have already solved the Stokes equation and therefore have obtained the velocity and pressure.

Finally, since the traction vector can be expressed as a function of the velocity basis functions on the edge i.e.

$$\vec{t} = \sum_{i=1}^m N_i^\vee \vec{t}_i$$

the left hand term yields an edge (1D) mass matrix \mathbb{M}' (see Section [E](#)).

Remark. In [STONE 27](#) an alternative to equation [9.64](#) is used. Although somewhat inefficient, the elemental matrices \mathbb{K} and \mathbb{G} and the corresponding body force rhs are built and the rhs of the traction equation is computed as follows:

$$\mathbb{M}' \cdot \vec{T} = -\mathbb{K} \cdot \vec{V} - \mathbb{G} \cdot \vec{P} + \vec{f}$$

where \vec{T} is the vector of assembled tractions which we want to compute and \vec{V} and \vec{P} are the solutions of the Stokes problem.

Remark. The assembled mass matrix is tri-diagonal and can be easily solved with a Conjugate Gradient method.

Remark. With a trapezoidal integration rule (i.e. Gauss-Lobatto - see [Section 4.2.7](#)) the matrix can even be diagonalised and the resulting matrix is simply diagonal, which results in a very cheap solve as mentioned in [Zhong et al. \(1993\) \[1409\]](#).

The CBF applied to the heat transport equation

We start from the strong form of the heat transfer equation (without the source terms for simplicity):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T$$

The weak form then writes:

$$\int_{\Omega} N^{\theta} \rho C_p \frac{\partial T}{\partial t} dV + \rho C_p \int_{\Omega} N^{\theta} \vec{v} \cdot \vec{\nabla} T dV = \int_{\Omega} N^{\theta} \vec{\nabla} \cdot k \vec{\nabla} T dV$$

Using once again integration by parts and divergence theorem:

$$\int_{\Omega} N \rho C_p \frac{\partial T}{\partial t} dV + \rho C_p \int_{\Omega} N \mathbf{v} \cdot \nabla T dV = \int_{\Gamma} N k \nabla T \cdot \mathbf{n} d\Gamma - \int_{\Omega} \nabla N \cdot k \nabla T dV$$

On the boundary we are interested in the heat flux $\mathbf{q} = -k \nabla T$

$$\int_{\Omega} N \rho C_p \frac{\partial T}{\partial t} dV + \rho C_p \int_{\Omega} N \mathbf{v} \cdot \nabla T dV = - \int_{\Gamma} N \mathbf{q} \cdot \mathbf{n} d\Gamma - \int_{\Omega} \nabla N \cdot k \nabla T dV$$

or,

$$\int_{\Gamma} N \mathbf{q} \cdot \mathbf{n} d\Gamma = - \int_{\Omega} N \rho C_p \frac{\partial T}{\partial t} dV - \rho C_p \int_{\Omega} N \mathbf{v} \cdot \nabla T dV - \int_{\Omega} \nabla N \cdot k \nabla T dV$$

Considering the normal heat flux $q_n = \mathbf{q} \cdot \mathbf{n}$ as an unknown living on the nodes on the boundary,

$$q_n = \sum_{i=1}^2 q_{n|i} N_i$$

so that the left hand term becomes a mass matrix for the basis functions living on the boundary. We have already covered the right hand side terms when building the FE system to solve the heat transport equation, so that in the end

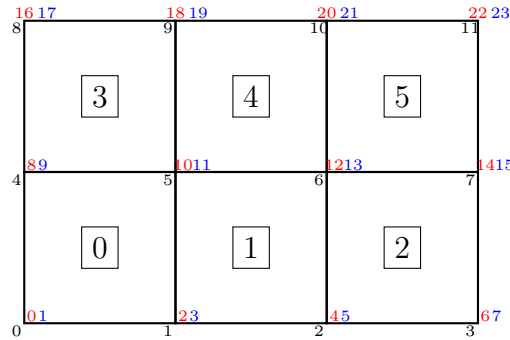
$$\mathbb{M}' \cdot \vec{Q}_n = -\mathbb{M} \cdot \frac{\partial \mathbf{T}}{\partial t} - K_a \cdot \mathbf{T} - K_d \cdot \mathbf{T}$$

where \vec{Q}_n is the assembled vector of normal heat flux components. Note that in all terms the assembly only takes place over the elements along the boundary.

Note that the resulting matrix is symmetric.

Some implementation details for the Stokes equation

What follows is relevant for [STONE](#) 27 which relies on Q_1 shape functions for the velocity. Let us start with a small example, a 3x2 element FE grid:



Red color corresponds to the dofs in the x direction, blue color indicates a dof in the y direction.

We have $nnp=12$, $nel=6$, $N_{femV}=24$. Let us assume that free slip boundary conditions are applied. The boundary conditions `fix_bc` array is then:

```
bc_fix=[T T T T T T T T T T T T T T T T T T T T T T T T]
```

Note that since corners belong to two edges, we effectively prescribed no-slip boundary conditions on those.

why does array contain only T??

We wish to compute the tractions on the boundaries, and more precisely for the dofs for which a Dirichlet velocity boundary condition has been prescribed. The number of (traction) unknowns N_{femTr} is then the number of T in the `bc_fix` array. In our specific case, we have $N_{femTr}=12$. This means that we need for each targeted dof to be able to find its identity/number between 0 and $N_{femTr}-1$. We therefore create the array `bc_nb` which is filled as follows:

finish

```
bc_nb=[T T T T T T T T T T T T T T T T T T T T T T T T]
```

This translates as follows in the code:

```
NfemTr=np.sum(bc_fix)
bc_nb=np.zeros(NfemV,dtype=np.int32)
counter=0
for i in range(0,NfemV):
    if (bc_fix[i]):
        bc_nb[i]=counter
        counter+=1
```

The algorithm is then as follows

- Prepare two arrays to store the matrix M_{cbf} and its right hand side rhs_{cbf}
- Loop over all elements
- For each element touching a boundary, compute the residual vector $R_{el} = -f_{el} + \mathbb{K}_{el}\mathcal{V}_{el} + \mathbb{G}_{el}\mathcal{P}_{el}$
- Loop over the four edges of the element using the connectivity array
- For each edge loop over the number of degrees of freedom (2 in 2D)
- For each edge assess whether the dofs on both ends are target dofs.

G If so, compute the mass matrix M_{edge} for this edge

H extract the 2 values off the element residual vector and assemble these in rhs_{cbf}

I Assemble M_{edge} into NfemTrxNfemTr matrix using bc_nb

```
M_cbf = np.zeros((NfemTr,NfemTr),np.float64) # A
rhs_cbf = np.zeros(NfemTr,np.float64)

for iel in range(0,nel): # B

    ... compute elemental residual ... # C

    #boundary 0-1 # D
    for i in range(0,ndofV): # E
        idof0=2*icon[0,iel]+i
        idof1=2*icon[1,iel]+i
        if (bc_fix[idof0] and bc_fix[idof1]): # F
            idofTr0=bc_nb[idof0]
            idofTr1=bc_nb[idof1]
            rhs_cbf[idofTr0]+=res_el[0+i] # H
            rhs_cbf[idofTr1]+=res_el[2+i]
            M_cbf[idofTr0,idofTr0]+=M_edge[0,0] #
            M_cbf[idofTr0,idofTr1]+=M_edge[0,1] # I
            M_cbf[idofTr1,idofTr0]+=M_edge[1,0] #
            M_cbf[idofTr1,idofTr1]+=M_edge[1,1] #

    #boundary 1-2 #D]

    ...

    #boundary 2-3 #D]

    ...

    #boundary 3-0 #D]

    ...
```


9.15 Computing gradients - the recovery process

recovery.tex

write about recovering accurate strain rate components and heat flux components on the nodes.

Let $\vec{g}(\vec{r})$ be the desired nodal field which we want to be the continuous (for example Q_1) representation of the field $\vec{\nabla} f^h$. Since the derivative of the basis function does not uniquely exist on the nodes we need to design an algorithm to do so. This problem is well known and has been investigated. The main standard techniques are listed hereafter.

refs!

 **Relevant Literature:** check OC Zienkiewicz, B Boroomand, and Jian Zhong Zhu. “Recovery procedures in error estimation and adaptivity: adaptivity in linear problems”. In: *Advances in Adaptive Computational Methods in Mechanics*. 1998, pp. 3–23

Global recovery

The global recovery approach is rather simple: we wish to find \vec{g}^h such that it satisfies

$$\int_{\Omega} \phi \vec{g}^h d\Omega = \int_{\Omega} \phi \vec{\nabla} f^h d\Omega \quad \forall \phi$$

We will then successively replace ϕ by all the basis functions N_i and since we have $g^h = \sum_j N_j g_j$ we then obtain

$$\sum_j \int N_i N_j d\Omega g_j = \int N_i \vec{\nabla} f^h d\Omega$$

or,

$$\mathbb{M} \cdot \vec{g} = \vec{f}$$

Local recovery - centroid average over patch

Local recovery - nodal average over patch

Let j be the node at which we want to compute \vec{g} . Then

$$\vec{g}_j = \vec{g}(\vec{r}_j) = \frac{\sum_{e \text{ adj. to } j} |\Omega_e| (\vec{\nabla} f)_e(\vec{r}_j)}{\sum |\Omega_e|}$$

where $|\Omega_e|$ is the volume of the element and $(\vec{\nabla} f)_e(\vec{r}_j)$ is the gradient of f as obtained with the basis functions inside element e and computed at location \vec{r}_j .

Local recovery - least squares over patch

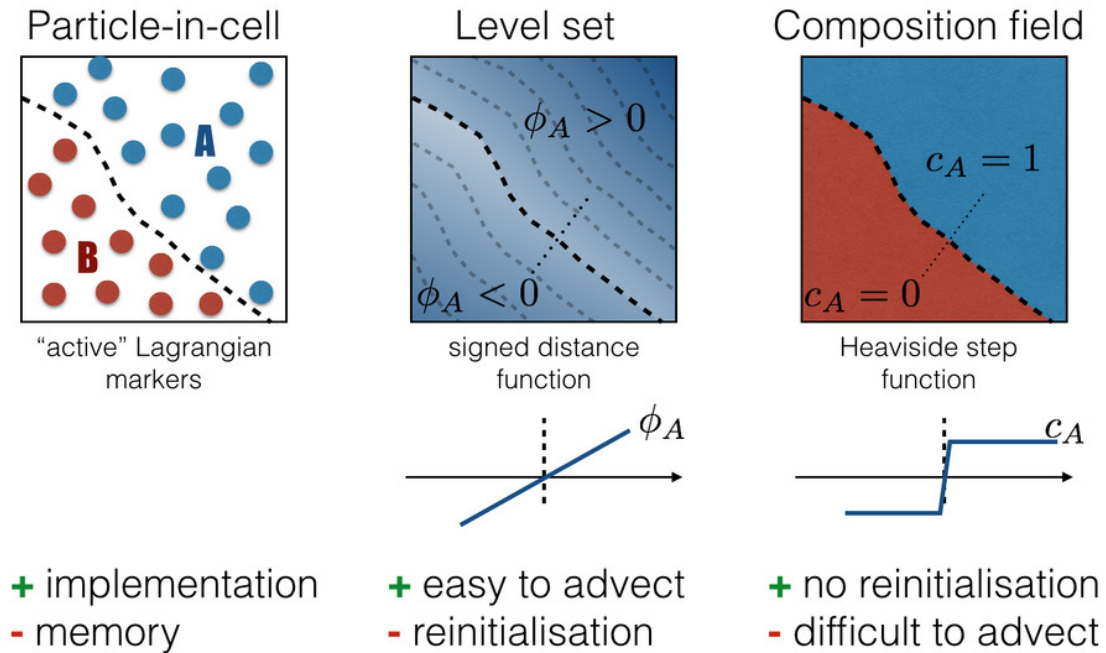
Link to pressure smoothing

When the penalty method is used to solve the Stokes equation, the pressure is then given by $p = -\lambda \vec{\nabla} \cdot \vec{v}$. As explained in section 7.4, the velocity is first obtained and the pressure is recovered by using this equation as a postprocessing step. Since the divergence cannot be computed easily at the nodes, the pressure is traditionally computed in the middle of the elements, yielding an elemental pressure field (remember, we are talking about $Q_1 P_0$ elements here – bi/tri-linear velocity, discontinuous constant pressure)

9.16 Tracking materials and/or interfaces

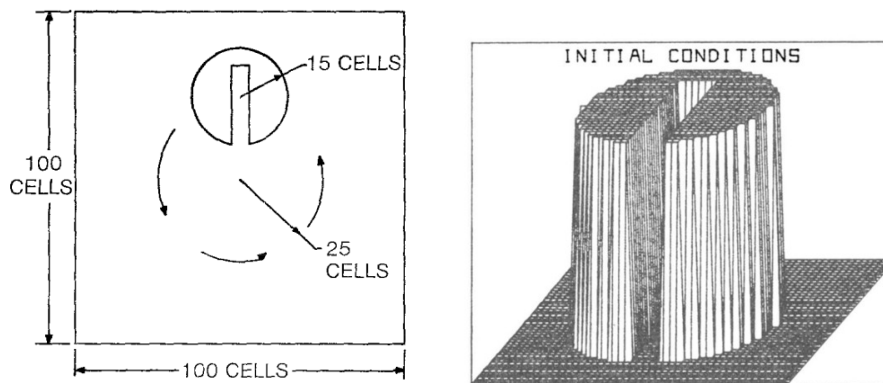
tracking.tex

Unless using a fully Lagrangian formulation, one needs an additional numerical method to represent/track the various materials present in an undeformable (Eulerian) mesh. The figure below (by B. Hillebrand) illustrates the three main methods used in geodynamics.



Note that what follows is applicable to FEM, FDM, etc ...

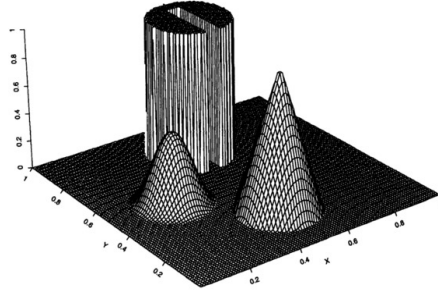
A typical test for advection algorithm is the Zalesak disk [1399]. It is a two dimensional test problem of solid body rotation with a constant angular velocity ω (in rad/sec):



Taken from [1399]. Left: Schematic representation of two dimensional solid body rotation problem. The field inside the cut out has value 3 and it is 1 outside.

The rotational speed is such that one full revolution is effected in 628 cycles. The width of the gap separating the two halves of the cylinder, as well as the maximum extent of the "bridge" connecting the two halves, is 5 cells. Right: Perspective view of initial conditions for the two dimensional! solid body rotation problem. Note that only a 50×50 portion of the mesh centered on the cylinder is displayed.

This benchmark is widely used in the literature [1193, 1220, 1307, 1000, 55, 1407]. Note that the Zalesak disc is often supplemented with a cone and a Gaussian features:



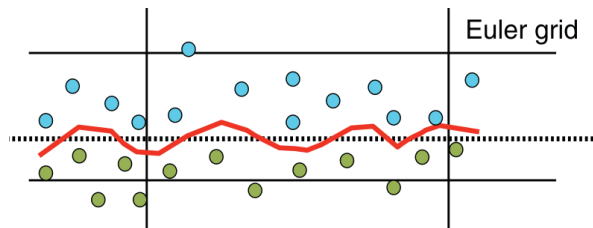
Taken from [777]. Initial data for solid rotation tests

The Particle-in-cell technique

pic.tex

Remark. *The terms 'particle' and 'marker' are commonly (and unfortunately) interchangeably used in the literature in the context of the particle-in-cell technique. However, one should be aware that the marker-and-cell (MAC) technique is something different: it was invented in the early 60's at the Los Alamos Laboratories by Harlow and Welch (1965) [548]. For more information on MAC see the review paper by McKee et al. (2008) [851]. Also, Tackley and King (2003) [1229] talk about the tracer-ratio method in the context of PIC...*

The Particle-in-cell method is by far the most widely used in computational geodynamics. In its most basic form it is a rather simple method to implement and this probably owes to its success and early adoption [1008] in non-parallel codes such as SOPALE [426], I2VIS [453] or CITCOMS [859] (Appendix ??). It has been implemented in ASPECT [438] and the inherent load balancing issues arising from the parallel implementation as well as from the use of Adaptive Mesh Refinement are discussed. It has also been implemented in the MILAMIN code [299] to study LLSVPs [918].



One of the main problems of the PIC method is the fact that the interface between the fluid is not tracked explicitly, and if one uses a random distribution of particles the black dotted line represents the 'real' interface between the fluids while the red line is likely to be the interface one would obtain based on the distribution of particles. Taken from Crameri *et al.* (2012) [285].

Samuel (2018) [1104] does a great job at explaining the core problem with PIC:

The method requires the method requires particle-mesh and mesh-particle mappings to be specified. These critical operations constitute a major source of inaccuracy in the PIC solution [891, 352, 1254]. Indeed, while the Lagrangian advection alone is not prone to significant numerical diffusion, particle-mesh mappings can introduce important amounts of dissipation. This is particularly true when the spatial distribution of particles is not homogeneous, leading to areas in the vicinity of gridpoints that are not sufficiently well sampled by particles, and other regions where the domain is oversampled by particles. This recurrent sampling problem develops in regions characterized by strong deformation, and concerns both compressible and incompressible flow [1337, 1021]. The non-homogeneous sampling has two main origins.

- The first one corresponds to inaccuracies in advecting the Lagrangian particles [868]. This aspect has drawn the attention of a few recent studies [1337, 1021], which have proposed the use of conservative schemes to map velocity components from the Eulerian grid to the Lagrangian particles during their advection. Such schemes have shown to significantly improve the accuracy of the interpolation, and result in a considerably more homogeneous spatial sampling.
- The second origin, which has received less attention, is related to the deforming nature of the flow [898], and is completely independent of the accuracy of the numerical methods for interpolating the velocities at particles' locations. In fact, for a given velocity field, particles should travel along their characteristics, and even in the case of incompressible flows, the distance between characteristics can vary in general, and can strongly diverge or converge in regions characterized by strong deformation. This naturally leads to the development of a non-homogeneous spatial distribution of the Lagrangian particles, even if the particles locations are perfectly known.

The basic methodology goes as follows:

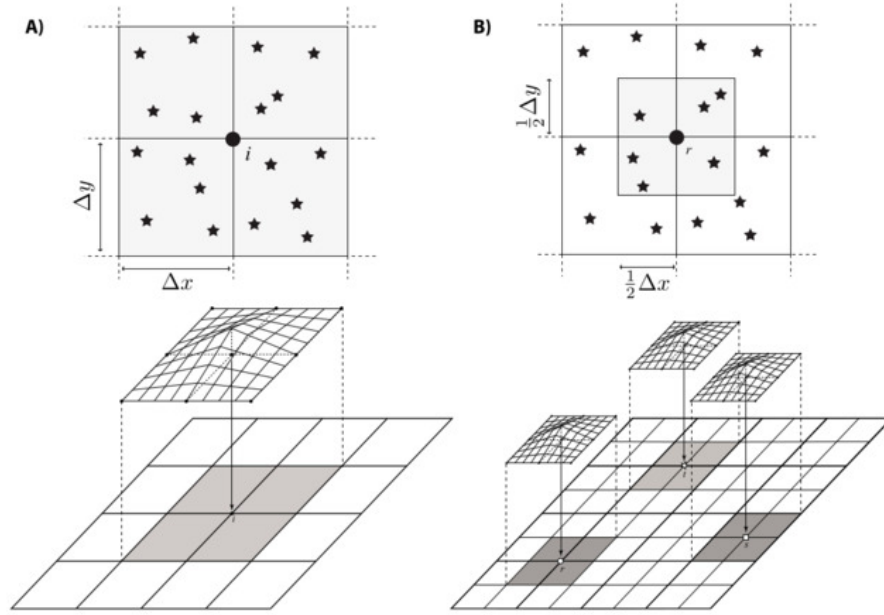
1. distribute particles in the domain,
2. assign a material identity (and/or any other quantity) to each particle,
3. project particle quantities on the nodes of the mesh,
4. solve the Stokes equations for a new velocity field,
5. interpolate the velocity onto the particles,
6. move the particles with their respective velocities,
7. go back to step 3.

As it turns out each step above needs to be carefully executed and is more difficult than it first looks.

Distributing particles in the domain Let us assume we wish to distribute N_p particles in the domain. How large must N_p be? To simplify, one end member could be 'as many particles as possible that fit in memory' while the other end member could be 'one per element/cell on average'. While the former does not necessarily guarantee a desired accuracy while being CPU and memory intensive, the latter will certainly lead to zones in the domain void of particles which will be problematic since the projection onto the mesh might yield zero values or very inaccurate values. How many particles (per element/cell) will be enough? Also, should the particles be randomly distributed in the domain or on some kind of regular grid? See [STONE 13](#).

Taken from Tackley and King (2003) [1229]: "Tracers are initialized on a regular grid with each tracer perturbed from its grid position by a random amount of up to \pm half a grid spacing, in order to eliminate artifacts due to tracer alignment."

Averaging and projection This is a very critical step. Unfortunately, there is no community-wide agreed-upon method. The problem at hand boils down to: at a given location (\vec{r}) in space I need a the value of a field which is carried by the particles. The first step is to find the particle(s) close to this point. If done naively, this is a very costly affair, and begs the question what 'close' means. Finding all particles within a radius R of point \vec{r} can be done very efficiently (e.g. with linked lists, Verlet lists, ...) but the choice of R proves to be critical: if too small, there may not be any particle inside the circle, and if too large there may be many particles inside the circle and the averaging over so many particles in space will prove to be over diffusive. In practice, the FD or FE mesh is used to provide an indication of R . In FDM, the four cells (or quarter cells) around a node represent the volume of space containing the particles whose properties are to be averaged [352] as illustrated in the following figure:



Taken from [352]. The "4-cell" and "1-cell" schemes for projecting properties defined on the markers (denoted by stars) onto a node (denoted by the solid circle). (A) The 4-cell scheme. The support of the interpolating function N_i associated with node i is indicated by the shaded region. Only markers within the support of node i contribute to the projection operation used to define the nodal value at i . The shape of the bilinear interpolation function for node i is indicated in the lower frame. (B) The 1-cell scheme. The thick lines in the lower frame indicate the grid used to discretize the Stokes equations, while the thin lines indicate the grid onto which marker properties are projected. The 1-cell scheme utilizes a compact support of size $\Delta x \times \Delta y$. The support for nodes r , s , t are indicated by the shaded regions. Only markers within the nodal support contribute to the projection operation for that node.

Given that the FEM requires to compute integrals over each element, one could assume that only the particles inside the element will contribute to the average values assigned to the quadrature points (which I coin 'elemental approach').

However, one could also decide to first average the properties onto the nodes before using these nodal values to assign values to the quadrature points (which I coin 'nodal approach'). In this case the FDM approach seen above could apply.

Finally, in both FDM and FEM bi/trilinear basis functions are used for the interpolation as they can be interpreted as weighing functions. Higher order basis functions could also be used but the standard Q_2 basis functions (Section 5.3) are 2-nd order polynomials which can take negative values (as opposed to the Q_1 basis functions which are strictly positive) and this can pose problems: in some cases, although all values to be averaged are positive, their weighed average can be negative. See Section 9.30 for concrete examples.

nodal approach

elemental approach (1) - piece-wise constant interpolation

What follows is written with simplicity in mind, although more mathematical formulations can be found in the literature [438].

Assuming that we have established a list of particles tracking a field $f(\vec{r})$ inside the element we must now compute their average value $\langle f \rangle$. The simplest approach which comes to mind is the arithmetic mean (*am*):

$$\langle f \rangle_{am} = \frac{\sum_{i=1}^n f_i}{n}$$

where n is the number of particles inside the element. In the case where f is the (mass) density ρ , it is indeed what should be used. However, turning now to viscosity η , we know that its value can vary by many orders of magnitude over very short distances. It is then likely that the average runs over values spanning values between 10^{18} Pa s and 10^{25} Pa s. As explained in [1124] the arithmetic averaging tends to 'favour' large values: if the sum runs over 10 particles, 9 carrying the value 10^{25} and 1 carrying the value 10^{19} , the average value is then

$$\langle \eta \rangle = \frac{9 \cdot 10^{25} + 1 \cdot 10^{19}}{10} \simeq 0.9 \cdot 10^{25}$$

which is much much closer to 10^{25} than to 10^{19} . Other averagings are then commonly used, namely the geometric mean (*gm*) and the harmonic mean (*hm*), defined as follows:

$$\langle f \rangle_{gm} = \left(\prod_i f_i \right)^{1/n} \quad \text{or,} \quad \log_{10} \langle f \rangle_{gm} = \frac{\sum_{i=1}^n \log_{10} f_i}{n}$$

and

$$\langle f \rangle_{hm} = \left(\frac{\sum_{i=1}^n \frac{1}{f_i}}{n} \right)^{-1} \quad \text{or,} \quad \frac{1}{\langle f \rangle_{hm}} = \frac{\sum_{i=1}^n \frac{1}{f_i}}{n}$$

The geometric mean can be seen as a form of arithmetic mean of \log_{10} values, while the harmonic mean can be seen as a form of arithmetic mean of the inverse values.

Looking back at the above example, the geometric mean of the viscosities is given by

$$\log \langle \eta \rangle_{gm} = \frac{9 \cdot 25 + 1 \cdot 19}{10} = 24.4 \quad \text{or,} \quad \langle \eta \rangle_{gm} \simeq 2.5 \cdot 10^{24}$$

and the harmonic mean:

$$\langle \eta \rangle_{hm} \simeq \left(\frac{1}{10 \cdot 10^{19}} \right)^{-1} = 10^{20}$$

We see that the harmonic mean tends to favour the small values. Also we recover the known property:

$$\langle f \rangle_{am} \geq \langle f \rangle_{gm} \geq \langle f \rangle_{hm} \quad (9.65)$$

Once a single average value has been computed for the whole element, then all quadrature points are assigned this value.

elemental approach (2) - Least Squares Interpolation One can revisit this topic on the grounds that with high(er) order elements optimal convergence is unlikely to be reached if viscosity (and density) are assumed to be constant inside each element (see Gassm  ller *et al.* (2019) [440]). One could therefore use the least-square method to arrive at a functional representation of the field inside the element which is as close as possible (in the least-squares sense, then) to the particle-based field.

Thielmann *et al.* (2014) [1254] use the Q_2P_{-1} element and introduce an element-wise interpolation scheme based on a least squares fitting of the particle properties and choose the functional to be a linear function to match the pressure space. They define the error ϵ such that

$$\epsilon^2 = \sum_{i=1}^n (\tilde{f}(x_i, y_i) - f_i)^2$$

with $\tilde{f}(x, y) = a + bx + cy$ and proceed to look for the minimum of ϵ^2 , i.e. $\vec{\nabla}(\epsilon^2) = 0$ in the $\{a, b, c\}$ space:

$$\begin{aligned}
0 = \frac{\partial \epsilon^2}{\partial a} &= 2 \sum_i (\tilde{f}(x_i, y_i) - f_i) \\
&= 2 \sum_i (a + bx_i + cy_i - f_i) \\
&= 2 \left[a \sum_i 1 + b \sum_i x_i + c \sum_i y_i - \sum_i f_i \right] \\
0 = \frac{\partial \epsilon^2}{\partial b} &= 2 \sum_i (\tilde{f}(x_i, y_i) - f_i) x_i \\
&= 2 \sum_i (a + bx_i + cy_i - f_i) x_i \\
&= 2 \left[a \sum_i x_i + b \sum_i x_i^2 + c \sum_i x_i y_i - \sum_i x_i f_i \right] \\
0 = \frac{\partial \epsilon^2}{\partial c} &= 2 \sum_i (\tilde{f}(x_i, y_i) - f_i) y_i \\
&= 2 \sum_i (a + bx_i + cy_i - f_i) y_i \\
&= 2 \left[a \sum_i y_i + b \sum_i x_i y_i + c \sum_i y_i^2 - \sum_i y_i f_i \right]
\end{aligned}$$

so

$$\begin{pmatrix} \sum_i 1 & \sum_i x_i & \sum_i y_i \\ \sum_i x_i & \sum_i x_i^2 & \sum_i x_i y_i \\ \sum_i y_i & \sum_i x_i y_i & \sum_i y_i^2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} \sum_i f_i \\ \sum_i x_i f_i \\ \sum_i y_i f_i \end{pmatrix}$$

This method can trivially be extended to three dimensions. It must also be noted that it is not cheap: for each element the matrix and rhs above must be formed and the system solved for a, b, c .

We could also then decide to use a bi-linear function \tilde{f} , i.e.

$$\tilde{f}(x, y) = a + bx + cy + dxy$$

which lies in the Q_1 space of Taylor-Hood quadrilateral elements. In this case the error is

$$\epsilon^2 = \sum_{i=1}^n (\tilde{f}(x_i, y_i) - f_i)^2 = \sum_{i=1}^n (a + bx_i + cy_i + dx_i y_i - f_i)^2$$

and one has to solve a 4×4 system this time:

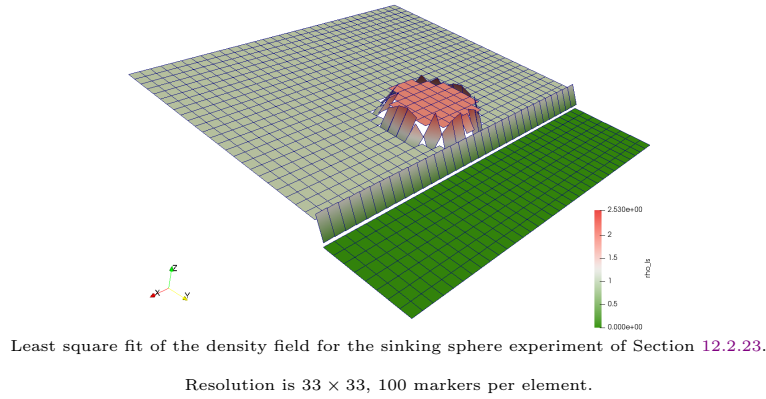
$$\begin{pmatrix} \sum_i 1 & \sum_i x_i & \sum_i y_i & \sum_i x_i y_i \\ \sum_i x_i & \sum_i x_i^2 & \sum_i x_i y_i & \sum_i x_i^2 y_i \\ \sum_i y_i & \sum_i x_i y_i & \sum_i y_i^2 & \sum_i x_i y_i^2 \\ \sum_i x_i y_i & \sum_i x_i^2 y_i & \sum_i y_i^2 & \sum_i x_i^2 y_i^2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = \begin{pmatrix} \sum_i f_i \\ \sum_i x_i f_i \\ \sum_i y_i f_i \\ \sum_i x_i y_i f_i \end{pmatrix}$$

which we write $\mathbf{A} \cdot \vec{c} = \mathbf{b}$. Note that the matrix \mathbf{A} is symmetric. We see that this is a potentially numerically problematic equation. Distances/coordinates in geodynamic calculations are of the order of 100-1000km and viscosities are between 10^{19} and 10^{26} Pa.s. The matrix would contain very large terms, which may compromise the accuracy of the system solve.

Once this linear system (or the previous one) has been solved we have obtained the coefficients a, b, c, d which allow us to compute \tilde{f} anywhere inside the element, and especially at the quadrature points. Once these coefficients have been obtained one can compute \tilde{f} anywhere in the element, and in particular at the quadrature points.

Remark. Using a different (bi)linear function \tilde{f} for each element means that it is likely to be discontinuous from one element to another in regions of high gradients.

There is however one drawback with this approach (linear or bi-linear alike): in the areas of steep gradients the computed coefficients can be such that the function \tilde{f} evaluated on a quadrature point is negative which 1) would be wrong but not numerically dramatic for density, 2) would be wrong and physically and numerically problematic for viscosity (a viscosity cannot be negative, and this would automatically destroy the SPD nature of the viscous block of the Stokes matrix).



This problem is discussed in Thielmann *et al.* (2014) in Section 3.2.1 and they call this "Over- and Under-shooting". A simple (iterative) fix is then designed which insures that the computed value is within user-defined acceptable bounds. This is also mentioned in [440] but the authors explain that this problem was not encountered in the context of the publication.

Remark. One could consider the above least-square approach with $\tilde{f} = a$, i.e. \tilde{f} is a zero-th order polynomial. In this case

$$\epsilon^2 = \sum_{i=1}^n (\tilde{f}(x_i, y_i) - f_i)^2 = \sum_{i=1}^n (a - f_i)^2$$

The gradient becomes

$$\vec{\nabla}(\epsilon^2) = \frac{d\epsilon^2}{da} = \sum_{i=1}^n 2(a - f_i) = 0$$

or $a = \frac{1}{n} \sum_i f_i$. We here recover the arithmetic averaging!

Remark. Two variants of the PIC methods have been proposed: the Deformable PIC (DPIC) by Samuel (2018) [1104], and the multiscale PIC in [31].

Remark. TO BE WRITTEN. A word about the tracer ratio method. [1229]. Trim *et al.* (2020) show a modified method with a tracer repositioning algorithm designed to promote even tracer coverage [1281].

Also look at Yang, Moresi, and Mansour [1378] and Bouffard, Labrosse, Choblet, Fournier, Aubert, and Tackley [119].

See STONE 67 for a concrete example of Particle-In-Cell use and a detailed explanation of its implementation. See also STONE 41 for an implementation of the least square method.

Interpolation of the velocity onto particles

Once the particle i has been localised inside a given element (Section 9.11) and its reduced coordinates (r, s, t) determined, the velocity at this location can be computed through the basis functions:

$$\vec{v}_i = \sum_{k=1}^m N_i(r, s, t) \vec{v}_k$$

This approach is not without problem: while the nodal velocities \vec{v}_k are such that²⁹ $\vec{\nabla} \cdot \vec{v} = 0$ (in the weak sense), the computed velocity \vec{v}_i is not necessarily divergence-free! In order to remedy this, a Conservative Velocity Interpolation (CVI) has been proposed in [1337]. Because the complete derivations for the CVI algorithm is quite large I have decided to make a new section about it (Section 9.31) rather than include it here.

Moving the particles This is discussed in the context of the Runge-Kutta Methods, see Section 9.10.

The level set function technique

lsf.tex

This method was developed in the 80's by Stanley Osher and James Sethian [809]

The Level-set Method (LSM), as it is commonly used in Computational Fluid Dynamics – and especially in Computational Geodynamics – represents a close curve Γ (say, in our case, the interface between two fluids or layers) by means of a function ϕ (called the level-set function, or LSF). Γ is then the zero level-set of ϕ :

$$\Gamma = \{(x, y) \mid \phi(x, y) = 0\} \quad (9.66)$$

The convention is that $\phi > 0$ inside the region delimited by Γ and $\phi < 0$ outside. The function value indicates on which side of the interface a point is located (negative or positive) and this is used to identify materials.

Furthermore, if the curve Γ moves with a velocity \vec{v} , then it satisfies the following equation:

$$\frac{\partial \phi}{\partial t} + \vec{v} \cdot \vec{\nabla} \phi = 0 \quad (9.67)$$

The level set function is generally chosen to be a signed distance function, i.e. $|\vec{\nabla} \phi| = 1$ everywhere and its value is also the distance to the interface. The function value indicates on which side of the interface a point is located (negative or positive) and this is used to identify materials.

As explained in [571], the level-set function ϕ is advected with the velocity \vec{v} which is obtained by solving the Stokes equations. This velocity does not guarantee that after an advection step the signed distance quality of the LSF is preserved. The LSF then needs to be corrected, which is also called reinitialisation. Finally, solving the advection equation must be done in an accurate manner both in time and space, so that so-called ENO (essentially non-oscillatory) schemes are often employed for the space derivative [962, 1103].

The level set method has not often been used in the geodynamics community with some notable exceptions. Bourguin et al use this method combined with Finite Differences to model lava flows [124, 123, 521, 499]. Braun *et al.* use a so-called particle based level set methodology in their FEM code in conjunction with AMR [136]. Zlotnik *et al.* coupled the X-FEM method with level set functions to model slab break-off and Rayleigh-Taylor Diapirism [1441]. This same particle level sets are studied by Samuel and Evonuk and applied to geophysical flows [1103]. In Suckale *et al.* (2010)

²⁹for incompressible flows, of course

[1218, 1217] the authors investigate simulating buoyancy-driven flow in the presence of large viscosity contrasts. Hale *et al.* (2010) [522] use the LSM in 3D and study the dynamics of slab tear faults. An overview of the method and applications can be found in [961].

Several improvements upon the original LSM have been proposed, such as for instance the conservative level set of [1407]. The most notable difference between CLS method originally proposed by Olsson *et al.* [959, 960] and standard LS method lies in the choice of LS function. Instead of the signed distance function, the CLS methods employ the Heaviside function $H(\phi)$

$$H(\phi) = \begin{cases} 1 & \phi > 0 \\ 1/2 & \phi = 0 \\ 0 & \phi < 0 \end{cases}$$


where ϕ is the signed distance function as in the LSM. In practice, a hyperbolic tangent function is used:

$$H(\phi) = \frac{1}{2}(1 + \tan(\phi/2\epsilon))$$

where ϵ defines the spreading width of H . In the case where there are only two fluids (i.e. a single level set is sufficient), the material properties such as density and viscosity are computed as follows:

$$\rho = \rho_1 + (\rho_2 - \rho_1)H(\phi)$$

$$\eta = \eta_1 + (\eta_2 - \eta_1)H(\phi)$$

 **Relevant Literature:** [1307, 1308, 873, 1306].

- Review of level-set methods [462]
- Interactive 3-D computation of fault surfaces using level sets [662]

The field/composition technique

This is the approach taken by the ASPECT developers [732, 560]. Each material i is represented by a compositional field c_i , which takes values between 0 and 1. Each compositional field is then advected with the (prescribed or computed) Stokes velocity [241]:

$$\frac{\partial c_i}{\partial t} + \mathbf{v} \cdot \nabla c_i = 0 \tag{9.68}$$

The value at a point (Finite element node or quadrature point) is 1 if it is in the domain covered by the material i , and 0 otherwise. In one dimension, each compositional field is a Heaviside function. This approach is somewhat similar to the LSM but the field is essentially discontinuous across the interface, which makes it very difficult to advect. On the plus side, compositional fields need not be reinitialised, as opposed to LSF's.

Accurate numerical advection is a notoriously difficult problem. Unless very specialised techniques are used it often yields undershoot ($c_i < 0$) and overshoot ($c_i > 0$), which ultimately yields mass conservation issues. Also, unless special care is taken, compositional fields tend to become more and more diffuse over time: the SUPG method (Section ??) and the entropy viscosity method [732, 1079] add small amounts of diffusion to dampen the under- and overshoots. This means that at a given point two or more compositions may have values, which require some form of averaging. If under- and overshoots are present, these averagings can become very problematic and even yield meaningless quantities (e.g. negative viscosities).

One rather old and popular filtering approach is the so-called Lenardic and Kaula (1993) [766] filter:

The filtering algorithm for two-component flow is as follows. An initial step distribution in C is assumed with $C = 0$ and $C = 1$ used to distinguish distinct materials. A high-order upwind solution scheme is applied to equation (2), with prescribed initial conditions, resulting in an uncorrected C field. The field is corrected via the following filtering algorithm:

1. The initial sum of all nodal C values is calculated and is assigned to the variable C_{sum0} .
2. Nodal C values below 0 are set to 0 and the peak value below 0 is assigned to the variable C_{min} .
3. Nodal C values above 1 are set to 1 and the peak value above 1 is assigned to the variable C_{max} .
4. Nodal C values less than or equal to the absolute value of C_{min} are set to 0.
5. Nodal C values greater than or equal to $2 - C_{max}$ are set to 1.
6. The sum of all nodal C values is calculated and assigned to the variable C_{sum1} .
7. The number of nodal C values not 1 or 0 is assigned to the variable NUM.
8. The variable DIST is defined as $(C_{sum0} - C_{sum1})/NUM$ and is added to all C values not 1 or 0.

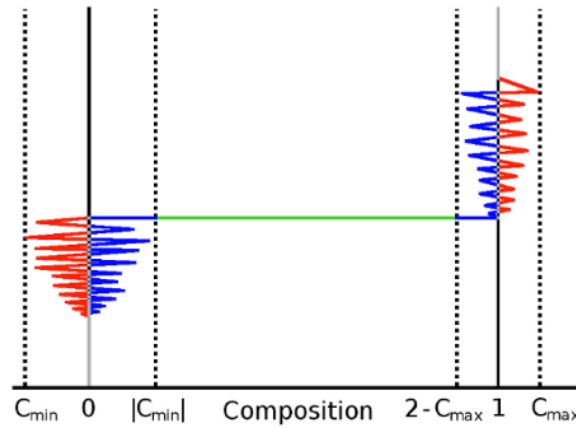
Taken from Lenardic and Kaula [766]

The basic idea of the filtering algorithm is to ensure that ϕ remains within the bounds $0 \leq \phi \leq 1$, and to minimize dispersion error. We refer the reader to Lenardic and Kaula (1993) for the detailed explanation and here give the outline of the algorithm for a discrete property field $\phi = \{\phi_i\}$.

Algorithm 10 A property filtering algorithm

- (1) Compute the initial sum S_0 of all values of ϕ .
 - (2) Find the minimal value ϕ_{min} below 0.
 - (3) Find the maximal value ϕ_{max} above 1.
 - (4) Set $\phi_i = 0$ if $\phi_i \leq |\phi_{min}|$.
 - (5) Set $\phi_i = 1$ if $\phi_i \geq 2 - \phi_{max}$.
 - (6) Compute the sum S_1 of all values of ϕ .
 - (7) Compute the number num of $0 < \phi_i < 1$.
 - (8) Add $dist = (S_1 - S_0)/num$ to all $0 < \phi_i < 1$.
-

From FENICS book



Filtering approach proposed by Lenardic and Kaula (1993). The composition field C is assumed to vary between 0 and 1. Grid points with C -values lower than 0 and greater than 1 are set to 0 and 1, respectively (red). C_{min} and C_{max} are the minimum and maximum spurious values observed. Grid points whose C -value is lower than $|C_{min}|$ or greater than $(2 - C_{max})$ are also set to 0 and 1, respectively (blue). The C -value of all grid points that do not exhibit spurious oscillations (green) is then corrected according to the difference between the original average composition and that computed after the resetting of the spurious values. Taken from Plesa *et al.* (2013) [1003].



Relevant Literature: [1333]

Entropy viscosity method [503]

write about DG approach

The Volume-of-Fluid method

The Volume-Of-Fluid (VOF) method is a fixed-grid approach based on the one-fluid model and considers that the various immiscible fluids (or ‘phases’) can be described as a single fluid whose local physical properties, namely density and viscosity, vary in space and time depending on the volume fraction C_i of each phase i [574, 1391].

The volume fraction of each fluid intrinsically obeys $\sum_{i=1}^n C_i = 1$ where n is the number of phases. Typically, $C_i = 1$ in grid cells filled only with fluid i , and $0 < C_i < 1$ in grid cells cross-cut by an interface. There are two main classes of VOF methods: methods that try to reconstruct exactly the interface between fluids (e.g. [1020]), which requires significant computational time, and methods that do not, such as in JADIM and OpenFOAM. With no interface reconstruction, the thickness of the interfacial region is defined by $0 < C_i < 1$, and typically occupies two to three grid cells.

Relevant Literature:

Hirt and Nichols [574] “Volume of fluid (VOF) method for the dynamics of free boundaries”

Dutta, Sarkar, and Mandal [355] “Ballooning versus curling of mantle plumes: views from numerical models”

Robey and Puckett [1079] “Implementation of a volume-of-fluid method in a finite element code with applications to thermochemical convection in a density stratified fluid in the earth’s mantle”

Louis-Napoléon, Gerbault, Bonometti, Thieulot, Martin, and Vanderhaeghe [811] “3D numerical modeling of crustal polydiapirs with Volume-Of-Fluid methods”

Louis-Napoléon, Bonometti, Gerbault, Martin, and Vanderhaeghe [810] “Models of convection and segregation in heterogeneous partially molten crustal roots with a VOF method—I: flow regimes”

See review of the method in Robey’s phd thesis [1078].

The method of characteristics

ask Arie to write something

[323]

The Marker Chain method

In two dimensions, the idea is quite simple: each interface is discretised by means of a number of Lagrangian points (which may or may not vary in time). The points are numbered and connected (think of the connectivity array of a 1D FEM code). In the case of small deformations, and in the absence of in/out-flow boundaries, the method is reasonably trivial to implement, and each couple of point defines a segment (and therefore its normal vector too) which can then be used to answer the question: “at this location, am I above or below this interface” or “am I this domain or outside this domain” (in the case that the interface does not reach any of the boundaries).

This method becomes somewhat impractical when large deformation occurs or, for example, when a domain splits into two (e.g. slab break off). One interface must then become two, and it requires an algorithm capable of detecting the breakup of the surface and capable of rebuilding/patching the new ones so that they can still be used further. Note that in case of large deformation some markers may get further and further apart from each other which makes for a poor representation of the surface. New markers should then be added but the question of when and where must then be addressed.

Also, switching to three dimensions can prove to be very difficult or simply very costly: the generation of the initial marker position is trivial but their connectivity can be complicated to establish at startup: for instance, a Stokes sphere will require a mesh made of triangles which maps exactly the surface of the sphere (see [1259, 904] for methods on how to efficiently produce such meshes). In

the case of more complex 3D geometries this may prove nearly impossible to do. So will the problem of splitting a surface into two (or merging two domains).

This method is usually coupled to Eulerian meshes (typically with FDM, but not only). It was used in [1366] in the context of salt domes analysis and later in [248, 245]. It is also used in [1309] but little details are given about the algorithms used to track and update the chain in the presence of such large deformation. It is also used (although coupled to level set functions) in the DOUAR code [136] (see Section ??). Having worked myself on this code and having had to produce complex initial triangulated surfaces for simulations (see for example [807]) it is easy to understand why later users of this code did implement the marker-in-cell technique. More recently, it is used to track the free surface position in a FDM code [353, 229].

Finally, Christensen [241] makes the following interesting comment: "One might assume that different methods of representing the discontinuity, for example, by a tracer chain [245] or a cloud of tracers, would solve these problems. However, the difficulties arise not only from the way in which material boundaries are represented. Physically, the rate of shear strain parallel to a rheological boundary is discontinuous. Within the finite element scheme such jump can only be realized at an element boundary. In an Eulerian scheme, where the discontinuity will crosscut the elements, the jump in strain rate must be approximated by a continuous variation, and effectively, the rheological properties on both sides of the discontinuity will be averaged in some way within the element."

It is also used in Tan & Gurnis (tagu07) [1234]: "The composition field is computed using the marker chain method [316, 1309]. The marker chain is advected using a fourth-order predictor-corrector scheme. If the distance between two adjacent markers is greater than a predefined threshold, a new marker is inserted in between them. The marker chain defines the material interface. Because of material entrainment, the length of the marker chain grows exponentially with time. The computational efficiency of the marker chain method severely deteriorates if there is substantial material entrainment, in which case we halt the computation. For some halted models, the marker chain is trimmed to remove excess entrainment, and the computation restarted in order to proceed further. The trimming of the marker chain introduces error in the composition field, but the magnitude of the error is estimated to be small and does not influence the stability of the chemical layer."


Literature: Lin & van Keken (2006) [791, 789, 790, 678, 919]

Hybrid methods

In Braun *et al.* [136] a level set method is presented which is based on a 3-D set of triangulated points, which makes it a hybrid between tracers and level set functions: in the DOUAR code (Appendix ??) the interface is then explicitly tracked by means of the tracers while the LSF is computed on the FE nodes. Although very promising in theory, this method proved to be difficult to use in practice since it requires a) a triangulation of the interfaces at $t = 0$ which is not trivial if the geometries are complex (think about a slab in 3D); b) the addition or removal of tracers because of the interface deformation and the patching of the triangulation; c) the calculation of the distance to the interfaces for each FE node based on the triangle normal vectors. This probably explains why the Particle-In-Cell method was later implemented in this code (pers. comm.). Note that another very similar approach is used in [1103].

Boundary fitted mesh

This method is rather simple to implement and works well for small deformations. It is for instance used by Frehner [417] (see online supplementary material) in which it is stated: "The numerical grid is set up in such a way that the interface between different material phases (two layers in this case) coincides with element boundaries. Hence, each element belongs to a unique material phase and no interpolation is necessary." With such a method, each element is initially attributed a material phase/number and its material properties do not change.

 Relevant Literature: three-dimensional front tracking method using a triangular mesh [1122].

9.17 Static condensation

static_condensation.tex

The idea behind static condensation is quite simple: in some cases, there are dofs belonging to an element which only belong to that element. For instance, the so-called MINI element ($P_1^+ \times P_1$) showcases a bubble function in the middle (see section 7.3). In the following, $\vec{\mathcal{V}}^\star$ corresponds to the list of such dofs inside an element. The discretised Stokes equations on any element looks like:

$$\begin{pmatrix} \mathbb{K} & L & \mathbb{G} \\ L^T & \mathbb{K}^\star & H \\ \mathbb{G}^T & H^T & 0 \end{pmatrix}_e \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{V}}^\star \\ \vec{\mathcal{P}} \end{pmatrix}_e = \begin{pmatrix} \vec{f} \\ \vec{f}^\star \\ \vec{h} \end{pmatrix}_e \quad (9.69)$$

This is only a re-writing of the elemental Stokes matrix where the matrix \mathbb{K} has been split in four parts. Note that the matrix \mathbb{K}^\star is diagonal. check

This can also be re-written in non-matrix form:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + L \cdot \vec{\mathcal{V}}^\star + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (9.70)$$

$$L^T \vec{\mathcal{V}} + \mathbb{K}^\star \cdot \vec{\mathcal{V}}^\star + H \cdot \vec{\mathcal{P}} = \vec{f}^\star \quad (9.71)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} + H^T \vec{\mathcal{V}}^\star = \vec{h} \quad (9.72)$$

The $\vec{\mathcal{V}}^\star$ in the second equation can be isolated:

$$\vec{\mathcal{V}}^\star = \mathbb{K}^{-\star} \cdot (\vec{f}^\star - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}})$$

and inserted in the first and third equations:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + L \left[\mathbb{K}^{-\star} (\vec{f}^\star - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}}) \right] + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (9.73)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} + H^T \left[\mathbb{K}^{-\star} (\vec{f}^\star - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}}) \right] = \vec{h} \quad (9.74)$$

or,

$$(\mathbb{K} - L \cdot \mathbb{K}^{-\star} \cdot L^T) \cdot \vec{\mathcal{V}} + (G - L \cdot \mathbb{K}^{-\star} \cdot H) \cdot \vec{\mathcal{P}} = \vec{f} - L \cdot \mathbb{K}^{-\star} \cdot \vec{f}^\star \quad (9.75)$$

$$(G^T - H^T \cdot \mathbb{K}^{-\star} \cdot L^T) \cdot \vec{\mathcal{V}} - (H^T \cdot \mathbb{K}^{-\star} \cdot H) \cdot \vec{\mathcal{P}} = \vec{h} - H^T \cdot \mathbb{K}^{-\star} \cdot \vec{f}^\star \quad (9.76)$$

i.e.

$$\underline{\mathbb{K}} \cdot \vec{\mathcal{V}} + \underline{\mathbb{G}} \cdot \vec{\mathcal{P}} = \underline{\vec{f}} \quad (9.77)$$

$$\underline{\mathbb{G}}^T \cdot \vec{\mathcal{V}} - \underline{\mathbb{C}} \cdot \vec{\mathcal{P}} = \underline{\vec{h}} \quad (9.78)$$

with

$$\underline{\mathbb{K}} = \mathbb{K} - L \cdot \mathbb{K}^{-\star} \cdot L^T \quad (9.79)$$

$$\underline{\mathbb{G}} = G - L \cdot \mathbb{K}^{-\star} \cdot H \quad (9.80)$$

$$\underline{\mathbb{C}} = H^T \cdot \mathbb{K}^{-\star} \cdot H \quad (9.81)$$

$$\underline{\vec{f}} = \vec{f} - L \cdot \mathbb{K}^{-\star} \cdot \vec{f}^\star \quad (9.82)$$

$$\underline{\vec{h}} = \vec{h} - H^T \cdot \mathbb{K}^{-\star} \cdot \vec{f}^\star \quad (9.83)$$

Note that $\underline{\mathbb{K}}$ is symmetric, and so is the Stokes matrix.

For instance, in the case of the MINI element, the dofs corresponding to the bubble could be eliminated at the elemental level, which would make the Stokes matrix smaller (see book by Braess [128]). However, it is then important to note that static condensation introduces a pressure-pressure term which was not there in the original formulation. This is also presented in the appendix of Karabelas *et al.* [670].

9.18 Measuring incompressibility

The velocity divergence error integrated over the whole element is given by

$$e_{div} = \int_{\Omega} (\vec{\nabla} \cdot \vec{v}^h - \underbrace{\vec{\nabla} \cdot \vec{v}}_{=0}) d\Omega = \int_{\Omega} (\vec{\nabla} \cdot \vec{v}^h) d\Omega \quad (9.84)$$

where Γ_e is the boundary of element e and \vec{n} is the unit outward normal of Γ_e .

Furthermore, one can show that [336]:

$$e_{div} = \int_{\Gamma_e} \vec{v}^h \cdot \vec{n} d\Gamma$$

The reason is as follows and is called the divergence theorem³⁰: suppose a volume V subset of \mathbb{R}^d which is compact and has a piecewise smooth boundary S , and if \vec{F} is a continuously differentiable vector field then

$$\int_V (\vec{\nabla} \cdot \vec{F}) dV = \int_S (\vec{F} \cdot \vec{n}) dS$$

The left side is a volume integral while the right side is a surface integral. Note that sometimes the notation $d\vec{S} = \vec{n} dS$ is used so that $\vec{F} \cdot \vec{n} dS = \vec{F} \cdot d\vec{S}$.

The average velocity divergence over an element can be defined as

$$\langle \vec{\nabla} \cdot \vec{v} \rangle_e = \frac{1}{V_e} \int_{\Omega_e} (\vec{\nabla} \cdot \vec{v}) d\Omega = \frac{1}{V_e} \int_{\Gamma_e} \vec{v} \cdot \vec{n} d\Gamma$$


Note that for elements using discontinuous pressures we shall recover a zero divergence element per element (local mass conservation) while for continuous pressure elements the mass conservation is guaranteed only globally (i.e. over the whole domain), see section 3.13.2 of [488].

Note that one could instead compute $\langle |\vec{\nabla} \cdot \vec{v}| \rangle_e$. Either volume or surface integral can be computed by means of an appropriate Gauss-Legendre quadrature algorithm.

³⁰https://en.wikipedia.org/wiki/Divergence_theorem

9.19 Picard and Newton

explain why our eqs are nonlinear

 **Relevant Literature** Quasi Newton methods [373]. Also check Christensen and Yuen [246] for a succinct N-R explanation in the context of stream fct formulation.

Picard iterations

Let us consider the following system of nonlinear algebraic equations:

$$\mathbb{A}(\vec{X}) \cdot \vec{X} = \vec{b}(\vec{X})$$

Both matrix and right hand side depend on the solution vector \vec{X} .

For many mildly nonlinear problems, a simple successive substitution iteration scheme (also called Picard method) will converge to the solution and it is given by the simple relationship:

$$\mathbb{A}(\vec{X}^n) \cdot \vec{X}^{n+1} = \vec{b}(\vec{X}^n)$$

where n is the iteration number. It is easy to implement:

1. guess \vec{X}^0 or use the solution from previous time step
2. compute \mathbb{A} and \vec{b} with current solution vector \vec{X}^{old}
3. solve system, obtain T^{new}
4. check for convergence (are \vec{X}^{old} and \vec{X}^{new} close enough?)
5. $\vec{X}^{old} \leftarrow \vec{X}^{new}$
6. go back to 2.

There are various ways to test whether iterations have converged. The simplest one is to look at $\|\vec{X}^{old} - \vec{X}^{new}\|$ (in the L_1 , L_2 or maximum norm) and assess whether this term is smaller than a given tolerance ϵ . However this approach poses a problem: in geodynamics, if two consecutively obtained temperatures do not change by more than a thousandth of a Kelvin (say $\epsilon = 10^{-3}\text{K}$) we could consider that iterations have converged but looking now at velocities which are of the order of a cm/year (i.e. $\sim 3 \cdot 10^{-11}\text{m/s}$) we would need a tolerance probably less than 10^{-13}m/s . We see that using absolute values for a convergence criterion is a potentially dangerous affair, which is why one uses a relative formulation (thereby making ϵ a dimensionless parameter):

$$\frac{\|\vec{X}^{old} - \vec{X}^{new}\|}{\|\vec{X}^{new}\|} < \epsilon$$

Another convergence criterion is proposed by Reddy (section 3.7.2) [1051]:

$$\left(\frac{(\vec{X}^{old} - \vec{X}^{new}) \cdot (\vec{X}^{old} - \vec{X}^{new})}{X^{new} \cdot X^{new}} \right)^{1/2} < \epsilon$$

Yet another convergence criterion is used in [1258]: the means $\langle \vec{X}^{old} \rangle$, $\langle \vec{X}^{new} \rangle$ as well as the variances σ^{old} and σ^{new} are computed, followed by the correlation factor R :

$$R = \frac{\langle (\vec{X}^{old} - \langle \vec{X}^{old} \rangle) \cdot (\vec{X}^{new} - \langle \vec{X}^{new} \rangle) \rangle}{\sqrt{\sigma^{old} \sigma^{new}}}$$

Since the correlation is normalised, it takes values between 0 (very dissimilar velocity fields) and 1 (very similar fields). The following convergence criterion is then used: $1 - R < \epsilon$.

write about nonlinear residual

Note that in some instances and improvement in convergence rate can be obtained by use of a relaxation formula where one first solves

$$\mathbb{A}(\vec{X}^n) \cdot \vec{X}^\star = \vec{b}(\vec{X}^n)$$

and then updates \vec{X}^n as follows:

$$\vec{X}^n = \gamma \vec{X}^n + (1 - \gamma) \vec{X}^\star \quad 0 < \gamma \leq 1$$

When $\gamma = 1$ we recover the standard Picard iterations formula above.

9.19.1 Defect correction formulation

Work in progress.

We start from the system to solve:

$$\mathbf{A}(\vec{X}) \cdot \vec{X} = \vec{b}(\vec{X})$$

with the associated residual vector \vec{F}

$$\vec{F}(\vec{X}) = \mathbf{A}(\vec{X}) \cdot \vec{X} - \vec{b}(\vec{X})$$

The Newton-Raphson algorithm consists of two steps:

1. solve $\mathbf{J}_k \cdot \delta \vec{X}_k = -\vec{F}(\vec{X}_k)$, or in the case of the incompressible Stokes equation FEM system:

$$\begin{pmatrix} \mathbf{J}_k^{\mathcal{V}\mathcal{V}} & \mathbf{J}_k^{\mathcal{V}\mathcal{P}} \\ \mathbf{J}_k^{\mathcal{P}\mathcal{V}} & 0 \end{pmatrix} \cdot \begin{pmatrix} \delta \vec{\mathcal{V}}_k \\ \delta \vec{\mathcal{P}}_k \end{pmatrix} = \begin{pmatrix} -\vec{F}_k^{\mathcal{V}} \\ -\vec{F}_k^{\mathcal{P}} \end{pmatrix}$$

2. update $\vec{X}_{k+1} = \vec{X}_k + \alpha_k \delta \vec{X}_k$

The defect correction Picard approach consists of neglecting the derivative terms present in the J terms (Eqs. 16,17,18 of [415]) so that

$$\mathbf{J}_k^{\mathcal{V}\mathcal{V}} \simeq \mathbb{K}_k \quad \mathbf{J}_k^{\mathcal{V}\mathcal{P}} \simeq \mathbb{G} \quad \mathbf{J}_k^{\mathcal{P}\mathcal{V}} \simeq \mathbb{G}^T$$

and step 1 of the above iterations become:

$$\begin{pmatrix} \mathbb{K}_k & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \delta \vec{\mathcal{V}}_k \\ \delta \vec{\mathcal{P}}_k \end{pmatrix} = \begin{pmatrix} -\vec{F}_k^{\mathcal{V}} \\ -\vec{F}_k^{\mathcal{P}} \end{pmatrix}$$

explain picard, defect picard, Newton, line search,

- VV Ermakov and Nikolai Nikolaevich Kalitkin. “The optimal step and regularization for Newton’s method”. In: *USSR Computational Mathematics and Mathematical Physics* 21.2 (1981), pp. 235–242. DOI: 10.1016/0041-5553(81)90022-7

- MS Engelman, Gilbert Strang, and K-J Bathe. “The application of quasi-Newton methods in fluid mechanics”. In: *International Journal for Numerical Methods in Engineering* 17.5 (1981), pp. 707–718
- Dana A Knoll and David E Keyes. “Jacobian-free Newton–Krylov methods: a survey of approaches and applications”. In: *Journal of Computational Physics* 193.2 (2004), pp. 357–397. DOI: 10.1016/j.jcp.2003.08.010
- He Yiqian and Yang Haitian. “Solving inverse couple-stress problems via an element-free Galerkin (EFG) method and Gauss–Newton algorithm”. In: *Finite Elements in Analysis and Design* 46.3 (2010), pp. 257–264. DOI: 10.1016/j.finel.2009.09.009
- Pierre Saramito. “A damped Newton algorithm for computing viscoplastic fluid flows”. In: *Journal of Non-Newtonian fluid mechanics* 238 (2016), pp. 6–15. DOI: 10.1016/j.jnnfm.2016.05.007
- M.R.T. Fraters, W. Bangerth, C. Thieulot, A.C. Glerum, and W. Spakman. “Efficient and Practical Newton Solvers for Nonlinear Stokes Systems in Geodynamic Problems”. In: *Geophy. J. Int.* 218 (2019), pp. 873–894. DOI: 10.1093/gji/ggz183
- Johann Rudi, Yu-hsuan Shih, and Georg Stadler. “Advanced Newton methods for geodynamical models of Stokes flow with viscoplastic rheologies”. In: *Geochemistry, Geophysics, Geosystems* 21.9 (2020), e2020GC009059. DOI: 10.1029/2020GC009059

9.20 Parallel or not?

Rationale

Let us assume that we want to run a simulation of the whole Earth mantle with a constant resolution of 5km. The volume of the mantle is

$$V_{mantle} = \frac{4}{3}\pi(R_{out}^3 - R_{in}^3) \simeq 10^{12} km^3$$

while the volume of an element is $V_e = 125 km^3$ (this is only an average since the tessellation of a hollow sphere with hexahedra yields elements which are not all similar [1259]). Consequently, the number of cells needed to discretise the mantle is

$$N_{el} = \frac{V_{mantle}}{V_e} \simeq 8 \times 10^9$$

We know that the matrix size is approx. 4 times the number of elements in 3D:

$$N \simeq 25 \times 10^9$$

Using between 9 and 125 particles per element (a very conservative number), the total number of particles is then

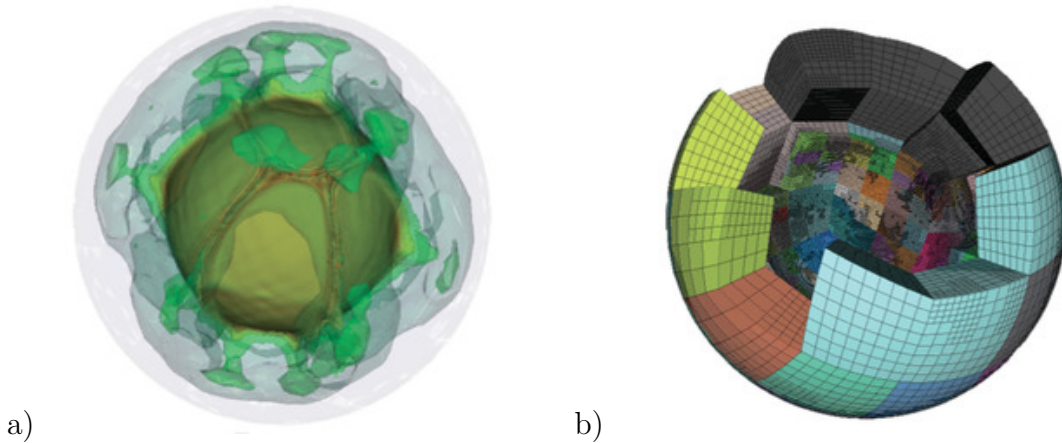
$$N_{particles} \geq 10^{10}$$

The unescapable conclusion is that high-resolution 3D calculations have a very large memory footprint and require extremely long computational times.

The only way to overcome this problem is by resorting to using supercomputers with many processors and large memory capacities.

The idea behind parallel programming is to have each processor carry out only a subset of the total number of operations required. In order to reduce the memory footprint on each processor, only a subset of the computational mesh is known by each: one speaks then of domain decomposition [1275].

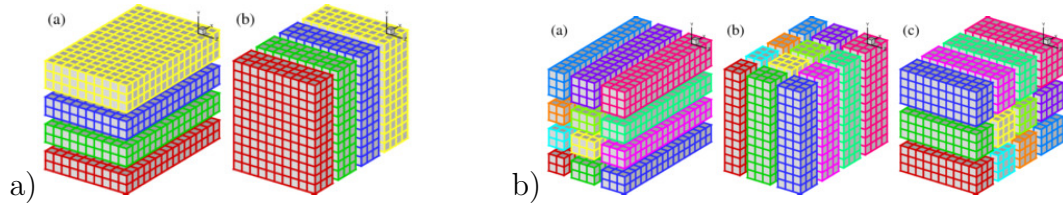
An example of such a large parallel calculation of 3D convection with domain decomposition in a spherical shell can be found in [732]:



a) Isocontours of the temperature field; b) Partitioning of the domain onto 512 proc. The mesh counts 1,424,176 cells. The solution has approximately 54 million unknowns (39 million vel., 1.7 million press., and 13 million temp.)

Basic approaches

In the past, many applications implemented the idea below on the left using 1D domain decomposition (also known as “slab decomposition”). In the following left figure, a 3D domain is arbitrarily chosen to be decomposed in Y and X directions. It can be seen that in state (a), any computations in the X-Z planes can be done in local memories while data along a Y mesh-line is distributed.



Left: 1D domain decomposition example using 4 processors: (a) decomposed in Y direction; (b) decomposed in X direction. Right: 2D domain decomposition example using a 4x3 processor grid.

A 1D decomposition, while quite simple, has some limitations, especially for large-scale simulations. Given a cubic mesh of size N^3 , one obvious constraint is that the maximum number of processors N_{proc} that can be used in a 1D decomposition is N as each slab has to contain at least one plane³¹. For a cubic mesh with $1000^3 = 10^9$ points, the constraint is $N_{proc} \leq 1000$. This is a serious limitation as most supercomputers today have more than 10,000 cores and some have more than 100,000. Large applications are also likely to hit the memory limit when each processor handles too much workload.

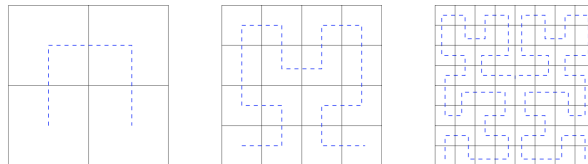
Load balancing

In computing, load balancing distributes workloads across multiple computing resources, such as computers, a computer cluster, network links, central processing units or disk drives. Load balancing aims to optimize resource use, maximize throughput, minimize response time, and avoid overload of any single resource.

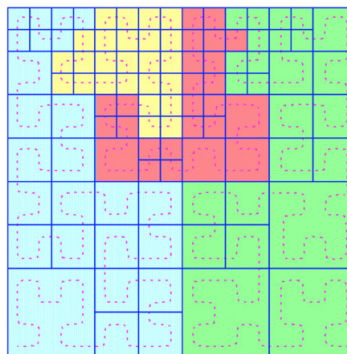
More concretely, the use of many processors in the case where the mesh is unstructured and highly irregular raises the question of how the partitioning is done so that each processor gets a similar workload, thereby minimising latency and optimising cpu time.

This is a difficult problem which is often addressed by means of graph partitioners, and which is made more complicated by the use of AMR.

A typical strategy goes as follows. One can prove that **space filling curves** can be generated for any regular subdivision of a square (see dashed line in the following figure).



It now the refined grid is the one shown on the following figure, one can use the curve to 'write' all elements in a line and then cut this curve in N_{proc} chunks (in this case 4, hence four colours).

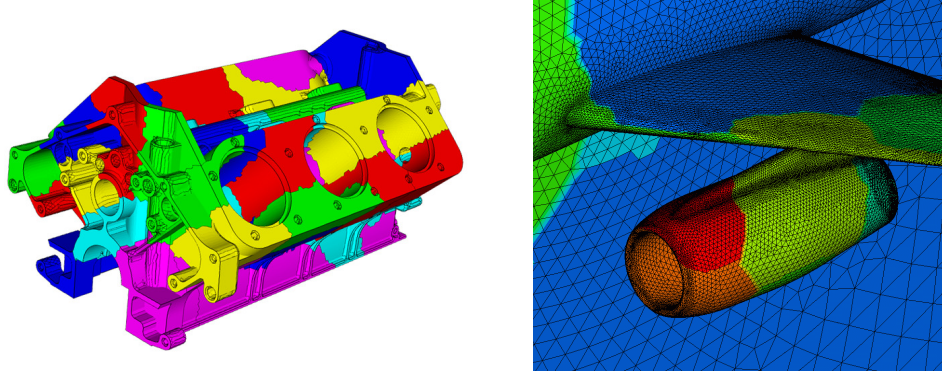


³¹This domain decomposition approach is the one carried out in the FANTOM code [1258] and this limitation was often encountered when using grids such as 160x160x23, thereby limiting the number of processors to about 100 processors.

The mesh counts in total 100 elements and each domain counts 25 elements. This decomposition looks very appealing but another aspect must be taken in account: communication.

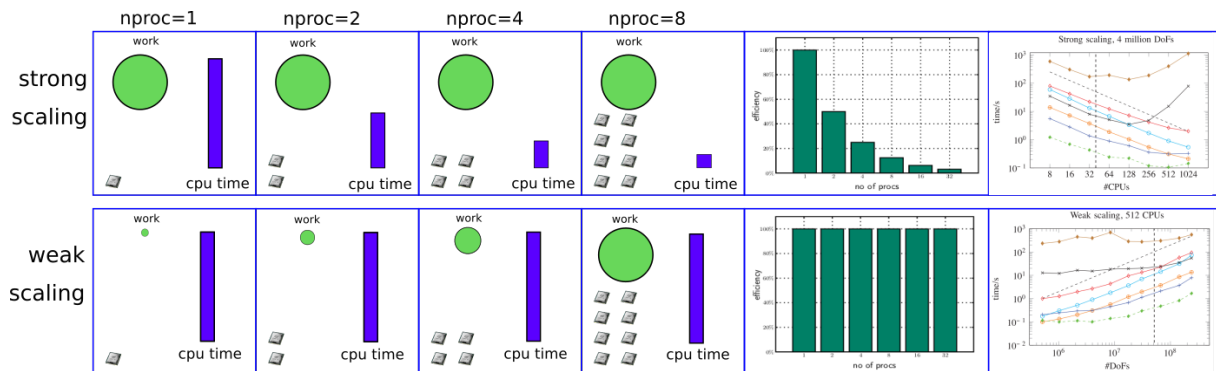
Ideally, one wishes to carry out the domain decomposition in a simple-enough manner so that it does not require too much work, and also in such a way that the surface between between all the domain is minimised, since it is related to the amount of communication across processors.

The following figures showcase examples of domain decomposition projected onto a complex 3D geometry.



Right: Domain decomposition for parallel processing of a wing-body-nacelle-pylon configuration. Each colour corresponds to a different processor.

Strong scaling vs weak scaling



9.21 Corner flow

cornerflow.tex

The mantle wedge comprised between the downgoing slab and the overriding plate has been extensively studied since very important geodynamical processes take place in it or right above it (slab dehydration and water transport, melting, over-riding plate deformation, vulcanism, ...).

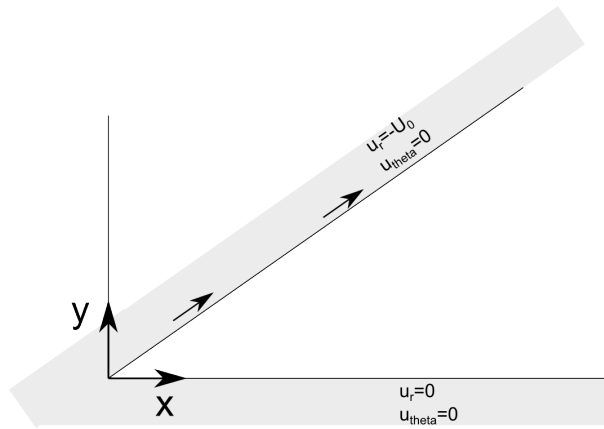
To first approximation one can approach the problem and simplify it greatly by assuming that both plates kinematic behaviour are independent of what happens in the wedge, that the wedge geometry does not change over time, that the problem is essentially 2D, and that the mantle extends very far away from the actual wedge (plates are infinite).

Under such assumptions, it is possible to derive an analytical solution for incompressible Stokes flow in the wedge as documented at p. 224 in Batchelor [52].

Literature: [1277]

FIND refs. check new version of Vol7 theoretical geophys

A corner flow setup is shown hereunder:



The solution to this problem is arrived at by means of the stream function Φ , defined as $u = -\partial\Phi/\partial y$ and $v = \partial\Phi/\partial x$, so that we automatically have $\vec{\nabla} \cdot \vec{v} = 0$. As shown in Section ??, the stream function Φ is then the solution to the biharmonic equation

$$\vec{\nabla}^2 \vec{\nabla}^2 \Phi = \vec{\nabla}^4 \Phi = 0$$

Considering the geometry of the problem has plates of infinite extent with constant relative velocity, the solution for velocity everywhere is expected to be independent of r . This means the equation is separable and we will use a solution of the form

$$\Phi(r, \theta) = R(r)f(\theta)$$

However, given the infinite extent of the domain, the velocity is expected to be independent of r , so we postulate $R(r) = r$ (look at the relationship between velocity components and stream function), or:

$$\Phi(r, \theta) = rf(\theta)$$

and we then have to solve

$$\Delta \left(\frac{1}{r} (f + f'') \right) = \frac{1}{r^3} (f + 2f'' + f''') = 0.$$

The solution of this equation for f is:

$$\begin{aligned} f(\theta) &= A \sin \theta + B \cos \theta + C\theta \sin \theta + D\theta \cos \theta \\ f'(\theta) &= A \cos \theta - B \sin \theta + C(\sin \theta + \theta \cos \theta) + D(\cos \theta - \theta \sin \theta) \end{aligned}$$

with

$$\begin{aligned}\mathbf{v}_r &= \frac{1}{r} \frac{\partial \Phi}{\partial \theta} = f'(\theta) \\ \mathbf{v}_\theta &= -\frac{\partial \Phi}{\partial r} = -f(\theta)\end{aligned}$$

A , B , C and D are four constants to be determined by means of the boundary conditions which are as follows:

$$\begin{aligned}\mathbf{v}_r(\theta = 0) &= 0 \\ \mathbf{v}_\theta(\theta = 0) &= 0 \\ \mathbf{v}_r(\theta = \theta_0) &= -U_0 \\ \mathbf{v}_\theta(\theta = \theta_0) &= 0\end{aligned}$$

or,

$$f'(0) = A + D = 0 \quad (9.85)$$

$$f(0) = B = 0 \quad (9.86)$$

$$f'(\theta_0) = -U_0 \quad (9.87)$$

$$f(\theta_0) = 0 \quad (9.88)$$

From the second equation it is trivial to see that $B = 0$, so that:

$$f(\theta) = A \sin \theta + C \theta \sin \theta + D \theta \cos \theta$$

$$f'(\theta) = A \cos \theta + C(\sin \theta + \theta \cos \theta) + D(\cos \theta - \theta \sin \theta)$$

From the first one we obtain $D = -A$ so that

$$f(\theta) = A(\sin \theta - \theta \cos \theta) + C \theta \sin \theta$$

$$f'(\theta) = A(\theta \sin \theta) + C(\sin \theta + \theta \cos \theta)$$

The last two boundary conditions yield:

$$0 = A(\sin \theta_0 - \theta_0 \cos \theta_0) + C \theta_0 \sin \theta_0$$

$$-U_0 = A(\theta_0 \sin \theta_0) + C(\sin \theta_0 + \theta_0 \cos \theta_0)$$

or,

$$A = -U_0 \frac{\theta_0 \sin \theta_0}{\theta_0^2 - \sin^2 \theta_0} \quad C = U_0 \frac{\sin \theta_0 - \theta_0 \cos \theta_0}{\theta_0^2 - \sin^2 \theta_0}$$

Finally:

$$(A, B, C, D) = (-\theta_0 \sin \theta_0, 0, \sin \theta_0 - \theta_0 \cos \theta_0, \theta_0 \sin \theta_0) \frac{U_0}{\theta_0^2 - \sin^2 \theta_0}$$


We have

$$\mathbf{e}_r = \cos \theta \mathbf{e}_x + \sin \theta \mathbf{e}_y \quad (9.89)$$

$$\mathbf{e}_\theta = -\sin \theta \mathbf{e}_x + \cos \theta \mathbf{e}_y \quad (9.90)$$

so that the velocity field can be expressed in Cartesian coordinates:

$$\begin{aligned}\mathbf{v} &= \mathbf{v}_r \mathbf{e}_r + \mathbf{v}_\theta \mathbf{e}_\theta \\ &= \mathbf{v}_r (\cos \theta \mathbf{e}_x + \sin \theta \mathbf{e}_y) + \mathbf{v}_\theta (-\sin \theta \mathbf{e}_x + \cos \theta \mathbf{e}_y) \\ &= (\mathbf{v}_r \cos \theta - \mathbf{v}_\theta \sin \theta) \mathbf{e}_x + (\mathbf{v}_r \sin \theta + \mathbf{v}_\theta \cos \theta) \mathbf{e}_y\end{aligned} \quad (9.91)$$

 **Relevant Literature:** Ribe [1063] present a simple model for the mantle flow induced by back arc spreading behind a subduction zone.

9.22 Surface processes

surfaceprocesses.tex

In 1D - simple nonlinear diffusion a la Burov & Cloetingh (1997)

This approach comes from Burov and Cloetingh [183] (1997). The tectonic-scale transport equations describe long term changes in topography $h(x, y, t)$ as a result of simultaneous short- and long-range mass transport processes [59, 721].

The short-range surface processes are represented by cumulative effects of hillslope processes (soil creep, rainsplash, slides) that remove material from uplifted areas down to the valleys. It is then assumed that the horizontal material flux \vec{q}_s is related to local slope $\vec{\nabla}h$ by $\vec{q}_s = -K_s \vec{\nabla}h$ where K_s is the effective diffusivity. Assumption of conservation of mass volume leads to the linear diffusion equation for erosion:

$$\frac{\partial h}{\partial t} = K_s \Delta h$$

This equation can be solved with constant-elevation (fixed h value) boundary conditions simulating local base levels of erosion.

Note that in practice the coefficient K_s might depend on slope and curvature, i.e.

$$\frac{\partial h}{\partial t} = K_s(x, y, h, \vec{\nabla}h) \Delta h$$

Following [474], Burov & Cloetingh use an empirical non linear expression $K_s = k_s(x)(\vec{\nabla}h)^n$.

In 1D - not so simple, a la Andr  s-Martinez *et al.* (2019)

This approach comes from Andr  s-Mart  nez, P  rez-Gussiny  , Armitage, and Morgan [25] (2019). The change in surface elevation rate due to surface processes is equal to the divergence of the sediment flux (assuming there is no density difference between the bedrock and sediment and ignoring the effects of compaction):

$$\frac{\partial h}{\partial t} = -\frac{\partial q_s}{\partial x}$$

where h is the topography, t is the time, q_s represents the sediment flux, and x is the horizontal coordinate.

The next step consists in a formulation for the sediment flux. Still following [25], in the subaerial environment, it is possible to define the sediment transport flux q_s in terms of the water flux q_w as

$$q_s = -(K + cq_w^n) \frac{\partial h}{\partial x}$$

where K is the slope diffusivity, c is the transport coefficient, and $n \geq 1$ is the power law that defines the type of relationship between the sediment transport and the water flux (Simpson & Schlunegger, 2003; Smith & Bretherton, 1972).

get these papers

This model accounts for hillslope diffusion processes where the topography will tend to a dispersive diffusion (Culling, 1960) and fluvial transport processes that result in concentrative diffusion due to water run off (Graf, 1984). For a simple parameterization we choose a linear relationship between sediment transport and water flux ($n = 1$).

The water flux can be related to the water discharge/effective rainfall α as

$$\frac{\partial}{\partial x}(\vec{n}q_w) = -\alpha$$

9.23 Geometric multigrid

The following is mostly borrowed from the Wikipedia page on multigrid methods³².

There are many types of (geometric) multigrid algorithms, but the common features are that a hierarchy of grids is considered. The important steps are:

- *Smoothing*: reducing high frequency errors, for example using a few iterations of the Gauss-Seidel method.
- *Residual Computation*: computing residual error after the smoothing operation(s).
- *Restriction*: downsampling the residual error to a coarser grid.
- *Interpolation or prolongation*: interpolating a correction computed on a coarser grid into a finer grid.
- *Correction*: Adding prolonged coarser grid solution onto the finer grid.

There are many choices of multigrid methods with varying trade-offs between speed of solving a single iteration and the rate of convergence with said iteration. The 3 main types are V-Cycle, F-Cycle, and W-Cycle.

Any geometric multigrid cycle iteration is performed on a hierarchy of grids and hence it can be coded using recursion. Since the function calls itself with smaller sized (coarser) parameters, the coarsest grid is where the recursion stops.

Note that the ratio of the number of nodes between two consecutive levels has to be constant between all the levels. Often powers of 2 are used (especially if the grids are based on quad/octrees) but it is not a requirement.

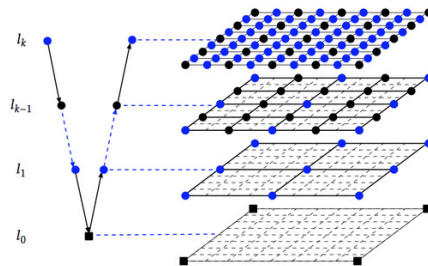


Image from <http://web.utk.edu/~wfeng1/research.html>

What follows is a pseudo-code example of a recursive V-Cycle Multigrid for solving the Poisson equation ($\nabla^2 \phi = f$) on a uniform grid of spacing h :

```
function phi = V_Cycle(phi,f,h)
% Pre-Smoothing
phi = smoothing(phi,f,h);
% Compute Residual Errors
r = residual(phi,f,h);
% Restriction
rhs = restriction(r);
eps = zeros(size(rhs));
% stop recursion at smallest grid size
if smallest_grid_size_is_achieved
```

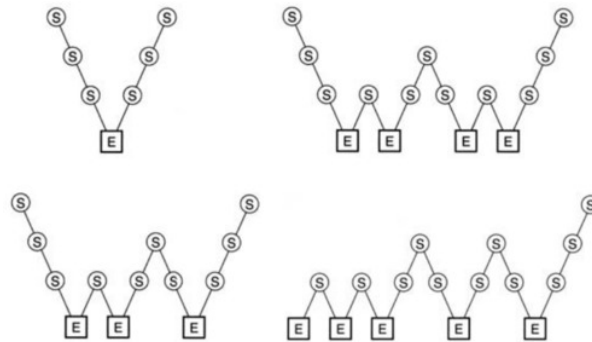
³²https://en.wikipedia.org/wiki/Multigrid_method


```

    eps = smoothing(eps,rhs,2*h);
else
    eps = V_Cycle(eps,rhs,2*h);
end
% Prolongation and Correction
phi = phi + prolongation(eps);
% Post-Smoothing
phi = smoothing(phi,f,h);
end


```

A multigrid method with an intentionally reduced tolerance can be used as an efficient preconditioner for an external iterative solver. The solution may still be obtained in $\mathcal{O}(N)$ time as well as in the case where the multigrid method is used as a solver. Multigrid preconditioning is used in practice even for linear systems, typically with one cycle per iteration.




Taken from [626]: Different types of multigrid cycle with four grid levels: (top left) V-cycle, (top right) W-cycle, (bottom left) F-cycle and (bottom right) full multigrid. ‘S’ denotes smoothing while ‘E’ denotes exact coarse-grid solution.

Check Kaus BEcker syllabus!

 **Relevant Literature:** [1001, 792, 1284, 626, 455, 845, 804, 1285, 908, 1381, 1351, 1246, 1383, 260] Book [153]

- ACuTEMan: A multigrid-based mantle convection simulation code and its optimization to the Earth Simulator, Kameyama (2005) [667]

9.24 Algebraic multigrid

 Relevant Literature: [930][946]

9.25 Computing depth

computing_depth.tex

In the case of a perfectly rectangular, cylindrical or spherical domain, computing the depth of any given point inside the domain is trivial. However, when the free surface becomes somewhat distorted, the concept of depth needs to be refined. What follows is an attempt at bringing clarity as to how to compute depth in all cases.

The depth $d(\mathbf{r})$ satisfies the equation:

$$\frac{\mathbf{g}}{|\mathbf{g}|} \cdot \nabla d = 1$$

with $d = 0$ at the surface.

This is a form of steady-state advection equation (the time derivative is zero, there is no diffusion, nor any source term).

Given the boundary conditions, one could solve this equation over the whole domain.

Note that in the case of a cartesian box, $\mathbf{g} = -g\mathbf{u}_z$, we need to solve

$$-\frac{\partial}{\partial z}d = 1$$

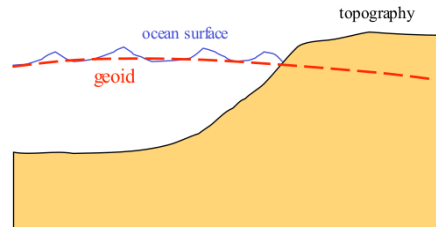
For a flat top surface at $d(z = L_z) = 0$ so that in the end

$$d(z) = L_z - z$$

9.26 The Geoid

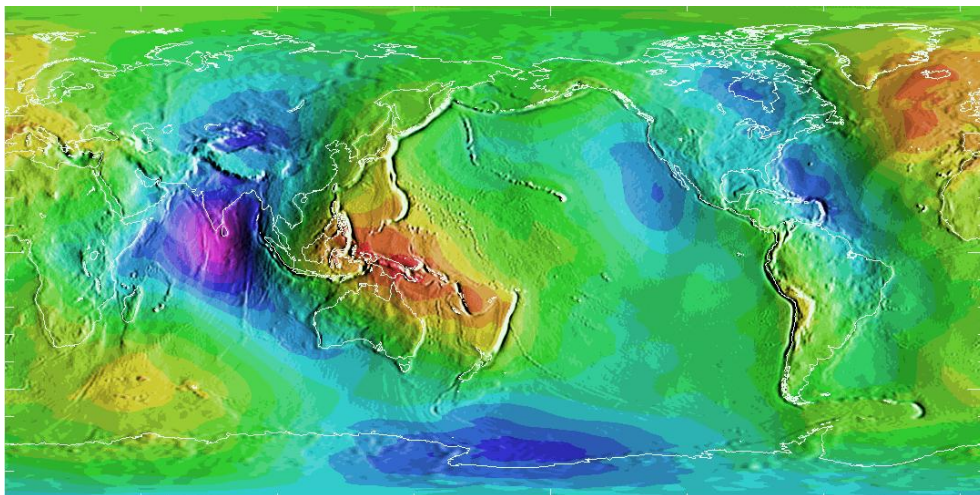
What is the geoid?

There is an infinity of equipotential surfaces of the gravitational potential U . However, there is a particular surface on the Earth that is "easy" to locate: the mean sea level. This is a somewhat arbitrary choice but it makes sense because the oceans are made of water (!): the surface of a fluid in equilibrium must follow an equipotential.



The geoid is usually defined in two ways:

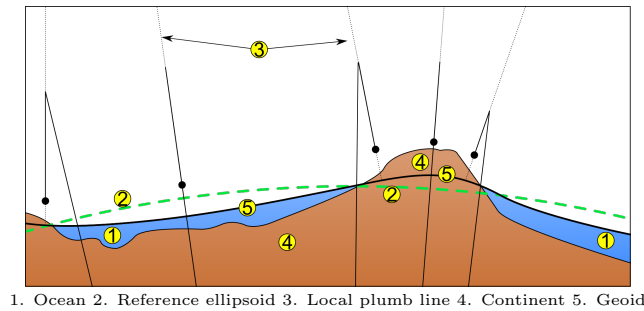
- it is the particular equipotential surface that coincides with the mean sea level (easy to define in the oceans -assuming no currents, waves,... - but harder on land since it is not the topographic surface).
- A gravitational equipotential surface. This means that everywhere at sea level experiences the same value of gravity potential, so there is no tendency for water to flow downhill since all points in the vicinity have the same value of gravity potential, pointed toward the center of the earth.



Data Max value: 85.4 meters, east of New Guinea. Data Min value:-107.0 meters, south of India. This image shows 15'x15' geoid undulations covering the planet Earth from the NIMA/GSFC WGS-84 EGM96 15' Geoid Height File. The undulations refer to the differences from the WGS-84(G873) reference ellipsoid. Map and description from National Geodetic Survey ³³

From Wikipedia: The geoid surface is irregular, unlike the reference ellipsoid (which is a mathematical idealized representation of the physical Earth), but is considerably smoother than Earth's physical surface. Although the physical Earth has excursions of +8,848 m (Mount Everest) and -11,034 m (Marianas Trench), the geoid's deviation from an ellipsoid ranges from +85 m (Iceland) to -106 m (southern India), less than 200 metre total.

³³https://www.usna.edu/Users/oceano/pguth/md_help/geology_course/geoid.htm



Taken from <https://en.wikipedia.org/wiki/Geoid>

the (reference) ellipsoid

First evidence that the Earth is round Erathostene (275-195 B.C.)

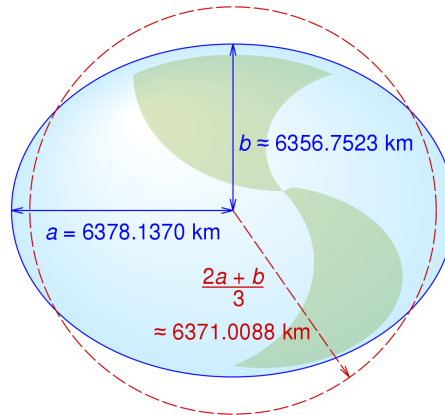
First hypothesis that the Earth is flattened at the poles: Newton

First measurement of the Earth's flattening at the poles: Clairaut (1736) and Bouguer (1743)

The shape of the Earth can be mathematically represented as an ellipsoid defined by:

- Semi-major axis = equatorial radius = a
- Semi-minor axis = polar radius = c
- Flattening (the relationship between equatorial and polar radius): $f = (a - c)/a$
- Eccentricity: $e^2 = 2f - f^2$

Many different reference ellipsoids have been defined and are in use. We define the *reference ellipsoid* = the ellipsoid that best fits the geoid. It is totally arbitrary, but practical. The most common reference ellipsoid is the WGS-84 one³⁴:



Taken from https://en.wikipedia.org/wiki/Reference_ellipsoid

Geoid undulations = differences, in meters, between the geoid reference ellipsoid (= geoid “height”). To clarify:

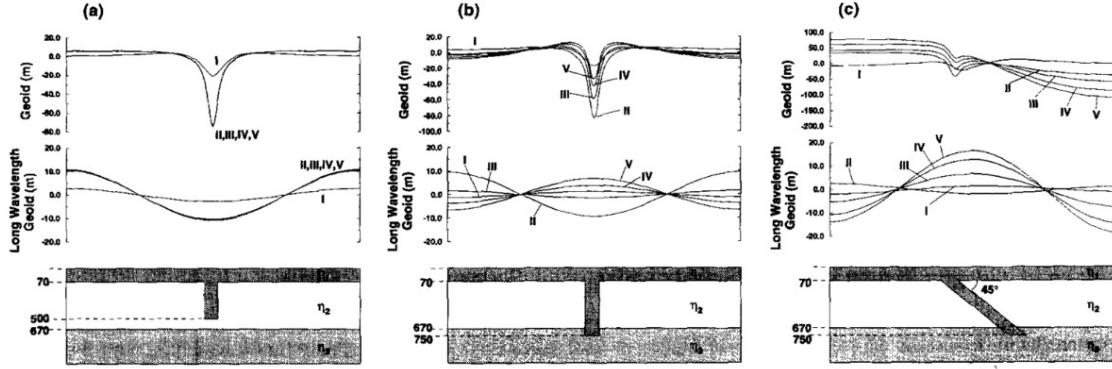
- Geoid = the equipotential surface of the Earth's gravity field that best fits (in a least squares sense) the mean sea level. The gravitational potential is constant on the geoid (by definition) but the gravitational acceleration is not!
- Reference Ellipsoid = the ellipsoid that best fits the geoid
- Geoid = the (actual) figure of the Earth
- Ellipsoid = the (theoretical) shape of the Earth

³⁴<https://confluence.qps.nl/qinsy/latest/en/world-geodetic-system-1984-wgs84-182618391.html>

How to compute it?

From Liu and Zhong [803] (2016): “The geoid is computed by ϕ/g , where ϕ is the surface gravitational potential anomaly and can be solved from the Poisson equation, $\nabla^2\phi = -4\pi\mathcal{G}\delta\rho$ where \mathcal{G} is the gravitational constant, and $\delta\rho$ includes both density variations in the mantle [...] and those associated with dynamic topographies at the surface and CMB. Dynamic topographies are determined from solving [the Stokes] equations under free-slip boundary conditions at the surface and CMB.”

Interesting modelling



Idealized 2D slab

calculations for each viscosity model: geoid and geoid filtered to pass only the longest wavelengths (~ 4000 km). (a) Cold slab extends to 500 km depth in the upper mantle, (b) Slab extends to 750 km so that it is partly supported by the high viscosity lower mantle at 670 km. (c) Slab tilted at 45° to the vertical extending to the top of the lower mantle. Taken from [900]

9.27 Mixing and stirring, the Lyapunov time/exponent

lyapunov.tex

Many approaches are taken in the literature when it comes to studying mixing/stirring in fluids, and in our case the mantle.

For example, Samuel and Tosi [1105] (2012) measure the convective stirring efficiency using two Lagrangian methods: the first determines the mixing time associated with different wavelengths of heterogeneity following the approach of Ferrachat and Ricard [392] (2001). The second determines the value of the maximum Finite Time Lyapunov Exponents (FTLE) as described in Farnetani and Samuel [386] (2003), and measures the rate at which heterogeneities are stretched by mantle motions.

In Tackley and Xie [1230] (2002) we read:

“ Two-dimensional simulations of simple mantle convection (Christensen [244], 1989; Kellogg and Turcotte [692], 1990; Schmalzl and Hansen [1119]) suggest that for whole-mantle convection the mantle should be homogenized in a time-scale of less than 1 Gyr, although unmixed islands may remain. Non-Newtonian rheology, which is thought to be important in the upper mantle [673], may somewhat inhibit mixing (Ten, Podladchikov, Yuen, Larsen, and Malevsky [1244], 1998). The rate at which blobs of anomalous (e.g. primitive) material are stretched and assimilated into the flow depends on their relative viscosity; very viscous blobs can survive intact for many mantle overturns (Manga [834], 1996).

In three-dimensional geometry with only poloidal flow, mixing may be significantly less efficient than in two dimensions (Schmalzl, Houseman, and Hansen [1121], 1995; Schmalzl, Houseman, and Hansen [1120], 1996) but, if the toroidal flow associated with plate motions is included (Gable, O’connell, and Travis [433], 1991), mixing can instead be more efficient (Ferrachat and Ricard [393] (1998)).

High viscosity in the deep mantle is not sufficient to maintain different reservoirs over geological time-scales (Ferrachat and Ricard [392], 2001; van Keken & Ballentine 1998, 1999), in contrast to predictions from earlier calculations at lower convective vigour (Gurnis and Davies [515], 1986). Part of the reason for this apparent discrepancy is that the latter study used a kinematically driven flow rather than buoyancy-driven flow. Since a viscosity jump does not affect densities, thermal buoyancy-driven flow has no problem crossing it, so substantial mass exchange occurs between upper and lower mantles. Buoyancy is thus necessary to maintain separate reservoirs. ”

In Gottschaldt, Walzer, Hendel, Stegman, Baumgardner, and Mühlhaus [475] (2006) we find

“ Heterogeneities in a convecting fluid are deformed by stirring and finally erased by diffusive mixing. Chemical diffusion in mantle rock acts on the scale of centimetres over the lifetime of the Earth, but our models resolve km-scales. Therefore this paper deals only with convective stirring. Some nice studies about mantle stirring and mixing in 2-d have been done (e.g. [244, 515, 694, 692, 865, 958, 957, 1242, 1243, 1244]). Unfortunately, results from studies in 2-d can be extrapolated to 3-d only in a limited manner. Tracers in 3-d poloidal stationary convection move on 2-d toruslike surfaces. Stirring is constrained to these surfaces. Cross-cell stirring becomes possible in time-dependent flows, but may not be very efficient [1120]. Large-scale stirring is enhanced by a toroidal component, but convectively isolated islands of laminar stirring may remain [393]. Convection in the Earth is time-dependent and the surface planform shows a strong toroidal component today. Since there is currently no general understanding about the stirring behaviour of the mantle, a case-by-case study of different models is necessary. ”

9.27.1 The Lyapunov exponent

Simply put, the Lyapunov time is the characteristic timescale on which a dynamical system is chaotic. It is defined as the inverse of a system's largest Lyapunov exponent.

The Lyapunov time mirrors the limits of the predictability of the system. By convention, it is defined as the time for the distance between nearby trajectories of the system to increase by a factor of e . However, measures in terms of 2-foldings and 10-foldings are sometimes found, since they correspond to the loss of one bit of information or one digit of precision respectively.

The Lyapunov exponent or Lyapunov characteristic exponent of a dynamical system is a quantity that characterizes the rate of separation of infinitesimally close trajectories. Quantitatively, two trajectories in phase space with initial separation $\delta\mathbf{Z}_0$ diverge (provided that the divergence can be treated within the linearized approximation) at a rate given by

$$|\delta\mathbf{Z}(t)| \approx e^{\lambda t} |\delta\mathbf{Z}_0|$$

where λ is the Lyapunov exponent.

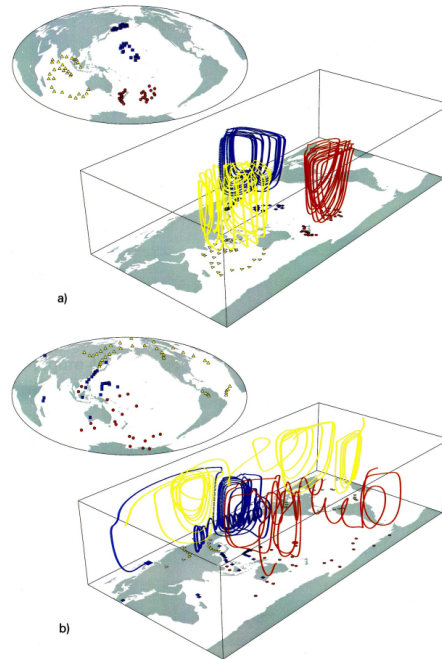
Measuring the Lyapunov exponent or time (or related quantities) is relevant in the context of mantle stirring. On the one hand it is argued that the mantle is convecting and very efficient at mixing resulting in a somewhat homogenous composition. On the other hand, there are modeling studies that suggest that whole-mantle convection can preserve heterogeneity in the presence of well-mixed mantle.

Mixing takes place by the repeated stretching and folding of interfaces. A measure of the mixing efficiency is the time evolution of the area of the mixing surface. Maximum efficiency of mixing is reached with turbulent mixing behavior where One can formally show whether mixing is laminar or turbulent by evaluating the Lyapunov exponents σ . These are of the form:

$$\sigma = \lim_{t \rightarrow \infty} \lim_{X \rightarrow 0} \left[\frac{1}{t} \ln \left(\frac{X(t)}{X(t=0)} \right) \right]$$

where $X(t)$ is the length of this segment at time t . Non-zero Lyapunov exponents indicate that stretching is exponential and the larger the exponent, the more efficient mixing is. However, the limits in the above equation are difficult to evaluate and the interpretation of the 'finite-time' Lyapunov exponent, where both limits are truncated, is difficult to formalize.

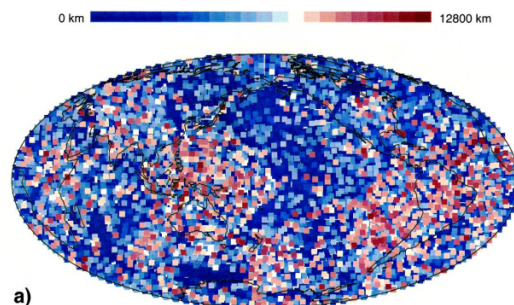
In Keken and Zhong [688] (1999) the authors use a steady state velocity pattern obtained for a model of present-day mantle convection. The velocity model is based on the solution of the Stokes equations in a 3D spherical model with variable rheology. To study mixing, they release particles in the velocity model and follow these by numerical integration.



a) The three particles in this plot were selected for their relatively regular pattern. b) Three other particles that traverse a large portion of the model. These particles feel the strong toroidal motion and their paths form corkscrew-like patterns. They indicate that certain parts of the model can exhibit strong mixing.

Taken from Keken and Zhong [688] (1999).

Rather than calculating the exponents explicitly, the authors use an approximation to the finite-time, finite-length Lyapunov exponent by evaluating the distance between two points that are closely spaced at time $t = 0$. For this they compute the advection of a large number of 10 km long line segments that were originally at 1500 km depth. The length of these segments is approximated by the distance between the endpoints and the results are summarized in the following figure:



Length of the line segment after 4 billion years. Approximately 14,000 line segments were released with regular spacing at 1500 km depth. The length of the segment is indicated by the colored symbols that are plotted at the initial position. The results indicate that there is a strong diversity in mixing behavior. In some regions (north Pacific, parts under the Indian/Australian plate) stretching is very limited, indicating laminar and consequently inefficient mixing. Regions that are under strong toroidal surface motion (western Pacific, Nazca and South America) show very efficient stretching of up to the maximum length of the diameter of the Earth. Taken from Keken and Zhong [688] (1999).

In Bello, Coltice, Rolf, and Tackley [70] (2014) the authors estimate for the first time the limit of predictability of Earth's mantle convection. Following the twin experiment method, we compute the Lyapunov time (i.e., e-folding time) for state of the art 3-D spherical convection models, varying rheology, and Rayleigh number.

Reconstruction of mantle convection and surface tectonics with (ensemble) Kalman filter: Bocher, Coltice, Fournier, and Tackley [99] (2016), Bocher, Fournier, and Coltice [100] (2018).

Investigating the initial condition of mantle models using data assimilation. PhD thesis. Price [1016] (2016).

Relevant Literature

Pierrehumbert [998] (1991) Colli, Bunge, and Schuberth [271] (2015),

9.27.2 configurational 'Shannon' entropy

Goltz and Böse [472] (2002), Camesasca, Kaufman, and Manas-Zloczower [204] (2006), Naliboff and Kellogg [926] (2007), van der Wiel et al. (2024).

9.27.3 Literature to sort out

Richter, Daly, and Nataf [1072] (1982), **hayk82** (1982),
Gurnis [513] (1986),
Christensen [244] (1989),
Schmalzl and Hansen [1119] (1994),
Keken and Zhong [688] (1999),
Farnetani, Legras, and Tackley [385] (2002),
Farnetani and Samuel [386] (2003),
Coltice and Schmalzl [273] (2006),
Huang and Davies [600] (2007),
Manga [833] (2010),
Samuel, Aleksandrov, and Deo [1102] (2011),
Wiel, Hinsbergen, Thieulot, and Spakman [1355] (2024) Thomas, Samuel, Farnetani, Aubert, and Chauvel [1266] (2024)

9.28 Phase transitions

The topic of phase transitions and their implementation in computational geodynamics is a very vast topic. It requires input from thermodynamics, geochemistry and petrology, and also requires dedicated algorithms which are quite complex.

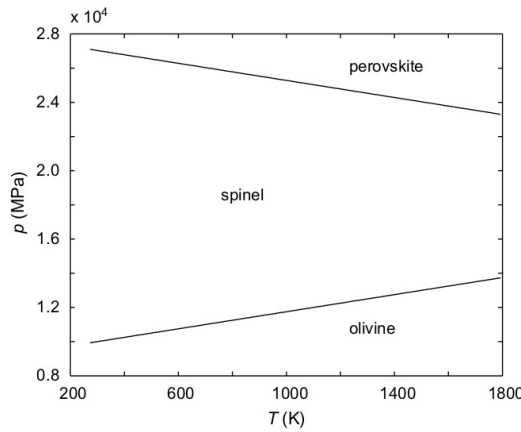
Let us start with simple examples from the literature:

- Zlotnik *et al.* (2007) [1440]. The equations that is used in this work are the standard incompressible Stokes equations by the authors chose to represent the density as a function of of temperature and pressure by the following expression:

$$\rho(T, p) = \rho_0[1 - \alpha(T - T_0)][1 + \beta(p - p_0)]$$

where α and β are, respectively, the thermal expansion and compressibility coefficients, and T_0 and p_0 are reference values at surface.

The authors then proceed to divide the phase diagram into three regions corresponding to three minerals: olivine, spinel-structured olivine, and perovskite:



Phase diagram indicating stable mineral phases in the temperature-pressure plane. The phase diagram is divided into three regions corresponding to three distinct minerals: olivine, spinel and perovskite. Taken from [1440].

They state that two major mineralogical phase transitions occur, one at 410 km depth and other at 660 km depth (other deeper transitions run outside the domain under study because their domain is 1000km deep). The density increases discontinuously across these phase transitions. In order to take into account the effect of these discontinuities, the density ρ_0 above is taken as a reference density plus an increment $\Delta\rho$:

$$\rho_0 = \rho_{olivine} + \Delta\rho$$

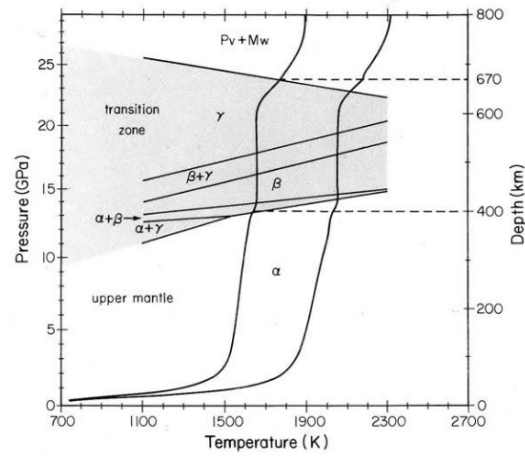
where

$$\Delta\rho = \begin{cases} 0 & \text{if } (T, p) \text{ is in the olivine region} \\ \Delta\rho_{es} & \text{if } (T, p) \text{ is in the spinel region} \\ \Delta\rho_{per} & \text{if } (T, p) \text{ is in the perovskite region} \end{cases}$$

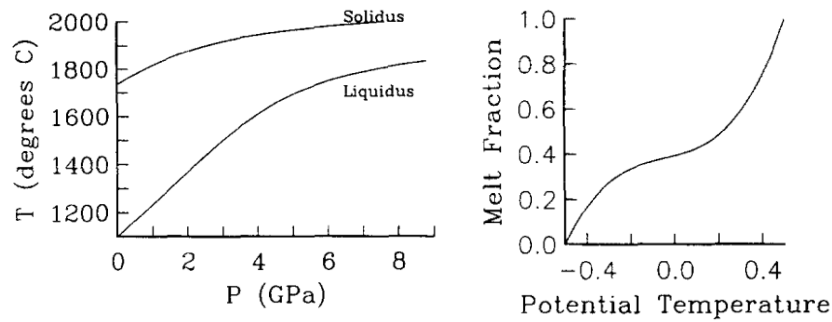
The authors unfortunately fail to report how the phase transitions affect the viscosity.

The obvious problem with this otherwise simple approach is that density varies in the domain but is not accompanied by a volume change so that it violates mass conservation.

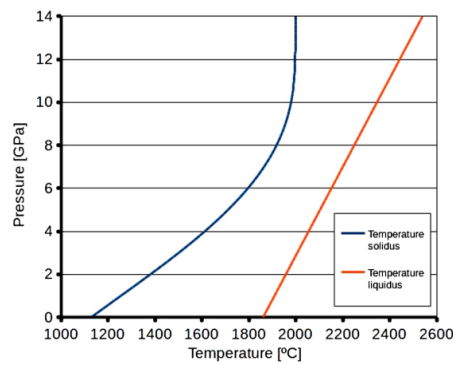
- the following phase diagram is taken from Peltier *et al.* (1997) [988].



Phase boundary for the $\alpha \rightarrow \beta \rightarrow \gamma$ transitions of Olivine



Left: Solidus and liquidus temperatures; Right: melt fraction as a function temperature for dry peridotite. Taken from King & Anderson (1995) [702], Both after McKenzie & Bickle [854].



Left: Solidus and liquidus curves. Taken from Lavecchia *et al.* (2017) [753].

📖 Relevant Literature: [1143]

9.29 Implementation of an elasto-viscous rheology

evrheo.tex

A viscoelastic material can behave both elastically and viscously. Its response to an applied stress is dependent on the material properties and can be found using the Maxwell time (t_M) of said material. This material constant is defined as the ratio of the material viscosity and shear modulus

$$t_M = \frac{\eta}{\mu}$$

where η is the viscosity and μ the shear modulus.

The total deviatoric strainrate tensor can be decomposed into an elastic component and a viscous component:

$$\dot{\epsilon}^d(\vec{v}) = \dot{\epsilon}_e^d(\vec{v}) + \dot{\epsilon}_v^d(\vec{v}) = \frac{\dot{\vec{\tau}}}{2\mu} + \frac{\vec{\tau}}{2\eta}$$

[From wikipedia] In continuum mechanics, objective stress rates are time derivatives of stress that do not depend on the frame of reference. Many constitutive equations are designed in the form of a relation between a stress-rate and a strain-rate (or the rate of deformation tensor). i The mechanical response of a material should not depend on the frame of reference. In other words, material constitutive equations should be frame indifferent (objective). If the stress and strain measures are material quantities then objectivity is automatically satisfied. However, if the quantities are spatial, then the objectivity of the stress-rate is not guaranteed even if the strain-rate is objective.

There are numerous objective stress rates in continuum mechanics - all of which can be shown to be special forms of Lie derivatives. Some of the widely used objective stress rates are [586]: a) the Truesdell rate of the Cauchy stress tensor, b) the Green-Naghdi rate of the Cauchy stress, and c) the Jaumann rate of the Cauchy stress.

The Jaumann rate of the Cauchy stress is a further specialization of the Lie derivative (Truesdell rate). This rate has the form

$$\dot{\vec{\tau}}^{t+\delta t} = \frac{D\vec{\tau}}{Dt} - (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t)) = \frac{\vec{\tau}^{t+\delta t} - \vec{\tau}^t}{\delta t} - (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t))$$

where D/Dt is the material derivative and $\dot{\omega}$ is the rotation rate -also called spin tensor- which is anti-symmetric and has zero trace - see Section 2.4.3:

$$\dot{\omega}(\vec{v}) = \frac{1}{2} \left(\vec{\nabla} \vec{v} - (\vec{\nabla} \vec{v})^T \right)$$

In the case of a Lagrangian description, we have [1328]

$$\dot{\vec{\tau}}^{t+\delta t} = \frac{\vec{\tau}^{t+\delta t} - \vec{\tau}^t}{\delta t} - (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t))$$

so that

$$\dot{\epsilon}^d(\vec{v}^{t+\delta t}) = \frac{\dot{\vec{\tau}}^{t+\delta t}}{2\mu} + \frac{\vec{\tau}^{t+\delta t}}{2\eta} = \frac{1}{2\mu} \left[\frac{\vec{\tau}^{t+\delta t} - \vec{\tau}^t}{\delta t} - (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t)) \right] + \frac{\vec{\tau}^{t+\delta t}}{2\eta}$$

Let us multiply this by $2\mu\delta t$ and transform the equations until a satisfying formulation is found:

$$2\mu\delta t \dot{\epsilon}^d(\vec{v}^{t+\delta t}) = \vec{\tau}^{t+\delta t} - \vec{\tau}^t - \delta t (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t)) + \frac{\mu\delta t}{\eta} \vec{\tau}^{t+\delta t}$$

$$2\mu\delta t \dot{\epsilon}^d(\vec{v}^{t+\delta t}) = \left(1 + \frac{\mu\delta t}{\eta} \right) \vec{\tau}^{t+\delta t} - \vec{\tau}^t - \delta t (\dot{\omega}(\vec{v}^t) \cdot \vec{\tau}^t - \vec{\tau}^t \cdot \dot{\omega}(\vec{v}^t))$$

$$\begin{aligned} \left(1 + \frac{\mu\delta t}{\eta}\right) \boldsymbol{\tau}^{t+\delta t} &= 2\mu\delta t \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + \boldsymbol{\tau}^t + \delta t(\dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t) \cdot \boldsymbol{\tau}^t - \boldsymbol{\tau}^t \cdot \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t)) \\ \boldsymbol{\tau}^{t+\delta t} &= 2\frac{\mu\delta t}{\left(1 + \frac{\mu\delta t}{\eta}\right)} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + \frac{1}{\left(1 + \frac{\mu\delta t}{\eta}\right)} \boldsymbol{\tau}^t + \frac{\delta t}{\left(1 + \frac{\mu\delta t}{\eta}\right)} (\dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t) \cdot \boldsymbol{\tau}^t - \boldsymbol{\tau}^t \cdot \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t)) \end{aligned}$$

We define:

$$\boxed{\eta_{eff} = \frac{\mu\delta t}{\left(1 + \frac{\mu\delta t}{\eta}\right)} = \frac{\eta\delta t}{\delta t + \eta/\mu} = \frac{\eta}{1 + t_M/\delta t}} \quad (9.92)$$

$$\boxed{Z = \frac{\eta_{eff}}{\mu\delta t} = \frac{\eta}{\mu\delta t + \eta}}$$

$$\boxed{\mathbf{J}^t = \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t) \cdot \boldsymbol{\tau}^t - \boldsymbol{\tau}^t \cdot \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t)}$$

so that we can write

$$\boldsymbol{\tau}^{t+\delta t} = 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + Z\boldsymbol{\tau}^t + Z\delta t \mathbf{J}^t$$

or,

$$\boldsymbol{\tau}^{t+\delta t} = 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + \underline{\boldsymbol{\tau}}^t \quad \text{with} \quad \underline{\boldsymbol{\tau}}^t = Z\boldsymbol{\tau}^t + Z\delta t \mathbf{J}^t$$

The total stress tensor is then

$$\boxed{\boldsymbol{\sigma}^{t+\delta t} = -p^{t+\delta t} \mathbf{1} + \boldsymbol{\tau}^{t+\delta t} = -p^{t+\delta t} \mathbf{1} + 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + Z\boldsymbol{\tau}^t + Z\delta t \mathbf{J}^t} \quad (9.93)$$

Remark. When $\mu \rightarrow \infty$ we have $\eta_{eff} \rightarrow \eta$ and $Z \rightarrow 0$ and we recover the Stokes equation for a purely viscous fluid.

Strong form

Let us now turn to the momentum conservation equation:

$$\begin{aligned} &\vec{\nabla} \cdot \boldsymbol{\sigma}^{t+\delta t} + \rho^{t+\delta t} \vec{g} = \vec{0} \\ \Rightarrow &\vec{\nabla} \cdot (-p^{t+\delta t} \mathbf{1} + \boldsymbol{\tau}^{t+\delta t}) + \rho^{t+\delta t} \vec{g} = \vec{0} \\ \Rightarrow &-\vec{\nabla} p^{t+\delta t} + \vec{\nabla} \cdot \boldsymbol{\tau}^{t+\delta t} + \rho^{t+\delta t} \vec{g} = \vec{0} \\ \Rightarrow &-\vec{\nabla} p^{t+\delta t} + \vec{\nabla} \cdot (2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + \underline{\boldsymbol{\tau}}^t) + \rho^{t+\delta t} \vec{g} = \vec{0} \end{aligned}$$

and finally

$$\boxed{-\vec{\nabla} p^{t+\delta t} + \vec{\nabla} \cdot 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) = -\rho^{t+\delta t} \vec{g} - \vec{\nabla} \cdot \underline{\boldsymbol{\tau}}^t}$$

Weak form

The mass conservation equation is still $\vec{\nabla} \cdot \vec{\mathbf{v}} = 0$ so we need not look into it since its weak form is in Section 7.5.

For the N_i^γ 's we can write:

$$\int_{\Omega_e} N_i^\gamma \vec{\nabla} \cdot \boldsymbol{\sigma}^{t+\delta t} d\Omega + \int_{\Omega_e} N_i^\gamma \rho \vec{g} d\Omega = \vec{0} \quad (9.94)$$

We can integrate by parts and drop the surface term³⁵ REVISIT and use Eq. (9.93):

$$\begin{aligned} \int_{\Omega_e} \vec{\nabla} N_i^\gamma \cdot \boldsymbol{\sigma}^{t+\delta t} d\Omega &= \int_{\Omega_e} N_i^\gamma \rho \vec{g} d\Omega \\ \int_{\Omega_e} \vec{\nabla} N_i^\gamma \cdot [-p^{t+\delta t} \mathbf{1} + 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t}) + \underline{\boldsymbol{\tau}}^t] d\Omega &= \int_{\Omega_e} N_i^\gamma \rho \vec{g} d\Omega \end{aligned} \quad (9.95)$$

$$\int_{\Omega_e} \vec{\nabla} N_i^\gamma \cdot [-p^{t+\delta t} \mathbf{1} + 2\eta_{eff} \dot{\boldsymbol{\epsilon}}^d(\vec{\mathbf{v}}^{t+\delta t})] d\Omega = \int_{\Omega_e} N_i^\gamma \rho \vec{g} d\Omega - \int_{\Omega_e} \vec{\nabla} N_i^\gamma \cdot \underline{\boldsymbol{\tau}}^t d\Omega \quad (9.96)$$

We see that the left hand term is virtually identical to the one in Section 7.5, although the viscosity has been replaced with the effective viscosity of Eq. (9.92). The headache will come from the right hand side term $\underline{\boldsymbol{\tau}}^t$, as we will see.

In two dimensions - Cartesian coordinates . The rotation rate tensor is given by:

$$\dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}) = \frac{1}{2} (\vec{\nabla} \vec{\mathbf{v}} - (\vec{\nabla} \vec{\mathbf{v}})^T) = \begin{pmatrix} 0 & \dot{\omega}_{xy} \\ -\dot{\omega}_{xy} & 0 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 0 & \frac{\partial v}{\partial x} - \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} - \frac{\partial u}{\partial y} & 0 \end{pmatrix}$$

so that the tensor \mathbf{J} can be computed explicitly:

$$\begin{aligned} \mathbf{J}^t &= \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t) \cdot \boldsymbol{\tau}^t - \boldsymbol{\tau}^t \cdot \dot{\boldsymbol{\omega}}(\vec{\mathbf{v}}^t) \\ &= \begin{pmatrix} 0 & \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \\ -\dot{\omega}_{xy}(\vec{\mathbf{v}}^t) & 0 \end{pmatrix} \cdot \begin{pmatrix} \tau_{xx}^t & \tau_{xy}^t \\ \tau_{xy}^t & \tau_{yy}^t \end{pmatrix} - \begin{pmatrix} \tau_{xx}^t & \tau_{xy}^t \\ \tau_{xy}^t & \tau_{yy}^t \end{pmatrix} \cdot \begin{pmatrix} 0 & \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \\ -\dot{\omega}_{xy}(\vec{\mathbf{v}}^t) & 0 \end{pmatrix} \\ &= \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \cdot \begin{pmatrix} \tau_{xx}^t & \tau_{xy}^t \\ \tau_{xy}^t & \tau_{yy}^t \end{pmatrix} - \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} \tau_{xx}^t & \tau_{xy}^t \\ \tau_{xy}^t & \tau_{yy}^t \end{pmatrix} \cdot \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \\ &= \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} \tau_{xy}^t & \tau_{yy}^t \\ -\tau_{xx}^t & -\tau_{xy}^t \end{pmatrix} - \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} -\tau_{xy}^t & \tau_{xx}^t \\ -\tau_{yy}^t & \tau_{xy}^t \end{pmatrix} \\ &= \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} 2\tau_{xy}^t & \tau_{yy}^t - \tau_{xx}^t \\ \tau_{yy}^t - \tau_{xx}^t & -2\tau_{xy}^t \end{pmatrix} \end{aligned} \quad (9.97)$$

so that the tensor equation $\underline{\boldsymbol{\tau}} = Z\boldsymbol{\tau} + Z\delta t\mathbf{J}$ can be reformulated as follows in a vector form:

$$\begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} = Z \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} + Z\delta t \begin{pmatrix} J_{xx}^t \\ J_{yy}^t \\ J_{xy}^t \end{pmatrix} = Z \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} + Z\delta t \dot{\omega}_{xy}^t \begin{pmatrix} 2\tau_{xy}^t \\ -2\tau_{xy}^t \\ \tau_{yy}^t - \tau_{xx}^t \end{pmatrix} \quad (9.98)$$

or,

$$\begin{aligned} \begin{pmatrix} \sigma_{xx}^{t+\delta t} \\ \sigma_{yy}^{t+\delta t} \\ \sigma_{xy}^{t+\delta t} \end{pmatrix} &= \begin{pmatrix} -p^{t+\delta t} \\ -p^{t+\delta t} \\ 0 \end{pmatrix} + 2\eta_{eff} \begin{pmatrix} \dot{\epsilon}_{xx}(\vec{\mathbf{v}}^{t+\delta t}) \\ \dot{\epsilon}_{yy}(\vec{\mathbf{v}}^{t+\delta t}) \\ \dot{\epsilon}_{xy}(\vec{\mathbf{v}}^{t+\delta t}) \end{pmatrix} + \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} \\ &= \begin{pmatrix} -p^{t+\delta t} \\ -p^{t+\delta t} \\ 0 \end{pmatrix} + 2\eta_{eff} \begin{pmatrix} \dot{\epsilon}_{xx}(\vec{\mathbf{v}}^{t+\delta t}) \\ \dot{\epsilon}_{yy}(\vec{\mathbf{v}}^{t+\delta t}) \\ \dot{\epsilon}_{xy}(\vec{\mathbf{v}}^{t+\delta t}) \end{pmatrix} + Z \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} + Z\delta t \dot{\omega}_{xy}^t \begin{pmatrix} 2\tau_{xy}^t \\ -2\tau_{xy}^t \\ \tau_{yy}^t - \tau_{xx}^t \end{pmatrix} \end{aligned} \quad (9.99)$$

...

$$\int_{\Omega_e} \mathbf{B}^T \cdot \begin{pmatrix} \sigma_{xx}^{t+\delta t} \\ \sigma_{yy}^{t+\delta t} \\ \sigma_{xy}^{t+\delta t} \end{pmatrix} d\Omega = \int_{\Omega_e} \vec{N}_b d\Omega \quad (9.100)$$

³⁵We will come back to this at a later stage

$$\int_{\Omega_e} \mathbf{B}^T \cdot \left[\begin{pmatrix} -p^{t+\delta t} \\ -p^{t+\delta t} \\ 0 \end{pmatrix} + 2\eta_{eff} \begin{pmatrix} \varepsilon_{xx}^d(\vec{\mathbf{v}}^{t+\delta t}) \\ \varepsilon_{yy}^d(\vec{\mathbf{v}}^{t+\delta t}) \\ \varepsilon_{xy}^d(\vec{\mathbf{v}}^{t+\delta t}) \end{pmatrix} \right] d\Omega = \int_{\Omega_e} \vec{N}_b d\Omega - \int_{\Omega_e} \mathbf{B}^T \cdot \left[Z \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} + Z\delta t \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} 2\tau_{xy}^t \\ -2\tau_{xy}^t \\ \tau_{yy}^t - \tau_{xx}^t \end{pmatrix} \right]$$

As seen in Section 7.5 the left hand side terms yield $\mathbb{K} \cdot \vec{V} + \mathbb{G} \cdot \vec{P}$. The buoyancy term in the rhs is also standard and yields \vec{f} . The discretised momentum equation then writes

$$\mathbb{K} \cdot \vec{V} + \mathbb{G} \cdot \vec{P} = \vec{f} + \vec{f}_{el}$$

and the last rhs term is

$$\boxed{\vec{f}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \left[Z \begin{pmatrix} \tau_{xx}^t \\ \tau_{yy}^t \\ \tau_{xy}^t \end{pmatrix} + Z\delta t \dot{\omega}_{xy}(\vec{\mathbf{v}}^t) \begin{pmatrix} 2\tau_{xy}^t \\ -2\tau_{xy}^t \\ \tau_{yy}^t - \tau_{xx}^t \end{pmatrix} \right] d\Omega}$$

with the matrix \mathbf{B} being given by

$$\mathbf{B} = \begin{pmatrix} \frac{\partial N_1^\vee}{\partial x} & 0 & 0 & \dots & \frac{\partial N_{mv}^\vee}{\partial x} & 0 & 0 \\ 0 & \frac{\partial N_1^\vee}{\partial y} & 0 & \dots & 0 & \frac{\partial N_{mv}^\vee}{\partial y} & 0 \\ \frac{\partial N_1^\vee}{\partial y} & \frac{\partial N_1^\vee}{\partial x} & 0 & \dots & \frac{\partial N_{mv}^\vee}{\partial x} & \frac{\partial N_{mv}^\vee}{\partial x} & 0 \end{pmatrix}$$

In three dimensions - Cartesian coordinates The spin tensor is given by

$$\dot{\omega}(\vec{\mathbf{v}}) = \frac{1}{2} \left(\vec{\nabla} \vec{\mathbf{v}} - (\vec{\nabla} \vec{\mathbf{v}})^T \right) = \begin{pmatrix} 0 & \dot{\omega}_{xy} & \dot{\omega}_{xz} \\ -\dot{\omega}_{xy} & 0 & \dot{\omega}_{yz} \\ -\dot{\omega}_{xz} & -\dot{\omega}_{yz} & 0 \end{pmatrix}$$

so that

$$\begin{aligned} \mathbf{J}^t &= \dot{\omega}(\vec{\mathbf{v}}^t) \cdot \boldsymbol{\tau}^t - \boldsymbol{\tau}^t \cdot \dot{\omega}(\vec{\mathbf{v}}^t) \\ &= \begin{pmatrix} 0 & \dot{\omega}_{xy} & \dot{\omega}_{xz} \\ -\dot{\omega}_{xy} & 0 & \dot{\omega}_{yz} \\ -\dot{\omega}_{xz} & -\dot{\omega}_{yz} & 0 \end{pmatrix} \cdot \begin{pmatrix} \tau_{xx} & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} \end{pmatrix} - \begin{pmatrix} \tau_{xx} & \tau_{xy} & \tau_{xz} \\ \tau_{xy} & \tau_{yy} & \tau_{yz} \\ \tau_{xz} & \tau_{yz} & \tau_{zz} \end{pmatrix} \cdot \begin{pmatrix} 0 & \dot{\omega}_{xy} & \dot{\omega}_{xz} \\ -\dot{\omega}_{xy} & 0 & \dot{\omega}_{yz} \\ -\dot{\omega}_{xz} & -\dot{\omega}_{yz} & 0 \end{pmatrix} \\ &= \end{aligned} \tag{9.101}$$

FINISH!!!

check appendix A of Loes' GR

9.30 Interpolation inside an element

The $n + 1$ Bernstein basis polynomials of degree n on the interval $[0, 1]$ are defined as ³⁶

$$b_{m,n}(x) = \binom{n}{m} x^m (1-x)^{n-m} \quad m = 0, 1, \dots, n$$

The first few Bernstein polynomials are

$$b_{0,0}(x) = 1 \quad (9.102)$$

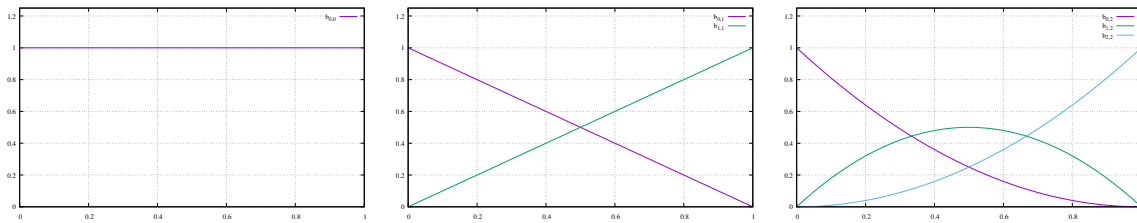
$$b_{0,1}(x) = 1 - x \quad (9.103)$$

$$b_{1,1}(x) = x \quad (9.103)$$

$$b_{0,2}(x) = (1 - x)^2$$

$$b_{1,2}(x) = 2x(1 - x)$$

$$b_{2,2}(x) = x^2 \quad (9.104)$$



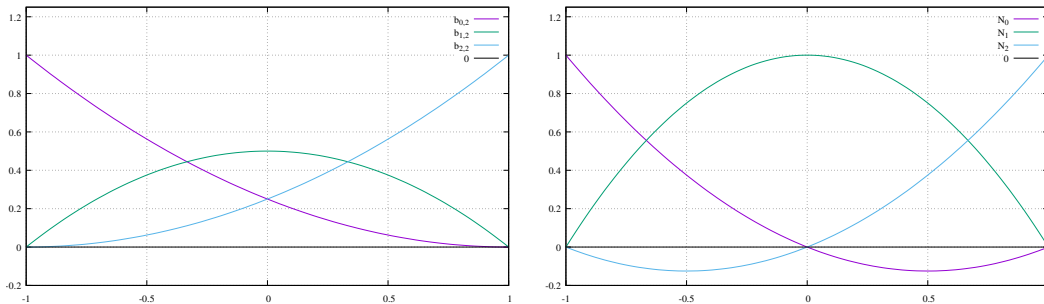
We see that the zero-th and first order polynomials are the same as the linear basis functions defined in Section 5.2. However the second order polynomials (and higher) differ from the second-order basis functions.

Also, the Bernstein polynomials have a lot of properties, but one that is of importance to us is the following: $b_{m,n}(x) \geq 0 \quad \forall x \in [0, 1]$, i.e. the polynomials are positive. This is however not true for basis functions for $n \geq 2$. Another important property shared with basis functions is that their sum over the interval is exactly 1, i.e. $\sum_m b_{m,n}(x) = 1$.

In order to facilitate the comparison between the 2nd-order basis functions and Bernstein polynomials, I will express the latter as a function of the reduced coordinate $r \in [-1, 1] = 2(x - 1/2)$ (or $x = (r + 1)/2$). We have then:

$$\begin{aligned} b_{0,2}(r) &= \frac{1}{4}(1 - r)^2 \\ b_{1,2}(r) &= \frac{1}{2}(1 - r^2) \\ b_{2,2}(r) &= \frac{1}{4}(1 + r)^2 \end{aligned} \quad (9.105)$$

Both 2nd-order Bernstein polynomials and basis functions are plotted here under:



³⁶https://en.wikipedia.org/wiki/Bernstein_polynomial

Having reached this point, the burning question is why should we care?

Example 1 In order to answer this question, let us carry out the following experiment: each node i in the element carries a field value f_i and for simplicity, we choose $f_0 = f(r = -1) = 1$, $f_1 = f(r = 0) = 0$, $f_2 = f(r = +1) = 0$. Then, we can compute the value of the field inside of the element as we usually do in the FE methodology:

$$f^h(r) = \sum_{i=0}^2 f_i N_i(r) = f_0 N_0(r) = N_0(r) = \frac{1}{2}r(r-1)$$

This means that although the field f is always positive (or null) inside the element its representation with the basis functions is negative over half (!) of the element (see purple curve on the right panel above). If we now turn to the Bernstein polynomials:

$$f^h(r) = \sum_{i=0}^2 b_{i,2} = b_{0,2} = \frac{1}{4}(1-r)^2$$

which is *always* positive over the interval $[-1, +1]$, and looking at the purple curve on the left panel above, we see that the value decreases monotonously when we go away from node 1, and reaches zero at the other end of the element.

Also:

$$\text{Shape function: } \int_{-1}^{+1} f^h(r) dr = \int_{-1}^{+1} \frac{1}{2}r(r-1) dr = \frac{1}{3}$$

$$\text{Bernstein polynomial: } \int_{-1}^{+1} f^h(r) dr = \int_{-1}^{+1} \frac{1}{4}(1-r)^2 dr = \frac{2}{3}$$

Analytical value for the integral can be obtained by splitting the integral as $\int_{-1}^0 + \int_0^{+1}$. The left part can be represented by a line of equation $-r$ and the right part simply by 0, so that the integral is equal to $1/2$. We then see that the Shape function-based interpolation underestimates the integral while the Bernstein polynomial-based interpolation overestimates it.

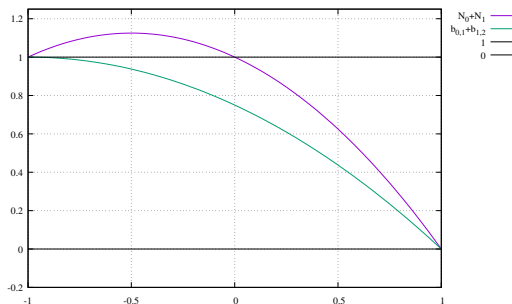
Example 2 We now choose $f_0 = f(r = -1) = 1$, $f_1 = f(r = 1) = 0$, $f_2 = f(r = +1) = 0$. Then

$$f^h(r) = \sum_{i=0}^2 f_i N_i(r) = f_0 N_0(r) f(1) N_1(r) = N_0(r) + N_1(r) = \frac{1}{2}r(r-1) + 1 - r^2 = -\frac{1}{2}r^2 - \frac{1}{2}r + 1$$

Looking now at the Bernstein polynomials:

$$f^h(r) = \sum_{i=0}^2 b_{i,2} = b_{0,2} + b_{1,2} = \frac{1}{4}(1-r)^2 + \frac{1}{2}(1-r^2) = -\frac{1}{4}r^2 - \frac{1}{2}r + \frac{3}{4}$$

If we now plot both approximations:



We see that in this case the basis function-based approximation yields values above 1 over half of the element while the Bernstein polynomial-based approximation remains between 0 and 1 as expected.

Approximation of polynomials Let us now explore another aspect of such an interpolation based on the Bernstein polynomials and assume that $f(r) = C$. Then

$$f^h(r) = \sum_{i=0}^2 f_i b_{i,2}(r) = C \sum_{i=0}^2 b_{i,2}(r) = C \cdot 1 = C$$

Such interpolation can exactly represent a constant field. Let us assume that $f(r) = ar + b$. Then

$$\begin{aligned} f^h(r) &= \sum_{i=0}^2 f(r_i) b_{i,2}(r) \\ &= \sum_{i=0}^2 (ar_i + b) b_{i,2}(r) \\ &= a \sum_{i=0}^2 r_i b_{i,2}(r) + b \sum_{i=1}^3 b_{i,2}(r) \\ &= a(-b_{0,2}(r) + b_{2,2}(r)) + b \cdot 1 \\ &= ar + b \end{aligned} \tag{9.106}$$

Such interpolation can exactly represent a linear field.

Let us assume that $f(r) = ar^2 + br + c$. Then

$$\begin{aligned} f^h(r) &= \sum_{i=0}^2 f(r_i) b_{i,2}(r) \\ &= \sum_{i=0}^2 (ar^2 + br + c) b_{i,2}(r) \\ &= \sum_{i=0}^2 r_i^2 b_{i,2}(r) + b \sum_{i=0}^2 r_i b_{i,2}(r) + c \sum_{i=0}^2 b_{i,2}(r) \\ &= a \sum_{i=0}^2 r_i^2 b_{i,2}(r) + br + c \\ &= a(b_{0,2}(r) + b_{2,2}(r)) + br + c \\ &= a \frac{1}{2} (1 + r^2) + br + c \end{aligned} \tag{9.107}$$

which is not equal to $f(r)$.

On the other hand it is trivial to show that

$$\begin{aligned}
f^h(r) &= \sum_{i=0}^2 f(r_i) N_i(r) \\
&= \sum_{i=0}^2 (ar_i^2 + br_i + c) N_i(r) \\
&= a \sum_{i=0}^2 r_i^2 N_i(r) + b \sum_{i=0}^2 r_i N_i(r) + c \sum_{i=0}^2 N_i(r) \\
&= a \sum_{i=0}^2 r_i^2 N_i(r) + b \sum_{i=0}^2 r_i N_i(r) + c \\
&= a(N_0(r) + N_2(r)) + b(-N_0(r) + N_2(r)) + c \\
&= ar^2 + br + c
\end{aligned} \tag{9.108}$$

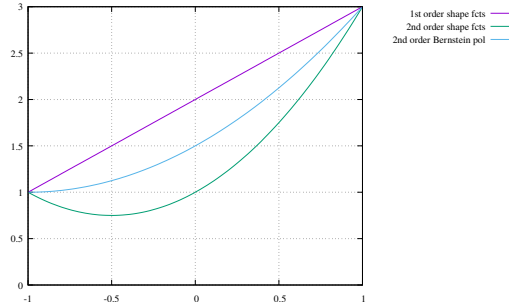
To hammer the point once more: let $f(r) = r^2 + r + 1$. Then

$$\begin{aligned}
f_{Q_1}^h &= f(-1) \frac{1}{2}(1-r) + f(+1) \frac{1}{2}(1+r) \\
&= \frac{1}{2}(1-r) + 3 \frac{1}{2}(1+r) \\
&= 2-r
\end{aligned} \tag{9.109}$$

$$f_{Q_2}^h = r^2 + r + 1 \tag{9.110}$$

$$f_{B_2}^h = \frac{1}{2}(1+r^2) + r + 1 \tag{9.111}$$

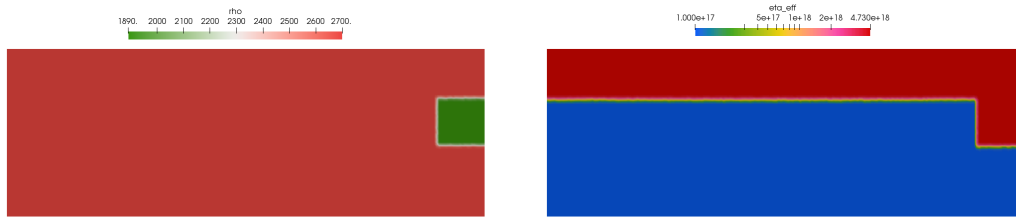
We see on the following figure that although Bernstein polynomials cannot represent $f(r)$ exactly they still do a better job than first order basis functions (Q_1 projection).



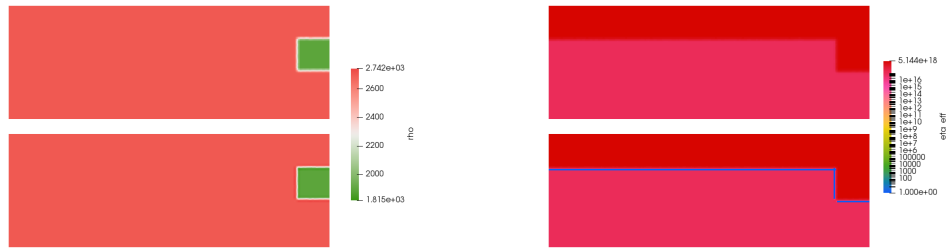
As a conclusion there is a trade-off: 2nd-order Bernstein polynomials *always* yield positive values when the field is positive (as opposed to 2nd-order basis functions) but they cannot represent exactly a 2nd-order polynomial field (while basis functions can).

The positivity can be really critical in geodynamical simulations: a negative density makes no sense, and a negative viscosity even less!

The 2nd-order Bernstein polynomials are used in Stone ???. The actual context of this stone is not important. Fields such as density and viscosity are known on the 9 nodes of the Q_2 element and need to be projected onto the 9 quadrature points. For instance, these nodal fields are given by:



The resulting fields on the quadrature points are shown:



The bottom row is obtained with the basis functions while the top row is obtained with the Bernstein polynomials as interpolants. The thin blue line actually indicated points with negative viscosity and on the left the colour bar shows densities below the value of 1890 (lowest density in the domain).

9.31 Conservative Velocity Interpolation (CVI)

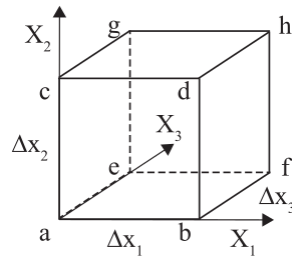
cvi.tex

To my knowledge the conservative velocity interpolation (CVI) was introduced to the computational geodynamics community in Wang, Agrusta, and Hunen [1337] (2015). As mentioned in the paper “An improved velocity interpolation scheme that conserves the divergence of the flow field has been developed by Jenny, Pope, Muradoglu, and Caughey [644] (2001) and the simplified scheme for incompressible flow (i.e., divergence free) has been demonstrated that it largely eliminates the spurious distribution of particles for 2D incompressible flow problem (see Meyer and Jenny [868] (2004)).”

Additional more recent publications on the topic of accurate marker advection: Sime, Maljaars, Wilson, and Keken [1168] (2021), Sime, Wilson, and Keken [1170] (2022).

9.31.1 A few remarks about Wang *et al.* (2015)

The article by Wang, Agrusta, and Hunen [1337] (2015) comes with supplementary material with more details on the derivation of the corrective velocities but that material is a Word document printed to pdf with an annoying layout of equations, different font sizes, lack of alignment, etc ... Also, Fig. 1 of the paper is reproduced here:



Why the authors chose to label nodes a,b,...h and not 1,2,...8 shall forever remain a mystery, but it is not as problematic as the labelling of the axes: indeed, if X_1 is the x -axis then X_3 should be the y -axis and X_2 the z -axis. That is quite illogical. Or is it a mistake in the drawing only? In any case this sheds some confusion on the equations presented in the paper so I have decided to carry out all the CVI derivations in this chapter.

Their paper does not seem to consider cases where the element is not a cuboid (so what about CitcomS, or ALE formulations?), nor does it address higher order elements. Finally many details of the setups in the paper are just not there and I had to email the author(s) multiple time regarding:

- the setup of the couette flow in section 3.1 is incomplete: for instance, size of the box ? velocity value ? exact formula for the vel field (couette flow, I know, but how thick are the layers before rotation)? etc ...

Wang answered me: “The box is a unit box (nondimensional 1*1). I attached the function for the analytical solution for the exact formula for the velocity field that you asked. I didn’t find the models file yet, so I can’t tell you what it is the value of the velocity. But I think it can be: 1m*1m box with 1m/s on the surface (V0). In Citcom, the timestep is chosen to let any material in one cell not to move more than half of the cell length (CFL=0.5). Then we have this parameter ”finetunedt” (< 1) to multiply it. I remember I usually use 0.9 or 0.7. So the CFL=0.45 or 0.35. Concerning the Couette flow we used a viscosity of 1e3, which make very sharp velocity contrast across the diagonal line.”

```
for (i=1;i<=E->lmesh.nno;i++)
{
```

```

x = E->X[1][i];
z = E->X[2][i];
eta1=E->control.testvelval[1];
eta2=E->control.testvelval[2];
alpha=E->control.testvelval[3]*PI/180; /*coordinate rotation angle */
V0=E->control.testvelval[4];
h=sqrt(2.0)*sin(alpha+PI/4); /*WHL: h (with analytical solution) is a function of the rotation angle */
V1=(x*sin(alpha)+z*cos(alpha))*2*V0*eta2/(eta1+eta2)/h;
V2=(x*sin(alpha)+z*cos(alpha))*2*V0*eta1/(eta1+eta2)/h+(eta2-eta1)*V0/(eta1+eta2);
if (x*sin(alpha)+z*cos(alpha)<0.5*h)
{
E->V[1][i]=V1*cos(alpha);
E->V[2][i]=-V1*sin(alpha);
}
else
{
E->V[1][i]=V2*cos(alpha);
E->V[2][i]=-V2*sin(alpha);
}
if (E->mesh.nsd == 3)
E->V[3][i]=0.;
}


```

- which advection scheme was used and I am worried that at no point in the publication the timestep size is either mentioned nor its importance discussed.

Wang answered: “About the timestep, my experience is that using smaller timestep would’t solve this kind of problem. Otherwise we probably would not need to use this new velocity interpolation. I could not remember that I tested the effects of timestep for this model. So it would be nice to know the result if you test it. The advection scheme is the 2nd Runge Kutta.”

- Agrusta wrote: ”here the input values for the couette flow: testvelval=100000,1,45,0.01 # eta1,eta2,angle,velocity. mesh = 33x33. initial tracers 100X100, random distribution”

Looking at their Fig. 2a,b we see black arrow tips in the blue region where velocity should be zero. Velocity is indeed zero and the authors confirmed that the arrow tips are an artefact of their visualisation software (!).

 **Relevant Literature:** McNally (2011) [858] proposed a divergence-free interpolation of vector fields from point values in the context of magnetohydrodynamics. Pusok, Kaus, and Popov [1021] (2016) has applied the CVI to staggered grid FDM.

9.31.2 In 2D with Q_1 basis functions - Naive approach

Let us start directly in reduced coordinates $(r, s) \in [-1 : 1]^2$ (i.e. the reference element). The velocity components inside of the element are given by:

$$\begin{aligned}
u^h(r, s) &= \sum_i \mathcal{N}_i(r, s) u_i \\
v^h(r, s) &= \sum_i \mathcal{N}_i(r, s) v_i
\end{aligned}$$

where \mathcal{N}_i are the four Q_1 basis functions defined as follows:

$$\begin{aligned}\mathcal{N}_1(r, s) &= \frac{1}{4}(1-r)(1-s) \\ \mathcal{N}_2(r, s) &= \frac{1}{4}(1+r)(1-s) \\ \mathcal{N}_3(r, s) &= \frac{1}{4}(1+r)(1+s) \\ \mathcal{N}_4(r, s) &= \frac{1}{4}(1-r)(1+s)\end{aligned}$$

The incompressibility constraint in the (r, s) -coordinate system reads

$$(\vec{\nabla} \cdot \vec{v})^h = \frac{\partial u^h}{\partial r} + \frac{\partial v^h}{\partial s} = \sum_i \left(\frac{\partial \mathcal{N}_i}{\partial r} u_i + \frac{\partial \mathcal{N}_i}{\partial s} v_i \right) = 0.$$

However, it is trivial to verify that the incompressibility condition is not and *can not* be verified for all values of $r, s \in [-1, 1]^2$. It would then make sense to think of a corrective term to the interpolation which would add just enough degrees of freedoms so as to insure an exact³⁷ incompressibility in the element. Let us then write:

$$\begin{aligned}u^h(r, s) &= \sum_i \mathcal{N}_i(r, s) u_i + (as + b)(1-r)(1+r) \\ v^h(r, s) &= \sum_i \mathcal{N}_i(r, s) v_i + (cr + d)(1-s)(1+s)\end{aligned}$$

Note that in this way the correction is zero on the $x = -1$ and $x = +1$ sides of the element for u , and likewise for v on the top and bottom sides (in other words the velocity remains continuous from one element to another). In this case,

$$\begin{aligned}\frac{\partial u^h}{\partial r} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + b)(-2r) \\ \frac{\partial v^h}{\partial s} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (cr + d)(-2s)\end{aligned}$$

We have introduced 4 coefficients (a, b, c, d) which remain to be determined. We start with:

$$\begin{aligned}\sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i &= -\frac{1}{4}(1-s)u_1 + \frac{1}{4}(1-s)u_2 + \frac{1}{4}(1+s)u_3 - \frac{1}{4}(1+s)u_4 \\ &= (1-s)\frac{u_2 - u_1}{4} + (1+s)\frac{u_3 - u_4}{4} \\ &= (1-s)u_{21} + (1+s)u_{34} \\ \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i &= -\frac{1}{4}(1-r)v_1 - \frac{1}{4}(1+r)v_2 + \frac{1}{4}(1+r)v_3 + \frac{1}{4}(1-r)v_4 \\ &= (1-r)\frac{v_4 - v_1}{4} + (1+r)\frac{v_3 - v_2}{4} \\ &= (1-r)v_{41} + (1+r)v_{32}\end{aligned}$$

where $u_{ij} = (u_i - u_j)/4$ and $v_{ij} = (v_i - v_j)/4$, so that in the end

$$\frac{\partial u^h}{\partial r} = (1-s)u_{21} + (1+s)u_{34} + (as + b)(-2r) \quad (9.112)$$

$$\frac{\partial v^h}{\partial s} = (1-r)v_{41} + (1+r)v_{32} + (cr + d)(-2s) \quad (9.113)$$

³⁷more on this later

The incompressibility condition is now:

$$(\vec{\nabla} \cdot \vec{v})^h = (1-s)u_{21} + (1+s)u_{34} + (as+b)(-2r) + (1-r)v_{41} + (1+r)v_{32} + (cr+d)(-2s) = 0$$

This can be rewritten as

$$(\vec{\nabla} \cdot \vec{v})^h = C_0 + C_1r + C_2s + C_3rs = 0$$

where the four C_i coefficients are functions of the velocities and the other coefficients. In order for this expression to be exactly zero *everywhere*, each C coefficient has to be independently zero.

$$\begin{aligned} C_0 \quad (.) \quad u_{21} + u_{34} + v_{41} + v_{32} &= 0 \\ C_1 \quad (r) \quad -v_{41} + v_{32} - 2b &= 0 \\ C_2 \quad (s) \quad -u_{21} + u_{34} - 2d &= 0 \\ C_3 \quad (rs) \quad -2a - 2c &= 0 \end{aligned}$$

The first line is simply the incompressibility condition expressed in the center of the element (i.e. $r = s = 0$), so we set it aside for now (I will come back to it later!) and focus on the remaining three.

At this stage it is important to note that in the absence of corrective terms (i.e. $a = b = c = d = 0$) then only $C_3 = 0$ and the divergence inside the element is a linear field.

We obtain

$$c = -a \quad b = \frac{1}{2}(-v_{41} + v_{32}) \quad d = \frac{1}{2}(-u_{21} + u_{34})$$

Since a and c are not otherwise constrained, we can set them to zero, and we then have:

$$b = \frac{1}{2}(v_{14} + v_{32}) \quad d = \frac{1}{2}(u_{12} + u_{34})$$

and finally

$$\begin{aligned} u^h(r, s) &= \sum_i \mathcal{N}_i(r, s)u_i + b(1-r)(1+r) = \sum_i \mathcal{N}_i(r, s)u_i + \frac{1}{2}(v_{14} + v_{32})(1-r)(1+r) \\ v^h(r, s) &= \sum_i \mathcal{N}_i(r, s)v_i + d(1-s)(1+s) = \sum_i \mathcal{N}_i(r, s)v_i + \frac{1}{2}(u_{12} + u_{34})(1-s)(1+s) \end{aligned}$$

By using these corrected interpolations for both components of the velocity then one ensures that a point-wise divergence free velocity field anywhere in the element. However, these derivations were carried out in the reference element. In fact they would work also for rectangular elements with minimal changes, but not for generic quadrilaterals.

To be clear, let us now compute the velocity divergence of the corrected velocity field above:

$$\begin{aligned} (\vec{\nabla} \cdot \vec{v})^h &= \frac{\partial u^h}{\partial r} + \frac{\partial v^h}{\partial s} \\ &= (1-s)u_{21} + (1+s)u_{34} + \frac{1}{2}(v_{14} + v_{32})(-2r) + (1-r)v_{41} + (1+r)v_{32} + \frac{1}{2}(u_{12} + u_{34})(-2s) \\ &= u_{21} + u_{34} + v_{41} + v_{32} - su_{21} + su_{34} - rv_{14} - rv_{32} - rv_{41} + rv_{32} - su_{12} - su_{34} \\ &= u_{21} + u_{34} + v_{41} + v_{32} \end{aligned} \tag{9.114}$$

A point must then be made crystal clear: the divergence is *not* zero. The quantity above is constant inside the element (it does not depend on r nor s). **All what the CVI algorithm does is to remove the spatial dependence of the velocity divergence inside the element.**

9.31.3 In 2D with Q_1 basis functions - better approach

We now consider a generic quadrilateral in the x, y -coordinate space and its equivalent in the reference space r, s . One can easily show that the gradient of a field f verifies

$$\begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} = \tilde{\mathbf{J}} \cdot \begin{pmatrix} \frac{\partial f}{\partial r} \\ \frac{\partial f}{\partial s} \end{pmatrix}$$

where $\tilde{\mathbf{J}}$ is the inverse of the Jacobian matrix. We then postulate again

$$\begin{aligned} u^h(r, s) &= \sum_i \mathcal{N}_i(r, s) u_i + (as + b)(1 - r)(1 + r) \\ v^h(r, s) &= \sum_i \mathcal{N}_i(r, s) v_i + (cr + d)(1 - s)(1 + s) \end{aligned}$$

In this case,

$$\begin{aligned} \frac{\partial u^h}{\partial r} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + b)(-2r) \\ \frac{\partial u^h}{\partial s} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + a(1 - r^2) \\ \frac{\partial v^h}{\partial r} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial r} v_i + c(1 - s^2) \\ \frac{\partial v^h}{\partial s} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (cr + d)(-2s) \end{aligned}$$

We have introduced 4 coefficients (a, b, c, d) which remain to be determined. In order to compute the velocity divergence inside the element we will need

$$\begin{aligned} \frac{\partial u}{\partial x} &= \tilde{J}_{xx} \frac{\partial u}{\partial r} + \tilde{J}_{xy} \frac{\partial u}{\partial s} \\ &= \tilde{J}_{xx} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + b)(-2r) \right) + \tilde{J}_{xy} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + a(1 - r^2) \right) \\ &= \tilde{J}_{xx} (-(1 - s)u_{12} + (1 + s)u_{34} + (as + b)(-2r)) \\ &\quad + \tilde{J}_{xy} (-(1 - r)u_{14} - (1 + r)u_{23} + a(1 - r^2)) \\ \frac{\partial v}{\partial y} &= \tilde{J}_{yx} (-(1 - s)v_{12} + (1 + s)v_{34} + c(1 - s^2)) \\ &\quad + \tilde{J}_{yy} (-(1 - r)v_{14} - (1 + r)v_{23} + (cr + d)(-2s)) \end{aligned}$$

where $u_{ij} = (u_i - u_j)/4$ and $v_{ij} = (v_i - v_j)/4$. The velocity divergence can be written as follows

$$\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = C_0 + C_1 r + C_2 s + C_3 rs + C_4 r^2 + C_5 s^2 = 0$$

with

$$\begin{aligned}
C_0 &= J_{xx}(-u_{12} + u_{34}) + J_{xy}(-u_{14} - u_{23}) + J_{yx}(-v_{12} + v_{34}) + J_{yy}(-v_{14} - v_{23}) \\
C_1 &= J_{xy}(u_{14} - u_{23}) + J_{yy}(v_{14} - v_{23}) - 2bJ_{xx} \\
C_2 &= J_{xx}(u_{12} + u_{34}) + J_{yx}(v_{12} + v_{34}) - 2dJ_{yy} \\
C_3 &= -2aJ_{xx} - 2cJ_{yy} \\
C_4 &= -aJ_{xy} \\
C_5 &= -cJ_{yx}
\end{aligned} \tag{9.115}$$

where the six C_i coefficients are functions of the velocities and the other coefficients. In order for this expression to be exactly null *everywhere*³⁸, each C coefficient has to be independently null.

This immediately yields $a = c = 0$ (since the components of the $\tilde{\mathbf{J}}$ tensor are not necessarily zero - and if J_{xy} and J_{yx} are zero then the equation for C_3 remains and we would still take $a = c = 0$ for simplicity) and the equation for C_3 is immediately satisfied. We then have:

$$\begin{aligned}
b &= \frac{1}{2J_{xx}}(J_{xy}(u_{14} - u_{23}) + J_{yy}(v_{14} - v_{23})) \\
d &= \frac{1}{2J_{yy}}(J_{xx}(u_{12} + u_{34}) + J_{yx}(v_{12} + v_{34}))
\end{aligned}$$

These expressions contain the same ingredients as before but also introduce more coupling between the velocity components. If the element is rectangular then $J_{xy} = J_{yx} = 0$ and

$$\begin{aligned}
b &= \frac{J_{yy}}{2J_{xx}}(v_{14} - v_{23}) \\
d &= \frac{J_{xx}}{2J_{yy}}(u_{12} + u_{34})
\end{aligned}$$

If the element is square then $J_{xx} = J_{yy} = 0$ so

$$\begin{aligned}
b &= \frac{1}{2}(v_{14} - v_{23}) \\
d &= \frac{1}{2}(u_{12} + u_{34})
\end{aligned}$$

and finally the velocity correction is

$$\begin{aligned}
\delta u &= \frac{1}{2}(v_{14} - v_{23})(1 - r)(1 + r) \\
\delta v &= \frac{1}{2}(u_{12} + u_{34})(1 - s)(1 + s)
\end{aligned} \tag{9.116}$$

9.31.4 Comparison with Wang *et al.* (2015) for 2D

Rather annoyingly Wang *et al.* (2015) use a reference element that is $[0, 1] \times [0, 1]$ as opposed to the standard $[-1, 1] \times [-1, 1]$:

³⁸We know by now that this is not possible

In a 4-node 2D rectangular cell system, bilinear interpolation provides a simple and quick interpolation scheme and is widely used. If we transform the rectangular cells into unit squares (Fig 1), the interpolation we used can be written as:

$$U_i^L(x_1, x_2) = \{(1-x_1)(1-x_2), x_1(1-x_2), (1-x_1)x_2, x_1x_2\} \cdot \{U_i^a, U_i^b, U_i^c, U_i^d\}, \quad (1)$$

where the two velocity components are interpolated independently as two separate scalars without considering the divergence of the vector field need to be 0.

$$\frac{\partial U_1}{\partial x_1} + \frac{\partial U_2}{\partial x_2} = 0 \quad (2)$$

The 2D divergence-free interpolation is achieved by adding correction items as follows [Meyer and Jenny, 2004]:

$$U_i = U_i^L + \Delta U_i, \quad (3)$$

$$\Delta U_1 = \frac{\Delta x_1}{2\Delta x_2} x_1(1-x_1)(U_2^a - U_2^b - U_2^c + U_2^d), \quad (4)$$

$$\Delta U_2 = \frac{\Delta x_2}{2\Delta x_1} x_2(1-x_2)(U_1^a - U_1^b - U_1^c + U_1^d). \quad (5)$$

Taken from the supplementary material of Wang *et al.* (2015).

Since basis functions must be 1 on their node, then the numbering must be as follows:

$$\begin{array}{ccc} \text{c--d} & & 4--3 \\ | \quad | & \Leftrightarrow & | \quad | \\ \text{a--b} & & 1--2 \end{array}$$

Setting $\Delta x_1 = \Delta x_2 = 1$, replacing a by 1, b by 2, c by 4 and d by 3, x_1 by r' and x_2 by s' , U_1 by u and U_2 by v , we arrive at (in order to render the notations a bit lighter I have set $U = U_1$ and $V = U_2$)

$$\begin{aligned} \Delta U &= \frac{1}{2} r' (1 - r') (v_1 - v_2 - v_4 + v_3) = \frac{1}{2} r' (1 - r') (4v_{14} - 4v_{23}) \\ \Delta V &= \frac{1}{2} s' (1 - s') (u_1 - u_2 - u_4 + u_3) = \frac{1}{2} s' (1 - s') (4u_{12} + 4u_{34}) \end{aligned}$$

Since $r = 2r' - 1$ and $s = 2s' - 1$ then we find that

$$\begin{aligned} \Delta U &= \frac{1}{2} (1 - r^2) (v_{14} - v_{23}) \\ \Delta V &= \frac{1}{2} (1 - s^2) (u_{12} + u_{34}) \end{aligned} \quad (9.117)$$

which is Eq. (9.116). In the case of the reference element then my velocity corrections are identical to theirs.

Let us look at the equations of the figure above. Since the authors state that they “transform the rectangular cells into unit squares” we do away with $\Delta x_1 = \Delta x_2 = 1$. Eqs. 3 and 1 together yield:

$$\begin{aligned} U &= (1-x_1)(1-x_2)U^a + x_1(1-x_2)U^b + (1-x_1)x_2U^c + x_1x_2U^d + \frac{1}{2}x_1(1-x_1)(V^a - V^b - V^c + V^d) \\ V &= (1-x_1)(1-x_2)V^a + x_1(1-x_2)V^b + (1-x_1)x_2V^c + x_1x_2V^d + \frac{1}{2}x_2(1-x_2)(U^a - U^b - U^c + U^d) \end{aligned}$$

Then

$$\begin{aligned} \frac{\partial U}{\partial x_1} &= -(1-x_2)U^a + (1-x_2)U^b - x_2U^c + x_2U^d + \frac{1}{2}(1-2x_1)(V^a - V^b - V^c + V^d) \\ \frac{\partial V}{\partial x_2} &= -(1-x_1)V^a - x_1V^b + (1-x_1)V^c + x_1V^d + \frac{1}{2}(1-2x_2)(U^a - U^b - U^c + U^d) \end{aligned}$$

So

$$\begin{aligned}
\frac{\partial U}{\partial x_1} + \frac{\partial V}{\partial x_2} &= -(1-x_2)U^a + (1-x_2)U^b - x_2U^c + x_2U^d + \frac{1}{2}(1-2x_1)(V^a - V^b - V^c + V^d) \\
&\quad -(1-x_1)V^a - x_1V^b + (1-x_1)V^c + x_1V^d + \frac{1}{2}(1-2x_2)(U^a - U^b - U^c + U^d) \\
&= -U^a + U^b + x_2(U^a - U^b - U^c + U^d) + \frac{1}{2}(V^a - V^b - V^c + V^d) - x_1(V^a - V^b - V^c + V^d) \\
&\quad -V^a + V^c + x_1(V^a - V^b - V^c + V^d) + \frac{1}{2}(U^a - U^b - U^c + U^d) - x_2(U^a - U^b - U^c + U^d) \\
&= -U^a + U^b + \frac{1}{2}(V^a - V^b - V^c + V^d) - V^a + V^c + \frac{1}{2}(U^a - U^b - U^c + U^d) \\
&\neq 0
\end{aligned} \tag{9.118}$$

Unfortunately, the authors seem to be under the impression that this quantity is zero since they talk of “2D divergence-free interpolation” and “the divergence of the vector field need to be zero”. Their own equations prove that this is not the case.

9.31.5 In 3D with Q_1 basis functions - Naive approach

In this case we are addressing the case of the divergence being as close to zero as possible in the reference element. We’ll treat the case of a generic hexahedron in the next section.

Let us start directly in reduced coordinates $(r, s, t) \in [-1 : 1]^3$:

$$\begin{aligned}
u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i \\
v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i \\
w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i
\end{aligned}$$

with

$$\begin{aligned}
\mathcal{N}_1 &= \frac{1}{8}(1-r)(1-s)(1-t) \\
\mathcal{N}_2 &= \frac{1}{8}(1+r)(1-s)(1-t) \\
\mathcal{N}_3 &= \frac{1}{8}(1+r)(1+s)(1-t) \\
\mathcal{N}_4 &= \frac{1}{8}(1-r)(1+s)(1-t) \\
\mathcal{N}_5 &= \frac{1}{8}(1-r)(1-s)(1+t) \\
\mathcal{N}_6 &= \frac{1}{8}(1+r)(1-s)(1+t) \\
\mathcal{N}_7 &= \frac{1}{8}(1+r)(1+s)(1+t) \\
\mathcal{N}_8 &= \frac{1}{8}(1-r)(1+s)(1+t)
\end{aligned}$$

The incompressibility constraint imposes:

$$\frac{\partial u^h}{\partial r} + \frac{\partial v^h}{\partial s} + \frac{\partial w^h}{\partial t} = 0 = \sum_i \left(\frac{\partial \mathcal{N}_i}{\partial r} u_i + \frac{\partial \mathcal{N}_i}{\partial s} v_i + \frac{\partial \mathcal{N}_i}{\partial t} w_i \right) = 0$$

However, once again it is trivial to verify that the incompressibility condition is not and can not be verified for all values of $r, s, t \in [-1, 1]^3$.

It would then make sense to think of a corrective term to the interpolation which would add just enough degrees of freedoms so as to insure an exact incompressibility in the element. Let us then write:

$$\begin{aligned} u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (as + bt + c)(1 - r)(1 + r) \\ v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (dr + et + f)(1 - s)(1 + s) \\ w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (gr + hs + i)(1 - t)(1 + t) \end{aligned}$$

We thereby make sure that the corrections are zero on the edges so that velocity remains continuous from one element to another. In this case,

$$\begin{aligned} \frac{\partial u^h}{\partial r} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + bt + c)(-2r) \\ \frac{\partial v^h}{\partial s} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (dr + et + f)(-2s) \\ \frac{\partial w^h}{\partial t} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial t} w_i + (gr + hs + i)(-2t) \end{aligned}$$

We have introduced 9 coefficients $(a, b, c, d, e, f, g, h, i)$ which remain to be determined. The incompressibility condition is now:

$$\sum_i \left(\frac{\partial \mathcal{N}_i}{\partial r} u_i + \frac{\partial \mathcal{N}_i}{\partial s} v_i + \frac{\partial \mathcal{N}_i}{\partial t} w_i \right) + (as + bt + c)(-2r) + (dr + et + f)(-2s) + (gr + hs + i)(-2t) = 0$$

This can be rewritten as

$$C_0 + C_1 r + C_2 s + C_3 t + C_4 rs + C_5 st + C_6 rt = 0$$

where the seven C_i coefficients are functions of the velocities and the other coefficients. In order for this expression to be exactly zero *everywhere*³⁹, each C coefficient has to be independently zero.

We start with:

$$\begin{aligned} 8 \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i &= (1 - s)(1 - t)(u_2 - u_1) + (1 + s)(1 - t)(u_3 - u_4) + (1 - s)(1 + t)(u_6 - u_5) + (1 + s)(1 + t)(u_7 - u_8) \\ 8 \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i &= (1 - r)(1 - t)(v_4 - v_1) + (1 + r)(1 - t)(v_3 - v_2) + (1 - r)(1 + t)(v_8 - v_5) + (1 + r)(1 + t)(v_7 - v_6) \\ 8 \sum_i \frac{\partial \mathcal{N}_i}{\partial t} w_i &= (1 - r)(1 - s)(w_5 - w_1) + (1 + r)(1 - s)(w_6 - w_2) + (1 + r)(1 + s)(w_7 - w_3) + (1 - r)(1 + s)(w_8 - w_4) \end{aligned}$$

Let us denote $u_{ij} = (u_i - v_j)/8$ (same for v, w), so that:

$$\begin{aligned} \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i &= (1 - s)(1 - t)u_{21} + (1 + s)(1 - t)u_{34} + (1 - s)(1 + t)u_{65} + (1 + s)(1 + t)u_{78} \\ \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i &= (1 - r)(1 - t)v_{41} + (1 + r)(1 - t)v_{32} + (1 - r)(1 + t)v_{85} + (1 + r)(1 + t)v_{76} \\ \sum_i \frac{\partial \mathcal{N}_i}{\partial t} w_i &= (1 - r)(1 - s)w_{51} + (1 + r)(1 - s)w_{62} + (1 + r)(1 + s)w_{73} + (1 - r)(1 + s)w_{84} \end{aligned}$$

³⁹By now we know this is not possible – see 2D

We finally arrive at:

$$\begin{aligned}
C_0 \quad (.) \quad & u_{21} + u_{34} + u_{65} + u_{78} + v_{41} + v_{32} + v_{85} + v_{76} + w_{51} + w_{62} + w_{73} + w_{84} = 0 \\
C_1 \quad (r) \quad & -v_{41} + v_{32} - v_{85} + v_{76} - w_{51} + w_{62} + w_{73} - w_{84} - 2c = 0 \\
C_2 \quad (s) \quad & -u_{21} + u_{34} - u_{65} + u_{78} - w_{51} - w_{62} + w_{73} + w_{84} - 2f = 0 \\
C_3 \quad (t) \quad & -u_{21} - u_{34} + u_{65} + u_{78} - v_{41} - v_{32} + v_{85} + v_{76} - 2i = 0 \\
C_4 \quad (rs) \quad & w_{51} - w_{62} + w_{73} - w_{84} - 2a - 2d = 0 \\
C_5 \quad (st) \quad & u_{21} - u_{34} - u_{65} + u_{78} - 2e - 2h = 0 \\
C_6 \quad (rt) \quad & v_{41} - v_{32} - v_{85} + v_{76} - 2b - 2g = 0
\end{aligned}$$

I leave C_0 alone but I still unfortunately end up with 6 equations and 9 unknowns a, b, c, d, e, f, g, h . Coming up with additional constraints is not trivial, so I will instead further assume $\alpha_r = b = a$, $\alpha_s = e = d$ and $\alpha_t = h = g$, and rename $\beta_r = c$, $\beta_s = f$ and $\beta_t = i$ so that I have now six unknowns $\alpha_r, \alpha_s, \alpha_t, \beta_r, \beta_s, \beta_t$ for six equations

$$\begin{aligned}
C_1 \quad (r) \quad & -v_{41} + v_{32} - v_{85} + v_{76} - w_{51} + w_{62} + w_{73} - w_{84} - 2\beta_r \\
C_2 \quad (s) \quad & -u_{21} + u_{34} - u_{65} + u_{78} - w_{51} - w_{62} + w_{73} + w_{84} - 2\beta_s \\
C_3 \quad (t) \quad & -u_{21} - u_{34} + u_{65} + u_{78} - v_{41} - v_{32} + v_{85} + v_{76} - 2\beta_t \\
C_4 \quad (rs) \quad & w_{51} - w_{62} + w_{73} - w_{84} - 2\alpha_r - 2\alpha_s \\
C_5 \quad (st) \quad & u_{21} - u_{34} - u_{65} + u_{78} - 2\alpha_s - 2\alpha_t \\
C_6 \quad (rt) \quad & v_{41} - v_{32} - v_{85} + v_{76} - 2\alpha_r - 2\alpha_t
\end{aligned}$$

This naturally yields:

$$\begin{aligned}
\beta_r &= \frac{1}{2}(-v_{41} + v_{32} - v_{85} + v_{76} - w_{51} + w_{62} + w_{73} - w_{84}) \\
&= \frac{1}{16}(v_1 - v_2 + v_3 - v_4 + v_5 - v_6 + v_7 - v_8 + w_1 - w_2 - w_3 + w_4 - w_5 + w_6 + w_7 - w_8) \\
\beta_s &= \frac{1}{2}(-u_{21} + u_{34} - u_{65} + u_{78} - w_{51} - w_{62} + w_{73} + w_{84}) \\
&= \frac{1}{16}(u_1 - u_2 + u_3 - u_4 + u_5 - u_6 + u_7 - u_8 + w_1 + w_2 - w_3 - w_4 - w_5 - w_6 + w_7 + w_8) \\
\beta_t &= \frac{1}{2}(-u_{21} - u_{34} + u_{65} + u_{78} - v_{41} - v_{32} + v_{85} + v_{76}) \\
&= \frac{1}{16}(u_1 - u_2 - u_3 + u_4 - u_5 + u_6 + u_7 - u_8 + v_1 + v_2 - v_3 - v_4 - v_5 - v_6 + v_7 + v_8)
\end{aligned}$$

and we need to solve

$$\begin{aligned}
\tilde{w} - 2\alpha_r - 2\alpha_s &= 0 \\
\tilde{u} - 2\alpha_s - 2\alpha_t &= 0 \\
\tilde{v} - 2\alpha_r - 2\alpha_t &= 0
\end{aligned}$$

where

$$\begin{aligned}
\tilde{u} &= u_{21} - u_{34} - u_{65} + u_{78} = \frac{1}{8}(-u_1 + u_2 - u_3 + u_4 + u_5 - u_6 + u_7 - u_8) \\
\tilde{v} &= v_{41} - v_{32} - v_{85} + v_{76} = \frac{1}{8}(-v_1 + v_2 - v_3 + v_4 + v_5 - v_6 + v_7 - v_8) \\
\tilde{w} &= w_{51} - w_{62} + w_{73} - w_{84} = \frac{1}{8}(-w_1 + w_2 - w_3 + w_4 + w_5 - w_6 + w_7 - w_8)
\end{aligned}$$

which yields:

$$\alpha_r = \frac{1}{4}(-\tilde{u} + \tilde{v} + \tilde{w}) \quad \alpha_s = \frac{1}{4}(\tilde{u} - \tilde{v} + \tilde{w}) \quad \alpha_t = \frac{1}{4}(\tilde{u} + \tilde{v} - \tilde{w})$$

So finally:

$$\begin{aligned} u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + [\alpha_r(s + t) + \beta_r](1 - r)(1 + r) \\ v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + [\alpha_s(r + t) + \beta_s](1 - s)(1 + s) \\ w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + [\alpha_t(r + s) + \beta_t](1 - t)(1 + t) \end{aligned}$$

9.31.6 In 3D with Q_1 basis functions - better approach

We start again from

$$\begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial z} \end{pmatrix} = \tilde{\mathbf{J}} \cdot \begin{pmatrix} \frac{\partial f}{\partial r} \\ \frac{\partial f}{\partial s} \\ \frac{\partial f}{\partial t} \end{pmatrix}$$

where $\tilde{\mathbf{J}}$ is the inverse of the Jacobian matrix \mathbf{J} . We then postulate

$$\begin{aligned} u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (as + bt + c)(1 - r)(1 + r) \\ v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (dr + et + f)(1 - s)(1 + s) \\ w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (gr + hs + i)(1 - t)(1 + t) \end{aligned}$$

so that:

$$\begin{aligned} \frac{\partial u}{\partial x} &= \tilde{J}_{xx} \frac{\partial u^h}{\partial r} + \tilde{J}_{xy} \frac{\partial u}{\partial s} + \tilde{J}_{xz} \frac{\partial u}{\partial t} \\ &= \tilde{J}_{xx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + bt + c)(-2r) \right] + \tilde{J}_{xy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + a(1 - r^2) \right] + \tilde{J}_{xz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} u_i + b(1 - r^2) \right] \\ \frac{\partial v}{\partial y} &= \tilde{J}_{yx} \frac{\partial v^h}{\partial r} + \tilde{J}_{yy} \frac{\partial v}{\partial s} + \tilde{J}_{yz} \frac{\partial v}{\partial t} \\ &= \tilde{J}_{yx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} v_i + d(1 - s^2) \right] + \tilde{J}_{yy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (dr + et + f)(-2s) \right] + \tilde{J}_{yz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} v_i + e(1 - s^2) \right] \\ \frac{\partial w}{\partial z} &= \tilde{J}_{zx} \frac{\partial w^h}{\partial r} + \tilde{J}_{zy} \frac{\partial w}{\partial s} + \tilde{J}_{zz} \frac{\partial w}{\partial t} \\ &= \tilde{J}_{zx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} w_i + g(1 - t^2) \right] + \tilde{J}_{zy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} w_i + h(1 - t^2) \right] + \tilde{J}_{zz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} w_i + (gr + hs + i)(-2t) \right] \end{aligned}$$

where for any function f :

$$\begin{aligned}
\sum_i \frac{\partial \mathcal{N}_i}{\partial r} f_i &= (1-s)(1-t)f_{21} + (1-s)(1+t)f_{65} + (1+s)(1-t)f_{34} + (1+s)(1+t)f_{78} \\
&= (f_{21} + f_{65} + f_{34} + f_{78}) \\
&+ (-f_{21} - f_{65} + f_{34} + f_{78})s \\
&+ (-f_{21} + f_{65} - f_{34} + f_{78})t \\
&+ (f_{21} - f_{65} - f_{34} + f_{78})st \\
&= f_{r1} + f_{r2}s + f_{r3}t + f_{r4}st \\
\sum_i \frac{\partial \mathcal{N}_i}{\partial s} f_i &= (1-r)(1-t)f_{41} + (1+r)(1-t)f_{32} + (1-r)(1+t)f_{85} + (1+r)(1+t)f_{76} \\
&= (f_{41} + f_{32} + f_{85} + f_{76}) \\
&+ (-f_{41} + f_{32} - f_{85} + f_{76})r \\
&+ (-f_{41} - f_{32} + f_{85} + f_{76})t \\
&+ (f_{41} - f_{32} - f_{85} + f_{76})rt \\
&= f_{s1} + f_{s2}r + f_{s3}t + f_{s4}rt \\
\sum_i \frac{\partial \mathcal{N}_i}{\partial t} f_i &= (1-r)(1-s)f_{51} + (1+r)(1-s)f_{62} + (1+r)(1+s)f_{73} + (1-r)(1+s)f_{84} \\
&= (f_{51} + f_{62} + f_{73} + f_{84}) \\
&+ (-f_{51} + f_{62} + f_{73} - f_{84})r \\
&+ (-f_{51} - f_{62} + f_{73} + f_{84})s \\
&+ (f_{51} - f_{62} + f_{73} - f_{84})rs \\
&= f_{t1} + f_{t2}r + f_{t3}s + f_{t4}rs
\end{aligned}$$

The velocity divergence is then

$$\begin{aligned}
& \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} \\
&= \tilde{J}_{xx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (as + bt + c)(-2r) \right] + \tilde{J}_{xy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + a(1 - r^2) \right] + \tilde{J}_{xz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} u_i + b(1 - r^2) \right] \\
&+ \tilde{J}_{yx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} v_i + d(1 - s^2) \right] + \tilde{J}_{yy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (dr + et + f)(-2s) \right] + \tilde{J}_{yz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} v_i + e(1 - s^2) \right] \\
&+ \tilde{J}_{zx} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial r} w_i + g(1 - t^2) \right] + \tilde{J}_{zy} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial s} w_i + h(1 - t^2) \right] + \tilde{J}_{zz} \left[\sum_i \frac{\partial \mathcal{N}_i}{\partial t} w_i + (gr + hs + i)(-2t) \right] \\
&= \tilde{J}_{xx} [u_{r1} + u_{r2}s + u_{r3}t + u_{r4}st + (as + bt + c)(-2r)] \\
&+ \tilde{J}_{xy} [u_{s1} + u_{s2}r + u_{s3}t + u_{s4}rt + a(1 - r^2)] \\
&+ \tilde{J}_{xz} [u_{t1} + u_{t2}r + u_{t3}s + u_{t4}rs + b(1 - r^2)] \\
&+ \tilde{J}_{yx} [v_{r1} + v_{r2}s + v_{r3}t + v_{r4}st + d(1 - s^2)] \\
&+ \tilde{J}_{yy} [v_{s1} + v_{s2}r + v_{s3}t + v_{s4}rt + (dr + et + f)(-2s)] \\
&+ \tilde{J}_{yz} [v_{t1} + v_{t2}r + v_{t3}s + v_{t4}rs + e(1 - s^2)] \\
&+ \tilde{J}_{zx} [w_{r1} + w_{r2}s + w_{r3}t + w_{r4}st + g(1 - t^2)] \\
&+ \tilde{J}_{zy} [w_{s1} + w_{s2}r + w_{s3}t + w_{s4}rt + h(1 - t^2)] \\
&+ \tilde{J}_{zz} [w_{t1} + w_{t2}r + w_{t3}s + w_{t4}rs + (gr + hs + i)(-2t)] \\
&= C_0 + C_1r + C_2s + C_3t + C_4rs + C_5st + C_6rt + C_7r^2 + C_8s^2 + C_9t^2 = 0
\end{aligned} \tag{9.119}$$

with:

$$\begin{aligned}
C_0 &= \tilde{J}_{xx}u_{r1} + \tilde{J}_{xy}u_{s1} + \tilde{J}_{xz}u_{t1} + \tilde{J}_{yx}v_{r1} + \tilde{J}_{yy}v_{s1} + \tilde{J}_{yz}v_{t1} + \tilde{J}_{zx}w_{r1} + \tilde{J}_{zy}w_{s1} + \tilde{J}_{zz}w_{t1} \\
&+ \tilde{J}_{xy}a + \tilde{J}_{xz}b + \tilde{J}_{yx}d + \tilde{J}_{yz}e + \tilde{J}_{zx}g + \tilde{J}_{zy}h \\
C_1 &= \tilde{J}_{xy}u_{s2} + \tilde{J}_{xz}u_{t2} + \tilde{J}_{yy}v_{s2} + \tilde{J}_{yz}v_{t2} + \tilde{J}_{zy}w_{s2} + \tilde{J}_{zz}w_{t2} - \tilde{J}_{xx}2c \\
C_2 &= \tilde{J}_{xx}u_{r2} + \tilde{J}_{xy}u_{s3} + \tilde{J}_{yx}v_{r2} + \tilde{J}_{yz}v_{t3} + \tilde{J}_{zx}w_{r2} + \tilde{J}_{zz}w_{t3} - \tilde{J}_{yy}2f \\
C_3 &= \tilde{J}_{xx}u_{r3} + \tilde{J}_{xy}u_{s3} + \tilde{J}_{yx}v_{r3} + \tilde{J}_{yy}v_{s3} + \tilde{J}_{zx}w_{r3} + \tilde{J}_{zy}w_{s3} - \tilde{J}_{zz}2i \\
C_4 &= \tilde{J}_{xz}u_{t4} + \tilde{J}_{yz}v_{t4} + \tilde{J}_{zz}w_{t4} - \tilde{J}_{xx}2a - \tilde{J}_{yy}2d \\
C_5 &= \tilde{J}_{xx}u_{r4} + \tilde{J}_{yx}v_{r4} + \tilde{J}_{zx}w_{r4} - \tilde{J}_{yy}2e - \tilde{J}_{zz}2h \\
C_6 &= \tilde{J}_{xy}u_{s4} + \tilde{J}_{yy}v_{s4} + \tilde{J}_{zy}w_{s4} - \tilde{J}_{xx}2b - \tilde{J}_{zz}2g \\
C_7 &= -\tilde{J}_{xy}a - \tilde{J}_{xz}b \\
C_8 &= -\tilde{J}_{yx}d - \tilde{J}_{yz}e \\
C_9 &= -\tilde{J}_{zx}g - \tilde{J}_{zy}h
\end{aligned}$$

Of course what we want is a point-wise zero velocity divergence so we would need $C_0 = C_1 = \dots C_9 = 0$. However we have 10 C coefficients/equations and only 9 variables $a, b, c, d, e, f, g, h, i$. We leave the C_0 equation alone and hope for the best (see 2D case). In other words we hope that if/when we have found $a, b, c, d, e, f, g, h, i$ so that $C_1 = \dots C_9 = 0$ then C_0 is 'small' (whatever that means). As mentioned earlier, the CVI only removes the spatial dependence of the velocity divergence inside an element, it does not zero it.

It is then trivial to obtain c, f, i from the equations of C_1, C_2, C_3 :

$$\begin{aligned}
C_1 = 0 &\Rightarrow \tilde{J}_{xy}u_{s2} + \tilde{J}_{xz}u_{t2} + \tilde{J}_{yy}v_{s2} + \tilde{J}_{yz}v_{t2} + \tilde{J}_{zy}w_{s2} + \tilde{J}_{zz}w_{t2} - \tilde{J}_{xx}2c = 0 \\
c &= \frac{1}{2\tilde{J}_{xx}}(\tilde{J}_{xy}u_{s2} + \tilde{J}_{xz}u_{t2} + \tilde{J}_{yy}v_{s2} + \tilde{J}_{yz}v_{t2} + \tilde{J}_{zy}w_{s2} + \tilde{J}_{zz}w_{t2}) \\
C_2 = 0 &\Rightarrow \tilde{J}_{xx}u_{r2} + \tilde{J}_{xz}u_{t3} + \tilde{J}_{yx}v_{r2} + \tilde{J}_{yz}v_{t3} + \tilde{J}_{zx}w_{r2} + \tilde{J}_{zz}w_{t3} - \tilde{J}_{yy}2f = 0 \\
f &= \frac{1}{2\tilde{J}_{yy}}(\tilde{J}_{xx}u_{r2} + \tilde{J}_{xz}u_{t3} + \tilde{J}_{yx}v_{r2} + \tilde{J}_{yz}v_{t3} + \tilde{J}_{zx}w_{r2} + \tilde{J}_{zz}w_{t3}) \\
C_3 = 0 &\Rightarrow \tilde{J}_{xx}u_{r3} + \tilde{J}_{xy}u_{s3} + \tilde{J}_{yx}v_{r3} + \tilde{J}_{yy}v_{s3} + \tilde{J}_{zx}w_{r3} + \tilde{J}_{zy}w_{s3} - \tilde{J}_{zz}2i = 0 \\
i &= \frac{1}{2\tilde{J}_{zz}}(\tilde{J}_{xx}u_{r3} + \tilde{J}_{xy}u_{s3} + \tilde{J}_{yx}v_{r3} + \tilde{J}_{yy}v_{s3} + \tilde{J}_{zx}w_{r3} + \tilde{J}_{zy}w_{s3})
\end{aligned}$$

Concerning a, b, d, e, g, h we are left with 6 equations for 6 unknowns, which can be cast as follows:

$$\begin{pmatrix} \tilde{J}_{xx} & & & & & \\ & \tilde{J}_{yy} & & & & \\ & & \tilde{J}_{yy} & & & \\ & & & \tilde{J}_{zz} & & \\ \tilde{J}_{xy} & \tilde{J}_{xx} & & & & \\ & \tilde{J}_{xz} & & & & \\ & & \tilde{J}_{yx} & \tilde{J}_{yz} & & \\ & & & & \tilde{J}_{zx} & \tilde{J}_{zy} \end{pmatrix} \begin{pmatrix} a \\ b \\ d \\ e \\ g \\ h \end{pmatrix} = \frac{1}{2} \begin{pmatrix} \tilde{J}_{xz}u_{t4} + \tilde{J}_{yz}v_{t4} + \tilde{J}_{zz}w_{t4} \\ \tilde{J}_{xx}u_{r4} + \tilde{J}_{yx}v_{r4} + \tilde{J}_{zx}w_{r4} \\ \tilde{J}_{xy}u_{s4} + \tilde{J}_{yy}v_{s4} + \tilde{J}_{zy}w_{s4} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

At this stage we can only hope that the system is not ill-posed and that a solution exists. Obviously solving a 6×6 linear system for every marker/particle/etc ... will turn out to be costly. Let's see if we cannot do better.

From the last three equations for C_7, C_8, C_9 we have

$$b = -\frac{\tilde{J}_{xy}}{\tilde{J}_{xz}}a \quad d = -\frac{\tilde{J}_{yz}}{\tilde{J}_{yx}}e \quad h = -\frac{\tilde{J}_{zx}}{\tilde{J}_{zy}}g$$

At this stage we have determined c, f, i entirely and have expressed b, d, h as functions of a, e, g . There only remain three unknowns a, e, g and the equations involving C_4, C_5, C_6 become:

$$\begin{aligned}
0 = C_4 &= \underbrace{\tilde{J}_{xz}u_{t4} + \tilde{J}_{yz}v_{t4} + \tilde{J}_{zz}w_{t4}}_{2T} - \tilde{J}_{xx}2a - \tilde{J}_{yy}2d = 2T - \tilde{J}_{xx}2a + \tilde{J}_{yy}2\frac{\tilde{J}_{yz}}{\tilde{J}_{yx}}e \\
0 = C_5 &= \underbrace{\tilde{J}_{xx}u_{r4} + \tilde{J}_{yx}v_{r4} + \tilde{J}_{zx}w_{r4}}_{2R} - \tilde{J}_{yy}2e - \tilde{J}_{zz}2h = 2R - \tilde{J}_{yy}2e + \tilde{J}_{zz}2\frac{\tilde{J}_{zx}}{\tilde{J}_{zy}}g \\
0 = C_6 &= \underbrace{\tilde{J}_{xy}u_{s4} + \tilde{J}_{yy}v_{s4} + \tilde{J}_{zy}w_{s4}}_{2S} - \tilde{J}_{xx}2b - \tilde{J}_{zz}2g = 2S + \tilde{J}_{xx}2\frac{\tilde{J}_{xy}}{\tilde{J}_{xz}}a - \tilde{J}_{zz}2g \quad (9.120)
\end{aligned}$$

This is much more manageable:

$$\begin{pmatrix} \tilde{J}_{xx} & -\tilde{J}_{yy}\tilde{J}_{yz}/\tilde{J}_{yx} & 0 \\ 0 & \tilde{J}_{yy} & -\tilde{J}_{zz}\tilde{J}_{zx}/\tilde{J}_{zy} \\ -\tilde{J}_{xx}\tilde{J}_{xy}/\tilde{J}_{xz} & 0 & \tilde{J}_{zz} \end{pmatrix} \cdot \begin{pmatrix} a \\ e \\ g \end{pmatrix} = \begin{pmatrix} T \\ R \\ S \end{pmatrix}$$

or,

$$\begin{pmatrix} A_{11} & A_{12} & 0 \\ 0 & A_{22} & A_{23} \\ A_{31} & 0 & A_{33} \end{pmatrix} \cdot \begin{pmatrix} a \\ e \\ g \end{pmatrix} = \begin{pmatrix} T \\ R \\ S \end{pmatrix}$$

The solution is not super-elegant, so I stop here and we might solve the 3×3 system on the fly.

Could there be a case where some off-diagonal \tilde{J} terms are zero and some are not?

Summary

$$\begin{aligned}
u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (as + bt + c)(1 - r)(1 + r) \\
v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (dr + et + f)(1 - s)(1 + s) \\
w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (gr + hs + i)(1 - t)(1 + t) \\
a &= \dots
\end{aligned} \tag{9.121}$$

$$\begin{aligned}
b &= -\frac{\tilde{J}_{xy}}{\tilde{J}_{xz}} a \\
c &= \frac{1}{2\tilde{J}_{xx}} (\tilde{J}_{xy} u_{s2} + \tilde{J}_{xz} u_{t2} + \tilde{J}_{yy} v_{s2} + \tilde{J}_{yz} v_{t2} + \tilde{J}_{zy} w_{s2} + \tilde{J}_{zz} w_{t2}) \\
d &= -\frac{\tilde{J}_{yz}}{\tilde{J}_{yx}} e \\
e &= \dots
\end{aligned} \tag{9.122}$$

$$\begin{aligned}
f &= \frac{1}{2\tilde{J}_{yy}} (\tilde{J}_{xx} u_{r2} + \tilde{J}_{xz} u_{t3} + \tilde{J}_{yx} v_{r2} + \tilde{J}_{yz} v_{t3} + \tilde{J}_{zx} w_{r2} + \tilde{J}_{zz} w_{t3}) \\
g &= \dots \\
h &= -\frac{\tilde{J}_{zx}}{\tilde{J}_{zy}} g \\
i &= \frac{1}{2\tilde{J}_{zz}} (\tilde{J}_{xx} u_{r3} + \tilde{J}_{xy} u_{s3} + \tilde{J}_{yx} v_{r3} + \tilde{J}_{yy} v_{s3} + \tilde{J}_{zx} w_{r3} + \tilde{J}_{zy} w_{s3})
\end{aligned} \tag{9.123}$$

Case of a regular grid made of cuboids In the case of a regular grid with nodes aligned with the x, y, z axis, the 6×6 system above is indefinite as $\tilde{J}_{xy} = \tilde{J}_{yx} = \tilde{J}_{xz} = \dots = 0$. Let us then rewrite the C equations again in this specific case:

$$\begin{aligned}
C_0 &= \tilde{J}_{xx} u_{r1} + \tilde{J}_{yy} v_{s1} + \tilde{J}_{zz} w_{t1} \\
C_1 &= \tilde{J}_{yy} v_{s2} + \tilde{J}_{zz} w_{t2} - \tilde{J}_{xx} 2c \\
C_2 &= \tilde{J}_{xx} u_{r2} + \tilde{J}_{zz} w_{t3} - \tilde{J}_{yy} 2f \\
C_3 &= \tilde{J}_{xx} u_{r3} + \tilde{J}_{yy} v_{s3} - \tilde{J}_{zz} 2i \\
C_4 &= \tilde{J}_{zz} w_{t4} - \tilde{J}_{xx} 2a - \tilde{J}_{yy} 2d \\
C_5 &= \tilde{J}_{xx} u_{r4} - \tilde{J}_{yy} 2e - \tilde{J}_{zz} 2h \\
C_6 &= \tilde{J}_{yy} v_{s4} - \tilde{J}_{xx} 2b - \tilde{J}_{zz} 2g \\
C_7 &= 0 \\
C_8 &= 0 \\
C_9 &= 0
\end{aligned}$$

We see that the condition $C_7 = C_8 = C_9 = 0$ are automatically satisfied. The c, f, i coefficients are obtained as in the general case above. We are left with the equations for C_4, C_5, C_6 (we leave the C_0 equation alone - note that it does not contain any coefficient a, b, c, \dots anymore anyways).

Also, elements are cuboids of size $h_x \times h_y \times h_z$, so that their Jacobian matrix is

$$\mathbf{J} = \begin{pmatrix} h_x/2 & 0 & 0 \\ 0 & h_y/2 & 0 \\ 0 & 0 & h_z/2 \end{pmatrix}$$

and its inverse:

$$\tilde{\mathbf{J}} = \begin{pmatrix} 2/h_x & 0 & 0 \\ 0 & 2/h_y & 0 \\ 0 & 0 & 2/h_z \end{pmatrix}$$

Then the C_4, C_5, C_6 equations become

$$\begin{aligned} 0 = C_4 &= \frac{2}{h_z} w_{t4} - \frac{2}{h_x} 2a - \frac{2}{h_y} 2d \\ 0 = C_5 &= \frac{2}{h_x} u_{r4} - \frac{2}{h_y} 2e - \frac{2}{h_z} 2h \\ 0 = C_6 &= \frac{2}{h_y} v_{s4} - \frac{2}{h_x} 2b - \frac{2}{h_z} 2g \end{aligned} \quad (9.124)$$

This is problematic since we are left with 6 unknowns and 3 equations So we should probably go back to the original definition of

$$\begin{aligned} u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (as + bt + c)(1 - r)(1 + r) \\ v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (dr + et + f)(1 - s)(1 + s) \\ w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (gr + hs + i)(1 - t)(1 + t) \end{aligned}$$

and simply choose 3 of the 6 coefficients a, b, d, e, g, h to be zero ? May be better, as proposed earlier: take $\alpha_r = a = b$, $\alpha_s = d = e$ and $\alpha_t = g = h$? Then, keeping only $\alpha_r, \alpha_s, \alpha_t$:

$$\begin{aligned} u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (\alpha_r(s + t) + c)(1 - r)(1 + r) \\ v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (\alpha_s(r + t) + f)(1 - s)(1 + s) \\ w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (\alpha_t(r + s) + i)(1 - t)(1 + t) \end{aligned}$$

The C_4, C_5, C_6 equations become

$$\begin{aligned} 0 = C_4 &= \frac{2}{h_z} w_{t4} - \frac{2}{h_x} 2\alpha_r - \frac{2}{h_y} 2\alpha_s \\ 0 = C_5 &= \frac{2}{h_x} u_{r4} - \frac{2}{h_y} 2\alpha_s - \frac{2}{h_z} 2\alpha_t \\ 0 = C_6 &= \frac{2}{h_y} v_{s4} - \frac{2}{h_x} 2\alpha_r - \frac{2}{h_z} 2\alpha_t \end{aligned}$$

and we have 3 equations and 3 unknowns:

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 2h_y/h_x & 0 & 2h_y/h_z \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ v_{s4} \end{pmatrix}$$

multiply last line by h_z/h_y :

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ h_z/h_y \cdot 2h_y/h_x & 0 & h_z/h_y \cdot 2h_y/h_z \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ h_z/h_y \cdot v_{s4} \end{pmatrix}$$

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 2h_z/h_x & 0 & 2 \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ h_z/h_y \cdot v_{s4} \end{pmatrix}$$

subtract line 3 from line 1 and put in in line 3:

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 0 & -2h_z/h_y & 2 \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ h_z/h_y \cdot v_{s4} - w_{t4} \end{pmatrix}$$

now multiply 3rd line by h_x/h_z

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 0 & h_x/h_z \cdot -2h_z/h_y & h_x/h_z \cdot 2 \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ h_x/h_z(h_z/h_y \cdot v_{s4} - w_{t4}) \end{pmatrix}$$

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 0 & -2h_x/h_y & 2h_x/h_z \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ h_x/h_y \cdot v_{s4} - h_x/h_z w_{t4} \end{pmatrix}$$

Add line 2 to line 3:

$$\begin{pmatrix} 2h_z/h_x & 2h_z/h_y & 0 \\ 0 & 2h_x/h_y & 2h_x/h_z \\ 0 & 0 & 4h_x/h_z \end{pmatrix} \cdot \begin{pmatrix} \alpha_r \\ \alpha_s \\ \alpha_t \end{pmatrix} = \begin{pmatrix} w_{t4} \\ u_{r4} \\ u_{r4} + h_x/h_y \cdot v_{s4} - h_x/h_z w_{t4} \end{pmatrix}$$

From the third line we obtain:

$$\alpha_t = \frac{1}{4} \frac{h_z}{h_x} \left(u_{r4} + \frac{h_x}{h_y} v_{s4} - \frac{h_x}{h_z} w_{t4} \right) = \frac{1}{4} \left(\frac{h_z}{h_x} u_{r4} + \frac{h_z}{h_y} v_{s4} - w_{t4} \right)$$

Then

$$2 \frac{h_x}{h_y} \alpha_s + 2 \frac{h_x}{h_z} \alpha_t = u_{r4}$$

$$\begin{aligned} \alpha_s &= \frac{h_y}{h_x} \left(\frac{1}{2} u_{r4} - \frac{h_x}{h_z} \alpha_t \right) \\ &= \frac{1}{2} \frac{h_y}{h_x} u_{r4} - \frac{h_y}{h_z} \alpha_t \\ &= \frac{1}{2} \frac{h_y}{h_x} u_{r4} - \frac{h_y}{h_z} \frac{1}{4} \left(\frac{h_z}{h_x} u_{r4} + \frac{h_z}{h_y} v_{s4} - w_{t4} \right) \\ &= \frac{1}{2} \frac{h_y}{h_x} u_{r4} - \frac{1}{4} \left(\frac{h_y}{h_x} u_{r4} + v_{s4} - \frac{h_y}{h_z} w_{t4} \right) \\ &= \frac{1}{4} \frac{h_y}{h_x} u_{r4} - \frac{1}{4} v_{s4} + \frac{1}{4} \frac{h_y}{h_z} w_{t4} \\ &= \frac{1}{4} \left(\frac{h_y}{h_x} u_{r4} - v_{s4} + \frac{h_y}{h_z} w_{t4} \right) \end{aligned}$$

and finally:

$$2 \frac{h_z}{h_x} \alpha_r + 2 \frac{h_z}{h_y} \alpha_s = w_{t4}$$

$$\begin{aligned}
\alpha_r &= \frac{h_x}{h_z} \left(\frac{1}{2} w_{t4} - \frac{h_z}{h_y} \alpha_s \right) \\
&= \frac{1}{2} \frac{h_x}{h_z} w_{t4} - \frac{h_x}{h_y} \alpha_s \\
&= \frac{1}{2} \frac{h_x}{h_z} w_{t4} - \frac{h_x}{h_y} \left(\frac{1}{4} \frac{h_y}{h_x} u_{r4} - \frac{1}{4} v_{s4} + \frac{1}{4} \frac{h_y}{h_z} w_{t4} \right) \\
&= \frac{1}{2} \frac{h_x}{h_z} w_{t4} - \left(\frac{1}{4} u_{r4} - \frac{1}{4} \frac{h_x}{h_y} v_{s4} + \frac{1}{4} \frac{h_x}{h_z} w_{t4} \right) \\
&= -\frac{1}{4} u_{r4} + \frac{1}{4} \frac{h_x}{h_y} v_{s4} + \frac{1}{4} \frac{h_x}{h_z} w_{t4} \\
&= \frac{1}{4} \left(-u_{r4} + \frac{h_x}{h_y} v_{s4} + \frac{h_x}{h_z} w_{t4} \right)
\end{aligned}$$

$$\begin{aligned}
\beta_r &= \frac{1}{2\tilde{J}_{xx}} (\tilde{J}_{yy} v_{s2} + \tilde{J}_{zz} w_{t2}) \\
&= \frac{h_x}{4} \left(\frac{2}{h_y} v_{s2} + \frac{2}{h_z} w_{t2} \right) \\
&= \frac{1}{2} \left(\frac{h_x}{h_y} v_{s2} + \frac{h_x}{h_z} w_{t2} \right) \\
\beta_s &= \frac{1}{2\tilde{J}_{yy}} (\tilde{J}_{xx} u_{r2} + \tilde{J}_{zz} w_{t3}) \\
&= \frac{h_y}{4} \left(\frac{2}{h_x} u_{r2} + \frac{2}{h_z} w_{t3} \right) \\
&= \frac{1}{2} \left(\frac{h_y}{h_x} u_{r2} + \frac{h_y}{h_z} w_{t3} \right) \\
\beta_t &= \frac{1}{2\tilde{J}_{zz}} (\tilde{J}_{xx} u_{r3} + \tilde{J}_{yy} v_{s3}) \\
&= \frac{h_z}{4} \left(\frac{2}{h_x} u_{r3} + \frac{2}{h_y} v_{s3} \right) \\
&= \frac{1}{2} \left(\frac{h_z}{h_x} u_{r3} + \frac{h_z}{h_y} v_{s3} \right)
\end{aligned}$$

To recap,

$$\begin{aligned}
u^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) u_i + (\alpha_r(s + t) + \beta_r)(1 - r)(1 + r) \\
v^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) v_i + (\alpha_s(r + t) + \beta_s)(1 - s)(1 + s) \\
w^h(r, s, t) &= \sum_i \mathcal{N}_i(r, s, t) w_i + (\alpha_t(r + s) + \beta_t)(1 - t)(1 + t) \\
\alpha_r &= \frac{1}{4} \left(-u_{r4} + \frac{h_x}{h_y} v_{s4} + \frac{h_x}{h_z} w_{t4} \right) \\
\alpha_s &= \frac{1}{4} \left(\frac{h_y}{h_x} u_{r4} - v_{s4} + \frac{h_y}{h_z} w_{t4} \right) \\
\alpha_t &= \frac{1}{4} \left(\frac{h_z}{h_x} u_{r4} + \frac{h_z}{h_y} v_{s4} - w_{t4} \right) \\
\beta_r &= \frac{1}{2} \left(\frac{h_x}{h_y} v_{s2} + \frac{h_x}{h_z} w_{t2} \right) \\
\beta_s &= \frac{1}{2} \left(\frac{h_y}{h_x} u_{r2} + \frac{h_y}{h_z} w_{t3} \right) \\
\beta_t &= \frac{1}{2} \left(\frac{h_z}{h_x} u_{r3} + \frac{h_z}{h_y} v_{s3} \right)
\end{aligned}$$

9.31.7 Comparison with Wang *et al.* (2015) for 3D

The following is taken from the supplementary material of Wang *et al.* (2015):

In Eqs (4) and (5), the correction item for each velocity component is calculated based on the other velocity component of the nodes. We extend this approach into 3D situation by adding a quadratic item of x_i to each velocity component:

$$\begin{aligned}
U_i^L(x_1, x_2, x_3) &= (1 - x_1)(1 - x_2) \left[(1 - x_3)U_i^a + x_3U_i^e \right] + \\
&\quad x_1(1 - x_2) \left[(1 - x_3)U_i^b + x_3U_i^f \right] + \\
&\quad (1 - x_1)x_2 \left[(1 - x_3)U_i^c + x_3U_i^g \right] + \\
&\quad x_1x_2 \left[(1 - x_3)U_i^d + x_3U_i^h \right]
\end{aligned} \tag{6}$$

$$\Delta U_1 = x_1(1 - x_1)(C_{10} + x_2C_{12}), \tag{7}$$

$$\Delta U_2 = x_2(1 - x_2)(C_{20} + x_3C_{23}), \tag{8}$$

$$\Delta U_3 = x_3(1 - x_3)(C_{30} + x_1C_{31}). \tag{9}$$

The coefficients of these item (C_{10} , C_{12} , C_{20} , C_{23} , C_{30} , C_{31}) in Eqs (7-9) is to be determined. They should satisfy the following divergence free condition for 3D incompressible flow field:

$$\frac{\partial U_1}{\partial x_1} + \frac{\partial U_2}{\partial x_2} + \frac{\partial U_3}{\partial x_3} = 0. \tag{10}$$

Thus, we take the first derivatives of U_i with respect to x_i based on Eqs (3, 6, 7, 8, 9):

$$\begin{aligned}
\Delta x_1 \frac{\partial U_1}{\partial x_1} &= (1 - x_2) * (1 - x_3) * [U_1^b - U_1^a] + x_2 * (1 - x_3) * [U_1^d - U_1^c] \\
&\quad + (1 - x_2) * x_3 * (U_1^f - U_1^e) + x_2 * x_3 * (U_1^h - U_1^g) \\
&\quad + (1 - 2x_1) * (C_{10} + C_{12} * x_2)
\end{aligned} \tag{11}$$

$$\begin{aligned}\Delta x_2 \frac{\partial U_2}{\partial x_2} = & (1 - x_1) * (1 - x_3) * [U_2^c - U_2^a] + x_1 * (1 - x_3) * [U_2^d - U_2^b] \\ & + (1 - x_1) * x_3 * (U_2^g - U_2^e) + x_1 * x_3 * (U_2^h - U_2^f) \\ & + (1 - 2x_2) * (C_{20} + C_{23} * x_3)\end{aligned}\quad (12)$$

$$\begin{aligned}\Delta x_3 \frac{\partial U_3}{\partial x_3} = & (1 - x_1) * (1 - x_2) * [U_3^e - U_3^a] + x_1 * (1 - x_2) * [U_3^f - U_3^b] \\ & + (1 - x_1) * x_2 * (U_3^g - U_3^c) + x_1 * x_2 * (U_3^h - U_3^d) \\ & + (1 - 2x_3) * (C_{30} + C_{31} * x_1)\end{aligned}\quad (13)$$

Substitute Eqs (11-13) in to Eq (10) and we have an identical equation with the six unknowns ($C_{10}, C_{12}, C_{20}, C_{23}, C_{30}, C_{31}$). As the result, the following items should have their coefficients to be zeros: $1, x_1, x_2, x_3, x_1x_2, x_2x_3, x_3x_1$, which lead to 7 equations:

$$\begin{aligned}1: & \frac{1}{\Delta x_1} [U_1^b - U_1^a + C_{10}] + \frac{1}{\Delta x_2} [U_2^c - U_2^a + C_{20}] + \frac{1}{\Delta x_3} [U_3^e - U_3^a + C_{30}] = 0 \\ x_1: & \frac{1}{\Delta x_1} [-2C_{10}] + \frac{1}{\Delta x_2} [U_2^a - U_2^c + U_2^d - U_2^b] + \frac{1}{\Delta x_3} [U_3^a - U_3^e + U_3^f - U_3^b + C_{31}] = 0 \\ x_2: & \frac{1}{\Delta x_1} [U_1^a - U_1^b + U_1^d - U_1^c + C_{12}] + \frac{1}{\Delta x_2} [-2C_{20}] + \frac{1}{\Delta x_3} [U_3^a - U_3^e + U_3^g - U_3^c] = 0 \\ x_3: & \frac{1}{\Delta x_1} [U_1^a - U_1^b + U_1^f - U_1^e] + \frac{1}{\Delta x_2} [U_2^a - U_2^c + U_2^g - U_2^e + C_{23}] + \frac{1}{\Delta x_3} [-2C_{30}] = 0 \\ x_1x_2: & \frac{1}{\Delta x_1} [-2C_{12}] + \frac{1}{\Delta x_3} [U_3^e - U_3^a + U_3^b - U_3^f + U_3^c - U_3^g + U_3^h - U_3^d] = 0 \\ x_2x_3: & \frac{1}{\Delta x_1} [U_1^b - U_1^a + U_1^d - U_1^c + U_1^f - U_1^e + U_1^h - U_1^g] + \frac{1}{\Delta x_2} [-2C_{23}] = 0 \\ x_3x_1: & \frac{1}{\Delta x_2} [U_2^c - U_2^a + U_2^b - U_2^d + U_2^e - U_2^g + U_2^h - U_2^f] + \frac{1}{\Delta x_3} [-2C_{31}] = 0\end{aligned}$$

From (10), we could also have

$$\begin{aligned}& \frac{(U_1^a + U_1^c + U_1^e + U_1^g - U_1^b - U_1^d - U_1^f - U_1^h)}{\Delta x_1} \\ & + \frac{(U_2^a + U_2^b + U_2^c + U_2^e - U_2^d - U_2^f - U_2^g - U_2^h)}{\Delta x_2} \\ & + \frac{(U_3^a + U_3^b + U_3^c + U_3^e - U_3^d - U_3^f - U_3^g - U_3^h)}{\Delta x_3} = 0\end{aligned}\quad (14)$$

which reduces 7 equations into 6 independent equations. Therefore, the coefficients in Eqs (7-9) as the six unknowns are determined as follows :

$$\begin{aligned}C_{12} &= \frac{\Delta x_1}{2\Delta x_3} [U_3^e - U_3^a + U_3^b - U_3^f + U_3^c - U_3^g + U_3^h - U_3^d], \\ C_{23} &= \frac{\Delta x_2}{2\Delta x_1} [U_1^b - U_1^a + U_1^c - U_1^d + U_1^e - U_1^f + U_1^h - U_1^g], \\ C_{31} &= \frac{\Delta x_3}{2\Delta x_2} [U_2^c - U_2^a + U_2^b - U_2^d + U_2^e - U_2^g + U_2^h - U_2^f], \\ C_{10} &= \frac{\Delta x_1}{2\Delta x_2} [U_2^a - U_2^c + U_2^d - U_2^b] + \frac{\Delta x_1}{2\Delta x_3} [U_3^a - U_3^e + U_3^f - U_3^b + C_{31}], \\ C_{20} &= \frac{\Delta x_2}{2\Delta x_1} [U_1^a - U_1^b + U_1^d - U_1^c + C_{12}] + \frac{\Delta x_2}{2\Delta x_3} [U_3^a - U_3^e + U_3^g - U_3^c], \\ C_{30} &= \frac{\Delta x_3}{2\Delta x_1} [U_1^a - U_1^b + U_1^f - U_1^e] + \frac{\Delta x_3}{2\Delta x_2} [U_2^a - U_2^c + U_2^g - U_2^e + C_{23}].\end{aligned}$$

Taken from the supplementary material of Wang *et al.* (2015).

In my opinion, it is quite unbelievable that such a document was accepted for publication (even as supplementary material). There is not much justification for why their equation 7 only contains x_2 and not also x_3 , same for the other two equations. Rather surprising is also the fact that although equations 3,6,7,8,9 do not contain any $\Delta x_{\{1,2,3\}}$ term then equations 11,12,13 do feature them. Nevertheless, we must make sense of this mess.

Since the authors state that they “transform the rectangular cells into unit squares” I do away with $\Delta x_1 = \Delta x_2 = \Delta x_3 = 1$ altogether. Also, U_1, U_2, U_3 have become U, V, W .

The polynomial representation of U, V, W on the element including the correction factors is

$$\begin{aligned}
U &= (1-x_1)(1-x_2)(1-x_3)U^a + (1-x_1)(1-x_2)x_3U^e \\
&+ x_1(1-x_2)(1-x_3)U^b + x_1(1-x_2)x_3U^f \\
&+ (1-x_1)x_2(1-x_3)U^c + (1-x_1)x_2x_3U^g \\
&+ x_1x_2(1-x_3)U^d + x_1x_2x_3U^h \\
&+ x_1(1-x_1)(C_{10} + x_2C_{12})
\end{aligned} \tag{9.125}$$

$$\begin{aligned}
V &= (1-x_1)(1-x_2)(1-x_3)V^a + (1-x_1)(1-x_2)x_3V^e \\
&+ x_1(1-x_2)(1-x_3)V^b + x_1(1-x_2)x_3V^f \\
&+ (1-x_1)x_2(1-x_3)V^c + (1-x_1)x_2x_3V^g \\
&+ x_1x_2(1-x_3)V^d + x_1x_2x_3V^h \\
&+ x_2(1-x_2)(C_{20} + x_3C_{23})
\end{aligned} \tag{9.126}$$

$$\begin{aligned}
W &= (1-x_1)(1-x_2)(1-x_3)W^a + (1-x_1)(1-x_2)x_3W^e \\
&+ x_1(1-x_2)(1-x_3)W^b + x_1(1-x_2)x_3W^f \\
&+ (1-x_1)x_2(1-x_3)W^c + (1-x_1)x_2x_3W^g \\
&+ x_1x_2(1-x_3)W^d + x_1x_2x_3W^h \\
&+ x_3(1-x_3)(C_{30} + x_1C_{31})
\end{aligned} \tag{9.127}$$

Then

$$\begin{aligned}
\frac{\partial U}{\partial x_1} &= -(1-x_2)(1-x_3)U^a - (1-x_2)x_3U^e + (1-x_2)(1-x_3)U^b + (1-x_2)x_3U^f \\
&+ -x_2(1-x_3)U^c - x_2x_3U^g + x_2(1-x_3)U^d + x_2x_3U^h \\
&+ (1-2x_1)(C_{10} + x_2C_{12})
\end{aligned} \tag{9.128}$$

$$\begin{aligned}
\frac{\partial V}{\partial x_2} &= -(1-x_1)(1-x_3)V^a - (1-x_1)x_3V^e - x_1(1-x_3)V^b - x_1x_3V^f \\
&+ (1-x_1)(1-x_3)V^c + (1-x_1)x_3V^g + x_1(1-x_3)V^d + x_1x_3V^h \\
&+ (1-2x_2)(C_{20} + x_3C_{23})
\end{aligned} \tag{9.129}$$

$$\begin{aligned}
\frac{\partial W}{\partial x_3} &= -(1-x_1)(1-x_2)W^a + (1-x_1)(1-x_2)W^e - x_1(1-x_2)W^b + x_1(1-x_2)W^f \\
&- (1-x_1)x_2W^c + (1-x_1)x_2W^g - x_1x_2W^d + x_1x_2W^h \\
&+ (1-2x_3)(C_{30} + x_1C_{31})
\end{aligned} \tag{9.130}$$

So the velocity divergence can be written

$$\frac{\partial U}{\partial x_1} + \frac{\partial V}{\partial x_2} + \frac{\partial W}{\partial x_3} = A + Bx_1 + Cx_2 + Dx_3 + Ex_1x_2 + Fx_2x_3 + Gx_3x_1 \tag{9.131}$$

with

$$A = -U^a + U^b + C_{10} - V^a + V^c + C_{20} - W^a + W^e + C_{30} \tag{9.132}$$

$$B = -2C_{10} + V^a - V^b - V^c + V^d + W^a - W^e - W^b + W^f + C_{31} \tag{9.133}$$

$$C = U^a - U^b - U^c + U^d + C_{12} - 2C_{20} + W^a - W^e - W^c + W^g \tag{9.134}$$

$$D = U^a - U^e - U^b + U^f + V^a - V^e - V^c + V^g + C_{23} - 2C_{30} \tag{9.135}$$

$$E = -2C_{12} - W^a + W^e + W^b - W^f + W^c - W^g - W^d + W^h \tag{9.136}$$

$$F = -U^a + U^e + U^b - U^f + U^c - U^g - U^d + U^h - 2C_{23} \tag{9.137}$$

$$G = -V^a + V^e + V^b - V^f + V^c - V^g - V^d + V^h - 2C_{31} \tag{9.138}$$

A term by term comparison of these equations shows that these are identical to the 7 equations in the supplementary material between Eq. 13 and Eq. 14 (why are these not numbered in the supplementary material?).

Ideally we wish to have all 7 coefficients A to G equal to zero. This leaves us with 7 equations involving 6 unknowns. In other words the system is over constrained and cannot be solved. However the authors seem to interpret this in the exact opposite way by offering yet one more constraint (Eq. 14) which a) is irrelevant b) is not justified (it is indeed related to Eq. 10 but only by taking all C_{ij} coefficients equal to zero and expressed for $x_1 = x_2 = x_3 = 1/2$). Funny enough, that constraint of Eq. 14 is not used further...

From $E = 0, F = 0, G = 0$ we get:

$$C_{12} = \frac{1}{2}(-W^a + W^e + W^b - W^f + W^c - W^g - W^d + W^h) \quad (9.139)$$

$$C_{23} = \frac{1}{2}(-U^a + U^e + U^b - U^f + U^c - U^g - U^d + U^h) \quad (9.140)$$

$$C_{31} = \frac{1}{2}(-V^a + V^e + V^b - V^f + V^c - V^g - V^d + V^h) \quad (9.141)$$

and from $B = 0, C = 0, D = 0$ we get

$$C_{10} = \frac{1}{2}(V^a - V^b - V^c + V^d + W^a - W^e - W^b + W^f + C_{31}) \quad (9.142)$$

$$C_{20} = \frac{1}{2}(U^a - U^b - U^c + U^d + C_{12} + W^a - W^e - W^c + W^g) \quad (9.143)$$

$$C_{30} = \frac{1}{2}(U^a - U^e - U^b + U^f + V^a - V^e - V^c + V^g + C_{23}) \quad (9.144)$$

These 6 expressions are identical to the ones in the paper. However, let us now turn to A :

$$A = -U^a + U^b + C_{10} - V^a + V^c + C_{20} - W^a + W^e + C_{30} \quad (9.145)$$

$$= -U^a + U^b + \frac{1}{2}(V^a - V^b - V^c + V^d + W^a - W^e - W^b + W^f + C_{31}) \quad (9.146)$$

$$-V^a + V^c + \frac{1}{2}(U^a - U^b - U^c + U^d + C_{12} + W^a - W^e - W^c + W^g) \quad (9.147)$$

$$-W^a + W^e + \frac{1}{2}(U^a - U^e - U^b + U^f + V^a - V^e - V^c + V^g + C_{23}) \quad (9.148)$$

$$= -U^a + U^b + \frac{1}{2}(V^a - V^b - V^c + V^d + W^a - W^e - W^b + W^f) \quad (9.149)$$

$$+ \frac{1}{2} \frac{1}{2}(-V^a + V^e + V^b - V^f + V^c - V^g - V^d + V^h) \quad (9.150)$$

$$-V^a + V^c + \frac{1}{2}(U^a - U^b - U^c + U^d + W^a - W^e - W^c + W^g) \quad (9.151)$$

$$+ \frac{1}{2} \frac{1}{2}(-W^a + W^e + W^b - W^f + W^c - W^g - W^d + W^h) \quad (9.152)$$

$$-W^a + W^e + \frac{1}{2}(U^a - U^e - U^b + U^f + V^a - V^e - V^c + V^g +) \quad (9.153)$$

$$+ \frac{1}{2} \frac{1}{2}(-U^a + U^e + U^b - U^f + U^c - U^g - U^d + U^h) \quad (9.154)$$

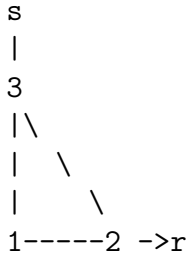
$$\neq 0 \quad (9.155)$$

(easy to prove: for example W^h appears only once)

Once again, we find that the divergence is not identically zero in the element, thereby refuting the statement “Adding these corrections does not improve the order of accuracy of the interpolation (it remains a second-order accurate scheme), but they ensure a divergence-free velocity field over the cell” on page 3 of the article. Their entire paper is based on a false premise.

9.31.8 In 2D with P_1 basis functions - what about triangles?

The reference linear element is:



The basis functions are

$$\begin{aligned}\mathcal{N}_1(r, s) &= 1 - r - s \\ \mathcal{N}_2(r, s) &= r \\ \mathcal{N}_3(r, s) &= s\end{aligned}\tag{9.156}$$

and the velocity vector is $\vec{v} = (u, v)$. Its representation inside the element is

$$\begin{aligned}u^h(r, s) &= \sum_i \mathcal{N}_i(r, s) u_i \\ v^h(r, s) &= \sum_i \mathcal{N}_i(r, s) v_i\end{aligned}$$

and the velocity divergence in the element is given by

$$(\vec{\nabla} \cdot \vec{v})^h = \frac{\partial u^h}{\partial r} + \frac{\partial v^h}{\partial s} = (-u_1 + u_2) + (-v_1 + v_3)$$

which is evidently not zero everywhere in the element. There is however a fundamental difference with regards to quadrilaterals for which the same quantity still contains r and s terms which opens the door to a correction in order to cancel them. In this case, not so much: this term is exactly the one we could not get rid off for quads!

The following consists of a few misguided attempts at designing a CVI scheme for triangles despite the above observation.

approach 1 As we have seen before the CVI approach consists in adding polynomial terms to the expressions of u^h and v^h . In what follows I assume that the additional terms are of the form (I here use only two basis functions per line, similarly to the quadrilateral counterpart):

$$u^h(r, s) = \sum_i N_i(r, s) u_i + f(r, s) r (1 - r - s)\tag{9.157}$$

$$v^h(r, s) = \sum_i N_i(r, s) v_i + g(r, s) s (1 - r - s)\tag{9.158}$$

Note that we thereby ensure that u is continuous across edges, and so is v .

The velocity divergence requirement is then

$$0 = \vec{\nabla} \cdot \vec{v}_h = -u_1 + u_2 + \partial_r f r (1 - r - s) + f(r, s) (1 - 2r - s)\tag{9.159}$$

$$-v_1 + v_3 + \partial_s g s (1 - r - s) + g(r, s) (1 - 2s - r)\tag{9.160}$$

- We start simple and postulate $f(r, s) = a$, $g(r, s) = b$, so then

$$0 = \vec{\nabla} \cdot \vec{v}_h = -u_1 + u_2 + a(1 - 2r - s) - v_1 + v_3 + b(1 - 2s - r) \quad (9.161)$$

$$= (-u_1 + u_2 - v_1 + v_3 + a + b) + (-2a - b)r + (-a - 2b)s \quad (9.162)$$

It is impossible to find a and b such that this expression is zero everywhere inside the element.

- We then turn to linear functions and postulate then $f(r, s) = a + br + cs$, $g(r, s) = d + er + fs$, so

$$\begin{aligned} 0 = \vec{\nabla} \cdot \vec{v}_h &= -u_1 + u_2 + \partial_r f r(1 - r - s) + f(r, s)(1 - 2r - s) \\ &\quad -v_1 + v_3 + \partial_s g s(1 - r - s) + g(r, s)(1 - 2s - r) \\ &= -u_1 + u_2 + br(1 - r - s) + (a + br + cs)(1 - 2r - s) \\ &\quad -v_1 + v_3 + fs(1 - r - s) + (d + er + fs)(1 - 2s - r) \\ &= -u_1 + u_2 - v_1 + v_3 + a + d \\ &\quad + (b - 2a + b - d + e)r \\ &\quad + (f - a + c - 2d + f)s \\ &\quad + (-b - 2b - e)r^2 \\ &\quad + (-f - c - 2f)s^2 \\ &\quad + (-b - f - b - 2c - 2e - f)rs \\ &= -u_1 + u_2 - v_1 + v_3 + a + d \\ &\quad + (2b - 2a - d + e)r \\ &\quad + (2f - a + c - 2d)s \\ &\quad + (-3b - e)r^2 \\ &\quad + (-3f - c)s^2 \\ &\quad + (-2b - 2f - 2c - 2e)rs \end{aligned}$$

Immediately $e = -3b$ and $c = -3f$. Inserting these in the last line yields $-2b - 2f - 2c - 2e = -2b - 2f + 6f + 6b = 4b + 4f = 0$, i.e. $b = -f$. Inserting these in the remaining lines:

$$\begin{aligned} a + d &= u_1 - u_2 + v_1 - v_3 \\ 2b - 2a - d + (-3b) &= 0 \\ 2(-b) - a + (3b) - 2d &= 0 \end{aligned}$$

or,

$$\begin{aligned} a + d &= u_1 - u_2 + v_1 - v_3 \\ -2a - b - d &= 0 \\ -a + b - 2d &= 0 \end{aligned}$$

or,

$$\begin{pmatrix} 1 & 0 & 1 \\ -2 & -1 & -1 \\ -1 & 1 & -2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ d \end{pmatrix} = \begin{pmatrix} u_1 - u_2 + v_1 - v_3 \\ 0 \\ 0 \end{pmatrix}$$

Determinant = $3 \cdot 2 \cdot 1 = 0$. Matrix is singular ... !!

- We now try bilinear functions and postulate $f(r, s) = a + br + cs + hrs$, $g(r, s) = d + er + fs + krs$, so then

$$\begin{aligned}
0 = \vec{\nabla} \cdot \vec{v}_h &= -u_1 + u_2 + \partial_r f r(1 - r - s) + f(r, s)(1 - 2r - s) \\
&\quad -v_1 + v_3 + \partial_s g s(1 - r - s) + g(r, s)(1 - 2s - r) \\
&= -u_1 + u_2 + (b + h s)r(1 - r - s) + (a + b r + c s + h r s)(1 - 2r - s) \\
&\quad -v_1 + v_3 + (f + k r)s(1 - r - s) + (d + e r + f s + k r s)(1 - 2s - r) \\
&= -u_1 + u_2 - v_1 + v_3 + a + d \\
&\quad + (b - 2a + b - d + e)r \\
&\quad + (f - a + c - 2d + f)s \\
&\quad + (-b - 2b - e)r^2 \\
&\quad + (-f - c - 2f)s^2 \\
&\quad + (-b - f - b - 2c - 2e - f + 2h + 2k)rs \\
&\quad + (-h - k - 2k - h)rs^2 \\
&\quad + (-h - k - 2h - k)r^2s \\
&= -u_1 + u_2 - v_1 + v_3 + a + d \\
&\quad + (b - 2a + b - d + e)r \\
&\quad + (f - a + c - 2d + f)s \\
&\quad + (-3b - e)r^2 \\
&\quad + (-3f - c)s^2 \\
&\quad + (-2b - 2f - 2c - 2e + 2h + 2k)rs \\
&\quad + (-2h - 3k)rs^2 \\
&\quad + (-3h - 2k)r^2s
\end{aligned} \tag{9.163}$$

Immediately we see that the last 2 lines yield $k = h = 0$ which are the coefficients in front of the new terms (with regards to linear f and g). This is a dead end too.

I *could* keep adding high order terms but I suspect it is a doomed effort and even if it would work, the cost would be prohibitive.

approach 2 This time I include all three basis functions r , s and $1 - r - s$, not just two. Then

$$u^h(r, s) = \sum_i N_i(r, s)u_i + f(r, s)rs(1 - r - s) \tag{9.164}$$

$$v^h(r, s) = \sum_i N_i(r, s)v_i + g(r, s)rs(1 - r - s) \tag{9.165}$$

$$0 = \vec{\nabla} \cdot \vec{v}^h = -u_1 + u_2 + \partial_r f rs(1 - r - s) + f(r, s)s(1 - 2r - s) \tag{9.166}$$

$$-v_1 + v_3 + \partial_s g rs(1 - r - s) + g(r, s)r(1 - 2s - r) \tag{9.167}$$

We postulate $f(r, s) = a$, $g(r, s) = b$, so then

$$0 = \vec{\nabla} \cdot \vec{v}^h = -u_1 + u_2 + as(1 - 2r - s) - v_1 + v_3 + br(1 - 2s - r) \tag{9.168}$$

$$= (-u_1 + u_2 - v_1 + v_3) + \dots \tag{9.169}$$

This is also a dead end and this will not change with high order terms in f and g . Because of the presence of all three basis functions in the additional terms we see that no coefficient will enter the parenthesis above and therefore it is doomed.

approach 3 We start from

$$\begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial u}{\partial y} \end{pmatrix} = \tilde{\mathbf{J}} \cdot \begin{pmatrix} \frac{\partial u}{\partial r} \\ \frac{\partial u}{\partial s} \end{pmatrix}$$

where $\tilde{\mathbf{J}}$ is the inverse of the Jacobian matrix. We then postulate again

$$\begin{aligned} u(r, s) &= \sum_i \mathcal{N}_i(r, s) u_i + (a_x + b_x r + c_x s + d_x r s + e_x r^2 + f_x s^2) \\ v(r, s) &= \sum_i \mathcal{N}_i(r, s) v_i + (a_y + b_y r + c_y s + d_y r s + e_y r^2 + f_y s^2) \end{aligned}$$

In this case,

$$\frac{\partial u}{\partial r} = \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (b_x + d_x s + 2e_x r) \quad (9.170)$$

$$\frac{\partial u}{\partial s} = \sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + (c_x + d_x r + 2f_x s) \quad (9.171)$$

$$\frac{\partial v}{\partial r} = \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (b_y + d_y s + 2e_y r) \quad (9.172)$$

$$\frac{\partial v}{\partial s} = \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (c_y + d_y r + 2f_y s) \quad (9.173)$$

We have

$$\begin{aligned} \frac{\partial u}{\partial x} &= \tilde{J}_{xx} \frac{\partial u}{\partial r} + \tilde{J}_{xy} \frac{\partial u}{\partial s} \\ &= \tilde{J}_{xx} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i + (b_x + d_x s + 2e_x r) \right) + \tilde{J}_{xy} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial s} u_i + (c_x + d_x r + 2f_x s) \right) \\ &= \tilde{J}_{xx} (-u_{12} + b_x + d_x s + 2e_x r) \\ &\quad + \tilde{J}_{xy} (-u_{13} + c_x + d_x r + 2f_x s) \end{aligned}$$

$$\begin{aligned} \frac{\partial v}{\partial y} &= \tilde{J}_{yx} \frac{\partial v}{\partial r} + \tilde{J}_{yy} \frac{\partial v}{\partial s} \\ &= \tilde{J}_{yx} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial r} v_i + (b_y + d_y s + 2e_y r) \right) + \tilde{J}_{yy} \left(\sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i + (c_y + d_y r + 2f_y s) \right) \\ &= \tilde{J}_{yx} (-u_{12} + b_y + d_y s + 2e_y r) \\ &\quad + \tilde{J}_{yy} (-v_{13} + c_y + d_y r + 2f_y s) \end{aligned}$$

where $u_{ij} = (u_i - u_j)$ and $v_{ij} = (v_i - v_j)$.

Then

$$\begin{aligned} \frac{\partial u^h}{\partial x} + \frac{\partial v^h}{\partial y} &= \tilde{J}_{xx} (-u_{12} + b_x + d_x s + 2e_x r) + \tilde{J}_{xy} (-u_{13} + c_x + d_x r + 2f_x s) \\ &\quad + \tilde{J}_{yx} (-u_{12} + b_y + d_y s + 2e_y r) + \tilde{J}_{yy} (-v_{13} + c_y + d_y r + 2f_y s) \end{aligned}$$

We see that yet again velocity components never multiply r nor s so that no space dependent correction can be designed.

9.31.9 In 2D with Q_2 basis functions - Naive approach

```

03===06===02
||    ||    ||
||    ||    ||
07===08===05
||    ||    ||
||    ||    ||
00===04===01

```

The basis functions are given by:

$$\begin{aligned}
 N_0(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s-1) \\
 N_1(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s-1) \\
 N_2(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s+1) \\
 N_3(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s+1) \\
 N_4(r, s) &= \frac{1}{2}(1-r^2)s(s-1) \\
 N_5(r, s) &= \frac{1}{2}r(r+1)(1-s^2) \\
 N_6(r, s) &= \frac{1}{2}(1-r^2)s(s+1) \\
 N_7(r, s) &= \frac{1}{2}r(r-1)(1-s^2) \\
 N_8(r, s) &= (1-r^2)(1-s^2)
 \end{aligned}$$

and their partial derivatives with respect to the reduced coordinates by

$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s-1) \\
\frac{\partial \mathcal{N}_1}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s-1) \\
\frac{\partial \mathcal{N}_2}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s+1) \\
\frac{\partial \mathcal{N}_3}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s+1) \\
\frac{\partial \mathcal{N}_4}{\partial r} &= (-2r)\frac{1}{2}s(s-1) \\
\frac{\partial \mathcal{N}_5}{\partial r} &= \frac{1}{2}(2r+1)(1-s^2) \\
\frac{\partial \mathcal{N}_6}{\partial r} &= (-2r)\frac{1}{2}s(s+1) \\
\frac{\partial \mathcal{N}_7}{\partial r} &= \frac{1}{2}(2r-1)(1-s^2) \\
\frac{\partial \mathcal{N}_8}{\partial r} &= (-2r)(1-s^2)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{N}_0}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s-1) \\
\frac{\partial \mathcal{N}_1}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s-1) \\
\frac{\partial \mathcal{N}_2}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s+1) \\
\frac{\partial \mathcal{N}_3}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s+1) \\
\frac{\partial \mathcal{N}_4}{\partial s} &= (1-r^2)\frac{1}{2}(2s-1) \\
\frac{\partial \mathcal{N}_5}{\partial s} &= \frac{1}{2}r(r+1)(-2s) \\
\frac{\partial \mathcal{N}_6}{\partial s} &= (1-r^2)\frac{1}{2}(2s+1) \\
\frac{\partial \mathcal{N}_7}{\partial s} &= \frac{1}{2}r(r-1)(-2s) \\
\frac{\partial \mathcal{N}_8}{\partial s} &= (1-r^2)(-2s)
\end{aligned}$$

We then have

$$\begin{aligned}
\frac{\partial u^h}{\partial r} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial r} u_i \\
&= \left[\frac{1}{2}(2r-1)\frac{1}{2}s(s-1) \right] u_0 + \left[\frac{1}{2}(2r+1)\frac{1}{2}s(s-1) \right] u_1 + \left[\frac{1}{2}(2r+1)\frac{1}{2}s(s+1) \right] u_2 + \left[\frac{1}{2}(2r-1)\frac{1}{2}s(s+1) \right] u_3 \\
&\quad + \left[(-2r)\frac{1}{2}s(s-1) \right] u_4 + \left[\frac{1}{2}(2r+1)(1-s^2) \right] u_5 + \left[(-2r)\frac{1}{2}s(s+1) \right] u_6 + \left[\frac{1}{2}(2r-1)(1-s^2) \right] u_7 \\
&\quad + \left[(-2r)(1-s^2) \right] u_8 \\
\frac{\partial v^h}{\partial s} &= \sum_i \frac{\partial \mathcal{N}_i}{\partial s} v_i \\
&= \left[\frac{1}{2}r(r-1)\frac{1}{2}(2s-1) \right] v_0 + \left[\frac{1}{2}r(r+1)\frac{1}{2}(2s-1) \right] v_1 + \left[\frac{1}{2}r(r+1)\frac{1}{2}(2s+1) \right] v_2 + \left[\frac{1}{2}r(r-1)\frac{1}{2}(2s+1) \right] v_3 \\
&\quad + \left[(1-r^2)\frac{1}{2}(2s-1) \right] v_4 + \left[\frac{1}{2}r(r+1)(-2s) \right] v_5 + \left[(1-r^2)\frac{1}{2}(2s+1) \right] v_6 + \left[\frac{1}{2}r(r-1)(-2s) \right] v_7 \\
&\quad + \left[(1-r^2)(-2s) \right] v_8
\end{aligned}$$

or, multiplying each side by 4:

$$\begin{aligned}
4\frac{\partial u^h}{\partial r} &= [(2r-1)s(s-1)] u_0 + [(2r+1)s(s-1)] u_1 + [(2r+1)s(s+1)] u_2 + [(2r-1)s(s+1)] u_3 \\
&\quad + [-4rs(s-1)] u_4 + [2(2r+1)(1-s^2)] u_5 + [-4rs(s+1)] u_6 + [2(2r-1)(1-s^2)] u_7 \\
&\quad + [-8r(1-s^2)] u_8 \\
&= [(2r-1)(s^2-s)] u_0 + [(2r+1)(s^2-s)] u_1 + [(2r+1)(s^2+s)] u_2 + [(2r-1)(s^2+s)] u_3 \\
&\quad + [-4rs(s-1)] u_4 + [2(2r+1)(1-s^2)] u_5 + [-4rs(s+1)] u_6 + [2(2r-1)(1-s^2)] u_7 \\
&\quad + [-8r(1-s^2)] u_8 \\
4\frac{\partial v^h}{\partial s} &= [r(r-1)(2s-1)] v_0 + [r(r+1)(2s-1)] v_1 + [r(r+1)(2s+1)] v_2 + [r(r-1)(2s+1)] v_3 \\
&\quad + [2(1-r^2)(2s-1)] v_4 + [-4rs(r+1)] v_5 + [2(1-r^2)(2s+1)] v_6 + [-4rs(r-1)] v_7 \\
&\quad + [-8s(1-r^2)] v_8 \\
&= [(r^2-r)(2s-1)] v_0 + [(r^2+r)(2s-1)] v_1 + [(r^2+r)(2s+1)] v_2 + [(r^2-r)(2s+1)] v_3 \\
&\quad + [2(1-r^2)(2s-1)] v_4 + [-4rs(r+1)] v_5 + [2(1-r^2)(2s+1)] v_6 + [-4rs(r-1)] v_7 \\
&\quad + [-8s(1-r^2)] v_8
\end{aligned}$$

We then have

$$\begin{aligned}
4(\vec{\nabla} \cdot \vec{v})^h &= 4\frac{\partial u^h}{\partial r} + 4\frac{\partial v^h}{\partial s} \\
&= (2u_5 - 2u_7 - 2v_4 + 2v_6) 1 \\
&\quad + (4u_5 + 4u_7 - 8u_8 + v_0 - v_1 + v_2 - v_3) r \\
&\quad + (u_0 - u_1 + u_2 - u_3 + 4v_4 + 4v_6 - 8v_8) s \\
&\quad + (-2u_0 - 2u_1 + 2u_2 + 2u_3 + 4u_4 - 4u_6 - 2v_0 + 2v_1 + 2v_2 - 2v_3 - 4v_5 + 4v_7) rs \\
&\quad + (-v_0 - v_1 + v_2 + v_3 + 2v_4 - 2v_6) r^2 \\
&\quad + (-u_0 + u_1 + u_2 - u_3 - 2u_5 + 2u_7) s^2 \\
&\quad + (2v_0 + 2v_1 + 2v_2 + 2v_3 - 4v_4 - 4v_5 - 4v_6 - 4v_7 + 8v_8) r^2 s \\
&\quad + (2u_0 + 2u_1 + 2u_2 + 2u_3 - 4u_4 - 4u_5 - 4u_6 - 4u_7 + 8u_8) rs^2
\end{aligned}$$

i.e.

$$(\vec{\nabla} \cdot \vec{v})^h = C_0 + C_1 r + C_2 s + C_3 r s + C_4 r^2 + C_5 s^2 + C_6 r^2 s + C_7 r s^2 \quad (9.174)$$

with

$$\begin{aligned} C_0 &= \frac{1}{4}(2u_5 - 2u_7 - 2v_4 + 2v_6) \\ C_1 &= \frac{1}{4}(4u_5 + 4u_7 - 8u_8 + v_0 - v_1 + v_2 - v_3) \\ C_2 &= \frac{1}{4}(u_0 - u_1 + u_2 - u_3 + 4v_4 + 4v_6 - 8v_8) \\ C_3 &= \frac{1}{4}(-2u_0 - 2u_1 + 2u_2 + 2u_3 + 4u_4 - 4u_6 - 2v_0 + 2v_1 + 2v_2 - 2v_3 - 4v_5 + 4v_7) \\ C_4 &= \frac{1}{4}(-v_0 - v_1 + v_2 + v_3 + 2v_4 - 2v_6) \\ C_5 &= \frac{1}{4}(-u_0 + u_1 + u_2 - u_3 - 2u_5 + 2u_7) \\ C_6 &= \frac{1}{4}(2v_0 + 2v_1 + 2v_2 + 2v_3 - 4v_4 - 4v_5 - 4v_6 - 4v_7 + 8v_8) \\ C_7 &= \frac{1}{4}(2u_0 + 2u_1 + 2u_2 + 2u_3 - 4u_4 - 4u_5 - 4u_6 - 4u_7 + 8u_8) \end{aligned}$$

Looking at C_0 , we see that it is effectively $(u_5 - u_7)/2 + (v_6 - v_4)/2$ which is the divergence expressed in the middle of the element using only the mid-edges velocity components (as in a staggered FD grid).

Looking now at C_4 we can write it

$$C_4 = \frac{1}{4}(-(v_0 - 2v_4 + v_1) + (v_3 - 2v_6 + v_2))$$

Since the reference element is of size 2×2 , then the distance between nodes 0 and 4, and 4 and 1 respectively is $h = 1$. We then recognise

$$\frac{v_0 - 2v_4 + v_1}{h^2} \sim v_4''$$

and likewise

$$\frac{v_3 - 2v_6 + v_2}{h^2} \sim v_6''$$

Can we recognize more FD stencils?

The divergence inside an element is a polynomial, and as before we then need to design a CVI so that we can get rid of the terms containing the C_{1-7} coefficients (while keeping C_0 as low as possible, although we don't have much control over this).

Because we need that the correction term are zero on the edges ($r = \pm 1$ and $s = \pm 1$), we postulate

$$\begin{aligned} \delta u(r, s) &= (1 - r^2)f(r, s) \\ \delta v(r, s) &= (1 - s^2)g(r, s) \end{aligned}$$

with

$$\begin{aligned} f(r, s) &= \sum_{i=0}^m \sum_{j=0}^n a_{ij} r^i s^j = a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{20}r^2 + a_{02}s^2 + a_{12}rs^2 + a_{21}r^2s + a_{22}r^2s^2 + \dots \\ g(r, s) &= \sum_{k=0}^p \sum_{l=0}^q b_{kl} r^k s^l = b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2 + b_{02}s^2 + b_{12}rs^2 + b_{21}r^2s + b_{22}r^2s^2 + \dots \end{aligned}$$

Then the partial derivatives of the velocity corrections are given by:

$$\begin{aligned}
\frac{\partial}{\partial r} \delta u(r, s) &= -2rf(r, s) + (1 - r^2) \frac{\partial f}{\partial r} \\
&= -2r(a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{20}r^2 + a_{02}s^2 + a_{12}rs^2 + a_{21}r^2s + a_{22}r^2s^2 + \dots) \\
&\quad + (1 - r^2)(a_{10} + a_{11}s + 2a_{20}r + a_{12}s^2 + 2a_{21}rs + 2a_{22}rs^2 + \dots) \\
\frac{\partial}{\partial s} \delta v(r, s) &= -2sg(r, s) + (1 - s^2) \frac{\partial g}{\partial s} \\
&= -2s(b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2 + b_{02}s^2 + b_{12}rs^2 + b_{21}r^2s + b_{22}r^2s^2 + \dots) \\
&\quad + (1 - s^2)(b_{01} + b_{11}r + 2b_{02}s + 2b_{12}rs + b_{21}r^2 + 2b_{22}r^2s + \dots)
\end{aligned}$$

We immediately see that a_{12} , a_{22} , a_{20} , a_{21} , b_{02} , b_{12} , b_{21} and b_{22} must be zero, as well as all higher order terms because these $r^\alpha s^\beta$ are not present in (9.174). Then

$$\begin{aligned}
f(r, s) &= a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{02}s^2 \\
g(r, s) &= b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2 \\
\delta u(r, s) &= (1 - r^2)(a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{02}s^2) \\
\delta v(r, s) &= (1 - s^2)(b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2) \\
\frac{\partial}{\partial r} \delta u(r, s) &= -2r(a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{02}s^2) + (1 - r^2)(a_{10} + a_{11}s) \\
\frac{\partial}{\partial s} \delta v(r, s) &= -2s(b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2) + (1 - s^2)(b_{01} + b_{11}r)
\end{aligned}$$

And we have 10 a_{ij} and b_{kl} coefficients to determine. Let us write the corrected velocity divergence:

$$\begin{aligned}
(\vec{\nabla} \cdot \vec{v})_{CVI}^h &= (\vec{\nabla} \cdot \vec{v})^h + \frac{\partial}{\partial r} \delta u(r, s) + \frac{\partial}{\partial s} \delta v(r, s) \\
&= C_0 + C_1r + C_2s + C_3rs + C_4r^2 + C_5s^2 + C_6r^2s + C_7rs^2 \\
&\quad - 2r(a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{02}s^2) + (1 - r^2)(a_{10} + a_{11}s) \\
&\quad - 2s(b_{00} + b_{10}r + b_{01}s + b_{11}rs + b_{20}r^2) + (1 - s^2)(b_{01} + b_{11}r)
\end{aligned}$$

If we want to cancel all first and second-order polynomial terms we need to have

$$\begin{aligned}
C_0 + a_{10} + b_{01} &= 0 \\
C_1 - 2a_{00} + b_{11} &= 0 \\
C_2 + a_{11} - 2b_{00} &= 0 \\
C_3 - 2a_{01} - 2b_{10} &= 0 \\
C_4 - 3a_{10} &= 0
\end{aligned} \tag{9.175}$$

$$\begin{aligned}
C_5 - 3b_{01} &= 0 \\
C_6 - 3a_{11} - 2b_{20} &= 0 \\
C_7 - 2a_{02} - 3b_{11} &= 0
\end{aligned} \tag{9.176}$$

In total there are 10 coefficients and 8 only equations. Interestingly, we see that this time around we also do not really stand a chance to actually have $C_0 + a_{10} + b_{01} = 0$ because a_{10} and b_{01} are actually given by (9.175) and (9.176):

$$\begin{aligned}
a_{10} &= C_4/3 \\
b_{01} &= C_5/3
\end{aligned}$$

I am then left with

$$C_1 - 2a_{00} + b_{11} = 0 \quad (9.177)$$

$$C_2 + a_{11} - 2b_{00} = 0 \quad (9.178)$$

$$C_3 - 2a_{01} - 2b_{10} = 0 \quad (9.179)$$

$$C_6 - 3a_{11} - 2b_{20} = 0 \quad (9.180)$$

$$C_7 - 2a_{02} - 3b_{11} = 0 \quad (9.181)$$

I now have 8 unknowns and 5 equations. Since the system is overconstrained, we could further zero b_{20} and a_{02} (thereby removing quadratic terms altogether from f and g). Then (9.180) and (9.181) give

$$a_{11} = C_6/3$$

$$b_{11} = C_7/3$$

and then (9.177) and (9.178) yield

$$a_{00} = \frac{1}{2}(C_1 + b_{11}) = \frac{1}{2}(C_1 + C_7/3)$$

$$b_{00} = \frac{1}{2}(C_2 + a_{11}) = \frac{1}{2}(C_2 + C_6/3)$$

Finally, we are left with (9.179) and we assume for simplicity $a_{01} = b_{10}$ so

$$a_{01} = b_{10} = C_3/4$$

It must be noted that this is only *one* possible approach.

In the end, choosing the a_{ij} 's and b_{kl} 's coefficients as obtained above will yield

$$\begin{aligned} (\vec{\nabla} \cdot \vec{v})_{CVI}^h &= C_0 + a_{10} + b_{01} \\ &= C_0 + \frac{C_4}{3} + \frac{C_5}{3} \\ &= \frac{1}{4}(2u_5 - 2u_7 - 2v_4 + 2v_6) + \frac{1}{3}\frac{1}{4}(-v_0 - v_1 + v_2 + v_3 + 2v_4 - 2v_6) + \frac{1}{3}\frac{1}{4}(-u_0 + u_1 + u_2 - u_3 - \\ &= \frac{1}{12}(6u_5 - 6u_7 - 6v_4 + 6v_6 - v_0 - v_1 + v_2 + v_3 + 2v_4 - 2v_6 - u_0 + u_1 + u_2 - u_3 - 2u_5 + 2u_7) \\ &= \frac{1}{12}(-(u_0 + u_3 + u_1 + u_2 - u_3 + 4u_5 - 4u_7 - v_0 - v_1 + v_2 + v_3 - 4v_4 + 4v_6)) \end{aligned}$$

Finish? What can we say there? what do we recognise?

Finally:

$$\begin{aligned} \delta u(r, s) &= (1 - r^2)(a_{00} + a_{10}r + a_{01}s + a_{11}rs + a_{02}s^2) \\ &= (1 - r^2)\left(\frac{1}{2}(C_1 + \frac{C_7}{3}) + \frac{C_4}{3}r + \frac{C_3}{4}s + \frac{C_6}{3}rs\right) \\ &= \frac{1}{12}(1 - r^2)(6C_1 + 2C_7 + 4C_4r + 3C_3s + 4C_6rs) \\ \delta v(r, s) &= (1 - s^2)(b_{00} + b_{10}r + b_{01}s + b_{11}rs) \\ &= (1 - s^2)\left(\frac{1}{2}(C_2 + \frac{C_6}{3}) + \frac{C_3}{4}r + \frac{C_5}{3}s + \frac{C_7}{3}rs\right) \\ &= \frac{1}{12}(1 - s^2)(6C_2 + 2C_6 + 3C_3r + 4C_5s + 4C_7rs) \end{aligned}$$

Let us verify one more time:

$$\begin{aligned}
12 \frac{\partial}{\partial r} \delta u(r, s) &= (-2r)(6C_1 + 2C_7 + 4C_4r + 3C_3s + 4C_6rs) + (1 - r^2)(4C_4 + 4C_6s) \\
12 \frac{\partial}{\partial s} \delta v(r, s) &= (-2s)(6C_2 + 2C_6 + 3C_3r + 4C_5s + 4C_7rs) + (1 - s^2)(4C_5 + 4C_7r)
\end{aligned}$$

so that

$$\begin{aligned}
\frac{\partial}{\partial r} \delta u(r, s) + \frac{\partial}{\partial s} \delta v(r, s) &= \frac{1}{12} (-12C_1r - 4C_7r - 8C_4r^2 - 6C_3rs - 8C_6r^2s + 4C_4 + 4C_6s - 4C_4r^2 - 4C_6r^2s) \\
&\quad + \frac{1}{12} (-12C_2s - 4C_6s - 6C_3rs - 8C_5s^2 - 8C_7rs^2 + 4C_5 + 4C_7r - 4C_5s^2 - 4C_7rs^2) \\
&= \frac{1}{12} (4C_4 + 4C_5 - 12C_1r - 12C_2s - 12C_3rs - 12C_4r^2 - 12C_5s^2 - 12C_6r^2s - 12C_7rs^2) \\
&= \frac{1}{3} C_4 + \frac{1}{3} C_5 - C_1r - C_2s - C_3rs - C_4r^2 - C_5s^2 - C_6r^2s - C_7rs^2
\end{aligned}$$

it adds up!

Recap:

$$\begin{aligned}
\delta u(r, s) &= \frac{1}{12}(1 - r^2)(6C_1 + 2C_7 + 4C_4r + 3C_3s + 4C_6rs) \\
\delta v(r, s) &= \frac{1}{12}(1 - s^2)(6C_2 + 2C_6 + 3C_3r + 4C_5s + 4C_7rs) \\
C_0 &= \frac{1}{4}(2u_5 - 2u_7 - 2v_4 + 2v_6) \\
C_1 &= \frac{1}{4}(4u_5 + 4u_7 - 8u_8 + v_0 - v_1 + v_2 - v_3) \\
C_2 &= \frac{1}{4}(u_0 - u_1 + u_2 - u_3 + 4v_4 + 4v_6 - 8v_8) \\
C_3 &= \frac{1}{4}(-2u_0 - 2u_1 + 2u_2 + 2u_3 + 4u_4 - 4u_6 - 2v_0 + 2v_1 + 2v_2 - 2v_3 - 4v_5 + 4v_7) \\
C_4 &= \frac{1}{4}(-v_0 - v_1 + v_2 + v_3 + 2v_4 - 2v_6) \\
C_5 &= \frac{1}{4}(-u_0 + u_1 + u_2 - u_3 - 2u_5 + 2u_7) \\
C_6 &= \frac{1}{4}(2v_0 + 2v_1 + 2v_2 + 2v_3 - 4v_4 - 4v_5 - 4v_6 - 4v_7 + 8v_8) \\
C_7 &= \frac{1}{4}(2u_0 + 2u_1 + 2u_2 + 2u_3 - 4u_4 - 4u_5 - 4u_6 - 4u_7 + 8u_8)
\end{aligned}$$

or

$$\begin{aligned}
\delta u(r, s) &= (1 - r^2)(a_{00} + a_{10}r + a_{01}s + a_{11}rs) \\
\delta v(r, s) &= (1 - s^2)(b_{00} + b_{10}r + b_{01}s + b_{11}rs) \\
a_{00} &= \frac{1}{2}(C_1 + C_7/3) \\
a_{01} &= C_3/4 \\
a_{10} &= C_4/3 \\
a_{11} &= C_6/3 \\
b_{00} &= \frac{1}{2}(C_2 + C_6/3) \\
b_{01} &= C_5/3 \\
b_{10} &= C_3/4 \\
b_{11} &= C_7/3
\end{aligned}$$

9.32 Computing field derivatives -WIP

One often needs the strain rate tensor in geodynamics for two main reasons: 1) it is a quantity which 'helps' with interpreting results 2) it is needed in the non-linear rheology, typically power-law.

Let us assume the scalar nodal field f (e.g., temperature, components of velocity, ...) has been obtained by solving a FE problem. Anywhere within an element the given finite element solution

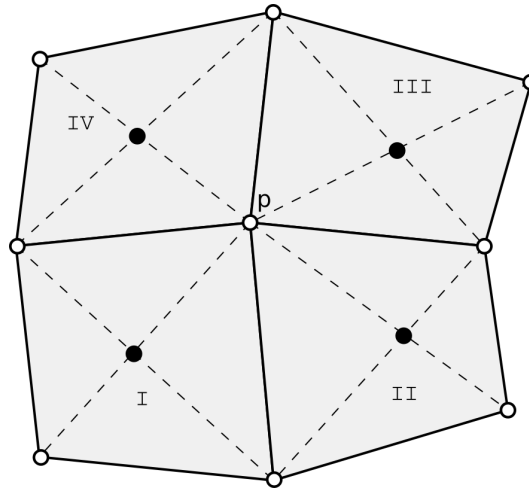
$$f^h(x, y) = \sum_k N_k(x, y) f_k$$

We wish to compute the field $\vec{g}^h = \vec{\nabla} f^h$ on the nodes with the highest accuracy.

For any point inside an element, this problem is trivial and we have

$$\vec{g}^h(x, y) = \sum_{k=1}^m \vec{\nabla} N_k(x, y) f_k \quad (9.182)$$

This method works adequately everywhere inside the element, but since the basis functions derivatives are not uniquely defined on the nodes this problem requires careful attention to arrive at the best result.



Relevant Literature:

- P. Labbé and A. Garon. “A robust implementation of Zienkiewicz and Shu’s local patch recovery method”. In: *Communications in Numerical Methods in Engineering* 11 (1995), pp. 427–434. DOI: 10.1002/cnm.1640110507
- O.C. Zienkiewicz and J.Z. Zhu. “The superconvergent patch recovery and a posteriori error estimates. Part 1: the recovery technique”. In: *Int. J. Num. Meth. Eng.* 33 (1992), pp. 1331–1364, O.C. Zienkiewicz and J.Z. Zhu. “The superconvergent patch recovery and a posteriori error estimates. Part 2: error estimates and adaptativity”. In: *Int. J. Num. Meth. Eng.* 33 (1992), pp. 1365–1382
- P. Ho-Liu, B.H. Hager, and A. Raefsky. “An improved method of Nusselt number calculation”. In: *Geophys. J. R. astr. Soc.* 88 (1987), pp. 205–215. DOI: 10.1111/j.1365-246X.1987.tb01375.x
- P.M. Gresho, R.L. Lee, R.L. Sani, M.K. Maslanik, and B.E. Eaton. “The consistent Galerkin FEM for computing derived boundary quantities in thermal and/or fluid problems”. In: *Int. J. Num. Meth. Fluids* 7 (1987), pp. 371–394

Centroid-to-node method ("method 1")

In this case the gradient is first computed in the 4-element patch to which node p belongs to as shown in Fig. (??). The gradient \vec{g} is computed at the centroid of each element of the patch and averaged out to yield \vec{g}_p .

$$\mathbf{g}_p^h = \frac{1}{4} \sum_e \left(\sum_{k=1}^m (\vec{\nabla} N_k f_k)_{\mathbf{r}=\mathbf{r}_c} \right)_e$$

where \mathbf{r}_c stands for the location of the centroid.

This technique is similar to the one of pressure smoothing showcased in Braun *et al.* (2008) [136] and is mentioned on p. 865 of Gresho & Sani [488] (section 4.2.6 ?).

Although very simple to implement, this approach is not without problem since the algorithm cannot be applied to the nodes on the boundary and an ad-hoc rule must be adopted for these.

Corner-to-node method ("method 2")

For each element of the patch the value of \vec{g} is computed at node p and the obtained values are then averaged out. At the time of writing, this is the technique implemented in ASPECT.

$$\mathbf{g}_p^h = \frac{1}{4} \sum_e \left(\sum_{k=1}^m (\vec{\nabla} N_k f_k)_{\mathbf{r}=\mathbf{r}_p} \right)_e$$

Consistent approach using basis functions ("method 3")

What follows is formulated in 2D Cartesian coordinates for simplicity. Let us start from the function g which is the gradient of the function f in the x -direction:

$$g(x, y) = \frac{\partial f}{\partial x}(x, y)$$

If we left-multiply this equation by a basis function $N_i(x, y)$ and integrate over an element, we arrive at

$$\int_{\Omega_e} N_i(x, y) g(x, y) dV = \int_{\Omega_e} N_i(x, y) \frac{\partial f}{\partial x}(x, y) dV \quad (9.183)$$

The function g is represented inside an element by

$$g^h(x, y) = \sum_j N_j(x, y) g_j = \vec{N} \cdot \vec{g}$$

where $\vec{g} = (g_1, g_2, \dots, g_{m_v})$ is the vector of nodal values for the element and \vec{N} is the vector of basis functions. Likewise we have:

$$\left. \frac{\partial f}{\partial x}(x, y) \right|^h = \frac{\partial \vec{N}}{\partial x} \cdot \vec{f}$$

where $\vec{f} = (f_1, f_2, \dots, f_{m_v})$ is the vector of f nodal values for the element. When we write (9.183) for $i = 1, 2, \dots, m_v$ we arrive at

$$\mathbf{M} \cdot \vec{g} = \mathbf{G}_x \cdot \vec{f}$$

where \mathbf{M} is the elemental mass matrix:

$$\mathbf{M} = \int_{\Omega_e} \vec{N}^T \vec{N} dV$$

and \mathbf{G}_x is the gradient matrix

$$\mathbf{G}_x = \int_{\Omega_e} \vec{N}^T \frac{\partial \vec{N}}{\partial x} dV$$

Both matrices are of size $m_v \times m_v$. After the assembly process we are now ready to solve the global system and obtain the derived nodal value at all nodes. This method is particularly interesting because it can use the existing algorithms already present in any FE which has solved the PDE to obtain f .

The nodal strain rate components are obtained as follows:

$$\dot{\epsilon}_{xx} = \mathbf{M}^{-1} \cdot \mathbf{G}_x \cdot \vec{U} \quad (9.184)$$

$$\dot{\epsilon}_{yy} = \mathbf{M}^{-1} \cdot \mathbf{G}_y \cdot \vec{V} \quad (9.185)$$

$$\dot{\epsilon}_{xy} = \frac{1}{2} \mathbf{M}^{-1} \cdot (\mathbf{G}_x \cdot \vec{V} + \mathbf{G}_y \cdot \vec{U}) \quad (9.186)$$

where \vec{U} is the vector of all nodal u and \vec{V} is the vector of all nodal v .

idea: try lumping M?

make link with consistent pressure recov for q1p0

9.33 Iterative solvers

In what follows, we want to solve the system of linear equations

$$\mathbf{A} \cdot \vec{x} = \vec{b} \quad (9.187)$$

for the vector \vec{x} . We denote the unique solution of this system by \vec{x}^* .

Note that in some cases the the known $n \times n$ matrix \mathbf{A} is symmetric (i.e., $\mathbf{A}^T = \mathbf{A}$), positive-definite (i.e. $\vec{x}^T \cdot \mathbf{A} \cdot \vec{x} > 0$ for all non-zero vectors \vec{x} in \mathbb{R}^n), and real, and \vec{b} is known as well (typically the \mathbb{K} matrix).

 Relevant Literature: Direct and Iterative Solvers [745]

Stationary iterative methods

Basic examples of stationary iterative methods use a splitting of the matrix \mathbf{A} such as

$$\mathbf{A} = \mathbf{D} + \mathbf{L} + \mathbf{U}$$

where \mathbf{D} is only the diagonal part of \mathbf{A} , \mathbf{L} is the strict lower triangular part of \mathbf{A} and \mathbf{U} is the strict upper triangular part of \mathbf{A} .

For instance:

$$\mathbf{A} = \begin{pmatrix} 1 & 5 & 8 \\ 6 & 4 & 2 \\ -1 & 7 & 5 \end{pmatrix} \quad \Rightarrow \quad \mathbf{D} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 5 \end{pmatrix} \quad \mathbf{L} = \begin{pmatrix} 0 & 0 & 0 \\ 6 & 0 & 0 \\ -1 & 7 & 0 \end{pmatrix} \quad \mathbf{U} = \begin{pmatrix} 0 & 5 & 8 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{pmatrix}$$

- Jacobi method⁴⁰: The solution is then obtained iteratively via

$$\mathbf{D} \cdot \vec{x}^{k+1} = -(\mathbf{L} + \mathbf{U}) \cdot \vec{x}^k + \vec{b} \quad k = 0, 1, \dots \quad (9.188)$$

where \vec{x}^k is the k -th approximation or iteration of \vec{x} and \vec{x}^0 is the initial guess (often taken to be zero). A sufficient (but not necessary) condition for the method to converge is that the matrix \mathbf{A} is strictly or irreducibly diagonally dominant.

- Gauss-Seidel method⁴¹: It is defined by the iteration

$$\mathbf{L}_* \cdot \vec{x}^{k+1} = -\mathbf{U} \cdot \vec{x}^k + \vec{b} \quad k = 0, 1, \dots \quad (9.189)$$

where $\mathbf{L}_* = \mathbf{L} + \mathbf{U}$.

There is actually a way to make the computation of $\vec{x}^{(k+1)}$ which uses the elements of $\vec{x}^{(k+1)}$ that have already been computed, and only the elements of $\vec{x}^{(k)}$ that have not been computed in the $k + 1$ iteration. This means that, unlike the Jacobi method, only one storage vector is required as elements can be overwritten as they are computed, which can be advantageous for very large problems.

Note that Gauss-Seidel is the same as SOR (successive over-relaxation) with $\omega = 1$.

- Successive over-relaxation method (SOR):

$$(\mathbf{D} + \omega \mathbf{L}) \cdot \vec{x}^{k+1} = -(\omega \mathbf{U} + (\omega - 1) \mathbf{D}) \cdot \vec{x}^k + \omega \vec{b} \quad k = 0, 1, \dots \quad (9.190)$$

⁴⁰https://en.wikipedia.org/wiki/Jacobi_method

⁴¹https://en.wikipedia.org/wiki/Gauss-Seidel_method

- Symmetric successive over-relaxation (SSOR)⁴²: The version of SOR for symmetric matrices \mathbf{A} , in which $\mathbf{U} = \mathbf{L}^T$ is given by the recursion

$$\vec{x}^{k+1} = \vec{x}^k - \gamma^k \mathbf{P}^{-1}(\mathbf{A} \cdot \vec{x}^{(k)} - \vec{b}) \quad k = 0, 1, \dots \quad (9.191)$$

with

$$\mathbf{P} = \left(\frac{\mathbf{D}}{\omega} + \mathbf{L} \right) \frac{\omega}{2 - \omega} \mathbf{D}^{-1} \cdot \left(\frac{\mathbf{D}}{\omega} + \mathbf{L} \right)$$

with $0 < \omega < 2$.

All these methods can be cast in a more general framework⁴³: The basic iterative methods work by splitting the matrix \mathbf{A} into \mathbf{M} - \mathbf{N} and here the matrix \mathbf{M} should be easily invertible. The iterative methods are now defined as

$$\mathbf{M} \cdot \vec{x}^{k+1} = \mathbf{N} \cdot \vec{x}^k + \vec{b} \quad (9.192)$$

with

- Richardson method: $\mathbf{M} = \frac{1}{\omega} \mathbf{I}$
- Jacobi method: $\mathbf{M} = \mathbf{D}$
- Damped Jacobi method: $\mathbf{M} = \frac{1}{\omega} \mathbf{D}$
- Gauss-Seidel method: $\mathbf{M} = \mathbf{D} + \mathbf{L}$
- Successive over-relaxation method: $\mathbf{M} = \frac{\mathbf{D}}{\omega} + \mathbf{L}$
- Symmetric successive over-relaxation: $\mathbf{M} = \left(\frac{\mathbf{D}}{\omega} + \mathbf{L} \right) \frac{\omega}{2 - \omega} \mathbf{D}^{-1} \cdot \left(\frac{\mathbf{D}}{\omega} + \mathbf{L} \right)$

and $\mathbf{N} = \mathbf{M} - \mathbf{A}$.

Krylov subspace methods

- Conjugate Gradient⁴⁴

It was first proposed by Hestenes and Stiefel in 1952 [567]. The method solves an SPD system $\mathbf{A} \cdot \vec{x} = \vec{b}$ of size n . In theory (i.e. exact arithmetic) it does so in n iterations. Each iteration requires a few inner products in \mathbb{R}^n and one matrix-vector multiplication. With roundoff error, CG can work poorly (or not at all), but for some \mathbf{A} (and \vec{b}), can get good approximate solution in $\ll n$ iterations.

As an iterative method, the conjugate gradient method monotonically (in the energy norm) improves approximations \vec{x}_k to the exact solution and may reach the required tolerance after a relatively small (compared to the problem size) number of iterations. The improvement is typically linear and its speed is determined by the condition number $\kappa(\mathbf{A})$ of the system matrix \mathbf{A} : the larger $\kappa(\mathbf{A})$ is, the slower the improvement.

If $\kappa(\mathbf{A})$ is large, preconditioning is commonly used to replace the original system $\mathbf{A} \cdot \vec{x} - \vec{b} = \vec{0}$ with $\mathbf{M}^{-1} \cdot (\mathbf{A} \cdot \vec{x} - \vec{b}) = \vec{0}$ such that $\kappa(\mathbf{M}^{-1} \cdot \mathbf{A})$ is smaller than $\kappa(\mathbf{A})$.

The resulting method is called the Preconditioned Conjugate Gradient method (PCG). An extreme case of preconditioner is $\mathbf{M} = \mathbf{A}^{-1}$ but it is a silly case since applying the preconditioner is as difficult as solving the system in the first place. In the end the goal is to find a matrix \mathbf{M} that is cheap to multiply, and is an approximate inverse of \mathbf{A} (or at least has a more clustered spectrum than \mathbf{A}).

⁴²https://en.wikipedia.org/wiki/Symmetric_successive_over-relaxation

⁴³https://en.wikipedia.org/wiki/Iterative_method

⁴⁴https://en.wikipedia.org/wiki/Conjugate_gradient_method

| | |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> $\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$ if \mathbf{r}_0 is sufficiently small, then return \mathbf{x}_0 as the result $\mathbf{p}_0 := \mathbf{r}_0$ $k := 0$ repeat $\alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$ $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$ $\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$ if \mathbf{r}_{k+1} is sufficiently small, then exit loop $\beta_k := \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}$ $\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$ $k := k + 1$ end repeat return \mathbf{x}_{k+1} as the result </pre> | <pre> $\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$ $\mathbf{z}_0 := \mathbf{M}^{-1} \mathbf{r}_0$ $\mathbf{p}_0 := \mathbf{z}_0$ $k := 0$ repeat $\alpha_k := \frac{\mathbf{r}_k^T \mathbf{z}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$ $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$ $\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$ if \mathbf{r}_{k+1} is sufficiently small then exit loop end if $\mathbf{z}_{k+1} := \mathbf{M}^{-1} \mathbf{r}_{k+1}$ $\beta_k := \frac{\mathbf{z}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{z}_k^T \mathbf{r}_k}$ $\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$ $k := k + 1$ end repeat The result is \mathbf{x}_{k+1} </pre> |
| <p>Let's put it all together into one piece now. The method of Conjugate Gradients is:</p> $d_{(0)} = r_{(0)} = b - Ax_{(0)},$ $\alpha_{(i)} = \frac{r_{(i)}^T r_{(i)}}{d_{(i)}^T A d_{(i)}} \quad (\text{by Equations 32 and 42}),$ $x_{(i+1)} = x_{(i)} + \alpha_{(i)} d_{(i)},$ $r_{(i+1)} = r_{(i)} - \alpha_{(i)} A d_{(i)},$ $\beta_{(i+1)} = \frac{r_{(i+1)}^T r_{(i+1)}}{r_{(i)}^T r_{(i)}},$ $d_{(i+1)} = r_{(i+1)} + \beta_{(i+1)} d_{(i)}.$ | |

Top: algorithms as obtained from Wikipedia (Left: CG; Right: PCG); Bottom: algorithm from Shewchuk (1994) [1156].

Also available on Wikipedia is a (naive) MATLAB implementation of the CG algorithm:


```

function x = conjgrad(A, b, x)
    r = b - A * x;
    p = r;
    rsold = r' * r;

    for i = 1:length(b)
        Ap = A * p;
        alpha = rsold / (p' * Ap);
        x = x + alpha * p;
        r = r - alpha * Ap;
        rsnew = r' * r;
        if sqrt(rsnew) < 1e-10
            break
        end
        p = r + (rsnew / rsold) * p;
        rsold = rsnew;
    end
end

```

We see that its implementation is actually rather simple and straightforward!

 **Relevant Literature:** Shewchuk, An Introduction to the Conjugate Gradient Method Without the Agonizing Pain [1156]. CG using mpi [1165]. Een kwart eeuw iteratieve methoden [1332].

The CG and PCG algorithms are used in Section 7.11.3. It is implemented in [STONE](#) 15,16,82.

- Biconjugate Gradient method ⁴⁵
- Biconjugate Gradient stabilised method ⁴⁶

⁴⁵https://en.wikipedia.org/wiki/Biconjugate_gradient_method

⁴⁶https://en.wikipedia.org/wiki/Biconjugate_gradient_stabilized_method

- MINRES: For iterative solution of symmetric systems $Ax = b$, the conjugate gradient method (CG) is commonly used when A is positive definite, while the minimum residual method (MINRES) is typically reserved for indefinite systems.
- Generalized minimal residual method (GMRES) ⁴⁷

⁴⁷https://en.wikipedia.org/wiki/Generalized_minimal_residual_method

9.34 Weak seeds in extension modelling

weakseeds.tex

This section was mostly written by I. van Zelst with some input by S. Buiter.

Numerical models that investigate dynamics of the lithosphere and upper mantle always start from an initial geometry with a set of prescribed mechanical and thermal conditions. This initial setup is usually a more-or-less standard representation of the lithosphere and asthenosphere, as defined from compilations of geological and geophysical observations and laboratory measurements. Deformation is driven by internal buoyancy forces and/or velocity or stress boundary conditions. However, unless an intrinsically unstable setup is defined or boundary conditions are discontinuous, deformation may take long model time to localize (up to millions of years). This is because these models need to build up numerical disturbance to create starting points for the deformation. In such models, deformation may in the first stages be accommodated by pure shear extension or shortening [1024, 901].

To avoid this long starting phase and, in addition, exert some control over the initial location of deformation (preferably away from the boundaries), modelers use different approaches to initiate and localize deformation.

One manner to localize deformation is by discontinuous boundary conditions, such as the so-called S-point velocity discontinuity at the bottom of the system (or the tip of a basal sheet) which is used in both numerical [143, 368, 1360, 61, 161, 1261, 137] and analogue studies [161, 872]. These models are usually on the scale of the (upper-) crust. S-point models are less flexible than upper-mantle scale models as they do not include feedback relations between deformation and the basal velocity field. Models of extension of continental lithosphere often use 'seeds' to initiate extension. Such seeds are usually small regions that are weaker than the surrounding crust and lithosphere. The use of seeds can be justified by considering the fact that in nature continental lithosphere is hardly ever (if at all) homogeneous in composition and stratification. In addition, extension often occurs in regions of former convergence, such as the opening of the North Atlantic Ocean that largely followed the old sutures of the Iapetus and Rheic Oceans [1364]. Analogues for numerical seeds can therefore be found in inherited faults, inherited crustal thickness changes, and/or plumes impacting the lithosphere. However, this immediately points out a problem with single-seed models as orogenic inheritance and mantle upwellings may be expected to occur over larger areas than a seed of some hundreds of meters to a few kilometers in width and height.

A literature survey shows that seeds in previous numerical studies differ in shape, size, orientation, mechanical and thermal properties, and depth in the models. Three types of weak seeds can be identified that have been used in previous models of (continental) extension:

Seeding through thermal effects A weak region can be achieved by a temperature anomaly in the crust or lithosphere [177], which is created by directly imposing a temperature difference, by assigning high radiogenic heat production, or modeling a thermal upwelling in the mantle below. The elevated temperature reduces viscosity values for models with a temperature-dependent viscosity. An advantage of using an imposed temperature anomaly is that it will dissipate with time, thus reducing the impact on later model stages [531]. Examples of thermal anomalies used to initiate extension are an elevated temperature at the base of the crust [531], an elevated temperature at the base of the lithosphere (100° in [184], up to 200° in [156]), a temperature anomaly imposed from the base of the lithosphere to the middle crust [232], and a 10mW m⁻² perturbation in basal heat flow [416]. In [184] the rifting is initiated by means of a thermal perturbation placed at the bottom of the mantle lithosphere with a maximum temperature T_2 exponentially which decays from the center to T_1 on the left and T_3 on the right.

Seeding by mechanical inhomogeneity A seed may be composed of a material with a lower rheological strength than the surroundings. A weak seed may, for example, have a lower imposed viscosity [764, 679, 878], a lower value for angle of internal friction [1023, 682, 1258, 483, 227], a lower value for cohesion [10], or a lower value for density [1269].

The seed may also be assigned different material properties, as, for example, a Von Mises seed in a frictional plastic material [614]. A frequently used approach is to assume that a region has already accumulated strain, leading to strain-weakening [754, 615, 995, 10, 9, 711, 11]. Previous studies have used a variety of shapes and sizes for weak seeds. Examples are square seeds, fault-shaped weak inclusions, and rectangular seeds with different aspect ratios: [615] use a 6×3 km seed, while the weak seed of [614] has a size of 12×10 km.

Instead on confining the seed to a geometrically simple region, randomly distributed seeds have also been used [1258], [1262], albeit for compression.

Seeding through geometrical discontinuity A seed is created by an abrupt variation in the thickness of the crust and/or lithosphere. A locally thinned crust could be thought to be caused by a previous rifting phase, whereas a thicker crust could represent preceding mountain building. Such crustal thickness variations effect not only mechanical strength, but may also impose a thermal anomaly. Burg & Schmalholz [178] implemented a Gaussian shaped mohorovičić discontinuity of 250m height as a representation of the weak zone resulting in a slightly thinner crust. A step change in crustal thickness alters the symmetry of the domain. Chenin *et al.* [230] implement a sinusoidal perturbation of the Moho.

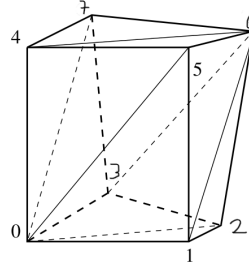
Only few studies have investigated how different methods of implementing a weak zone can affect the results of a model. [356] found that a single seed produces a symmetric narrow rift, an initial shear zone tends to produce an asymmetric rift, and multiple seeds promote a wide rift. Note however that this behavior will be affected by rheological stratification, as not all systems can evolve in a wide rift mode [615, 163]. [356] found that a seed needs to be 10 times weaker than the surrounding material in order to localize strain. To initiate shear bands with a Coulomb dip angle ($45 \pm \phi/2$, where ϕ is the angle of internal friction), seeds need to be well resolved (5-10 elements, [679]).

This variety in the shape, size, orientation, mechanical and thermal properties, and depth of the seed(s) begs the question if these different approaches to initiate extension could have an effect on model evolution? As such variations might not be removed by subsequent deformation stages, the initiation effects could propagate into later model evolution. In addition, weak seeds introduce a weakness into the extensional system that may potentially be long-lasting. For instance, seeds with a weakness defined by material or strain-weakened properties stay in the model and may control deformation also in later model stages. The heat associated with thermal seeds will diffuse away, but additional heat has been introduced into the initial system and the setup will therefore differ from models with mechanically weak seeds.

9.35 Computing the volume of a hexahedron

volume_hexahedron.tex

What follows is based on the report "Efficient Computation of Volume of Hexahedral Cells" by J. Grandy (1997) [481]. We assume the following internal numbering of the hexahedron, which is different than the one in the paper:



Modified from [481]

in ./images/hexahedron/

If the hexahedron is such that some or all the opposite faces are planes parallel to each other than the volume can be arrived at very simply⁴⁸.

The real catch here is that the four nodes which make up a face are not necessarily co-planar!

The volume is then computed as follows

$$\begin{aligned}
 V = & [(\vec{r}_6 - \vec{r}_1) + (\vec{r}_7 - \vec{r}_0), (\vec{r}_6 - \vec{r}_3), (\vec{r}_2 - \vec{r}_0)] \\
 & + [(\vec{r}_7 - \vec{r}_0), (\vec{r}_6 - \vec{r}_3) + (\vec{r}_5 - \vec{r}_0), (\vec{r}_6 - \vec{r}_4)] \\
 & + [(\vec{r}_6 - \vec{r}_1), (\vec{r}_5 - \vec{r}_0), (\vec{r}_6 - \vec{r}_4) + (\vec{r}_2 - \vec{r}_0)] \\
 & / 12
 \end{aligned} \tag{9.193}$$

where $[\cdot]$ is the triple product:

$$[\vec{A}, \vec{B}, \vec{C}] = \begin{vmatrix} A_x & B_x & C_x \\ A_y & B_y & C_y \\ A_z & B_z & C_z \end{vmatrix}$$

It is implemented and used in Stone 98. The code is shown hereunder:

```
def hexahedron_volume (x,y,z):
    val = ( triple_product ( x[6]-x[1]+x[7]-x[0] , y[6]-y[1]+y[7]-y[0] , z[6]-z[1]+z[7]-z[0] , \
                                x[6]-x[3] , y[6]-y[3] , z[6]-z[3] ,
                                \
                                x[2]-x[0] , y[2]-y[0] , z[2]-z[0]
                            )\
        + triple_product ( x[7]-x[0] , y[7]-y[0] , z[7]-z[0] ,
                                x[6]-x[3]+x[5]-x[0] , y[6]-y[3]+y[5]-y[0] , z[6]-z[3]+z[5]-z[0] , \
                                x[6]-x[4] , y[6]-y[4] , z[6]-z[4]
                            )\
        + triple_product ( x[6]-x[1] , y[6]-y[1] , z[6]-z[1] ,
                                x[5]-x[0] , y[5]-y[0] , z[5]-z[0] ,
                                x[6]-x[4]+x[2]-x[0] , y[6]-y[4]+y[2]-y[0] , z[6]-z[4]+z[2]-z[0]
                            ) )/12.
    return val
```

⁴⁸<https://en.wikipedia.org/wiki/Cuboid>

with

```
def triple_product (Ax,Ay,Az,Bx,By,Bz,Cx,Cy,Cz) :  
    val = Ax * ( By * Cz - Bz * Cy ) \  
        - Ay * ( Bx * Cz - Bz * Cx ) \  
        + Az * ( Bx * Cy - By * Cx )  
    return val
```

 Relevant Literature: [349]

9.36 Bandwidth reduction, matrix reordering

The need for reordering

The profile (or envelope) of a symmetric matrix determines how close its non-zero elements are to the diagonal:

$$\text{profile} = \sum_{i=1}^n (i - \min(ne(i)))$$

The bandwidth is the largest deviation:

$$\text{bandwidth} = \max(i - \min(ne(i)))$$

The cost for a band cholesky factorisation with bandwidth p is $n(p^2 + 3p)$ flops assuming $n \gg p$. (source?)

In conclusion: reducing bandwidth means a factor solve if Cholesky factorisation is used.

A simple example

Let us consider a structurally symmetric matrix \mathbf{M} . We wish to reduce its bandwidth by permuting rows and columns such as to move all the nonzero elements of \mathbf{M} in a band as close as possible to the diagonal. We then talk about *Bandwidth Reduction*.

We know that the solution a linear system remains unchanged if lines or columns of the matrix (and corresponding rhs) are permuted.

For example⁴⁹, let us consider the 5×5 matrix \mathbf{M} :

$$\mathbf{M} = \begin{pmatrix} 1 & . & . & . & 1 & . & . & . \\ . & 1 & 1 & . & . & 1 & . & 1 \\ . & 1 & 1 & . & 1 & . & . & . \\ . & . & . & 1 & . & . & 1 & . \\ 1 & . & 1 & . & 1 & . & . & . \\ . & 1 & . & . & . & 1 & . & 1 \\ . & . & . & 1 & . & . & 1 & . \\ . & 1 & . & . & . & 1 & . & 1 \end{pmatrix}$$

Simply through row and column permutations it can be rewritten

$$\mathbf{M}' = \begin{pmatrix} 1 & 1 & . & . & . & . & . & . \\ 1 & 1 & . & . & . & . & . & . \\ . & . & 1 & 1 & 1 & . & . & . \\ . & . & 1 & 1 & 1 & . & . & . \\ . & . & 1 & 1 & 1 & 1 & . & . \\ . & . & . & . & 1 & 1 & 1 & . \\ . & . & . & . & . & 1 & 1 & 1 \\ . & . & . & . & . & . & 1 & 1 \end{pmatrix}$$

The different existing algorithms

- The simplest bandwidth reduction method is the Cuthill-McKee algorithm (1969) [297]
- Reverse Cuthill-McKee algorithm (1976) [460]

⁴⁹Taken from <http://ciprian-zavoianu.blogspot.com/2009/01/project-bandwidth-reduction.html>

- The Gibbs-Poole-Stockmeyer and Gibbs-King algorithm is an alternative, often superior, profile reduction method (1976) [461]
- Sloan algorithm [1175, 1174]

Implementation in python

The documentation for the reverse Cuthill-McKee algorithm is available online⁵⁰, but it is quite useless as to how the result of the function call should be used. I therefore provide here a small python program in /images/reordering which builds matrix M and returns M' .

This header is necessary:

```
from scipy.sparse import csr_matrix
from scipy.sparse.csgraph import reverse_cuthill_mckee
```

After the matrix is filled, the reverse Cuthill-McKee algorithm is used to compute the permutation order of the rows and columns:

```
perm = reverse_cuthill_mckee(sparse_matrix, symmetric_mode=True)
```

The result is an array of size n which indicates the new order of rows and columns:

```
[6 3 7 5 1 2 4 0]
```

All that we have to do then is to use this array to rebuild the matrix:

```
sparse_matrix=sparse_matrix[np.ix_(perm,perm)]
```

In case a right hand side vector exists, it must be reordered too in order to match the matrix:

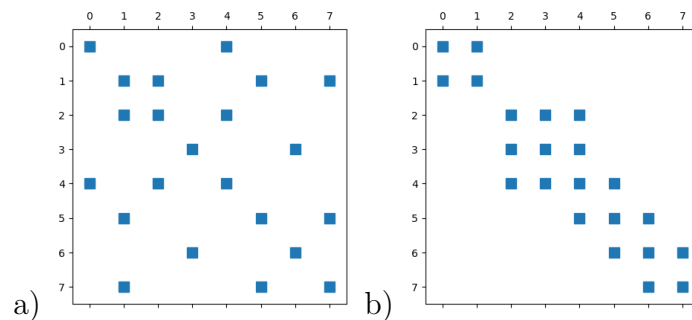
```
rhs=rhs[np.ix_(perm)]
```

Assuming a solution of the reordered matrix and rhs has been obtained, it must be reordered back. We therefore create the inverse permutation array:

```
perm_inv=np.empty(n,dtype=np.int32)
for i in range(0,n):
    perm_inv[perm[i]]=i
```

and use it as follows:

```
sol=sol[np.ix_(perm_inv)]
```



⁵⁰https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csgraph.reverse_cuthill_mckee.html

9.37 Scaling between dimensioned and dimensionless quantities

All quantities with an r subscript denote real quantities, while quantities with an m subscript denote dimensionless model quantities, i.e. values that are used as inputs to the code.

Let us define four coefficients K_L , K_T , K_M and K_θ as follows:

$$K_L = \frac{L_m}{L_r} \quad K_T = \frac{T_m}{T_r} \quad K_M = \frac{M_m}{M_r} \quad K_\theta = \frac{\theta_m}{\theta_r} \quad (9.194)$$

where $L_{\{m,r\}}$, $T_{\{m,r\}}$, $M_{\{m,r\}}$ and $\theta_{\{m,r\}}$ are respectively lengths, times, masses and temperatures, and the coefficients K bear the following dimensions:

$$K_M \rightarrow kg^{-1}, \text{ or } M^{-1}$$

$$K_L \rightarrow m^{-1}, \text{ or } L^{-1}$$

$$K_T \rightarrow s^{-1}, \text{ or } T^{-1}$$

$$K_\theta \rightarrow K^{-1}, \text{ or } \theta^{-1}$$

Velocity The dimensions of a velocity are: $m.s^{-1}$, or LT^{-1} , so that

$$\mathbf{v}_m = \frac{L_m}{T_m} = \frac{K_L L_r}{K_T T_r} = K_L K_T^{-1} \frac{L_r}{T_r} = K_L K_T^{-1} \mathbf{v}_r \quad (9.195)$$

Viscosity The dimensions of a viscosity are $Pa.s$, or $ML^{-1}T^{-1}$, so that

$$\mu_m = \frac{M_m}{L_m T_m} = \frac{K_M}{K_L K_T} \frac{M_r}{L_r T_r} = \frac{K_M}{K_L K_T} \mu_r \quad (9.196)$$

Cohesion The dimensions of a cohesion are those of a stress, i.e. Pa , or $ML^{-1}T^{-2}$, so that

$$c_m = \frac{M_m}{L_m T_m^2} = \frac{K_M}{K_L K_T^2} \frac{M_r}{L_r T_r^2} = \frac{K_M}{K_L K_T^2} c_r$$

Density The dimensions of a density are $kg.m^{-3}$, or ML^{-3} , so that

$$\rho_m = \frac{M_m}{L_m^3} = \frac{K_M}{K_L^3} \frac{M_r}{L_r^3} = \frac{K_M}{K_L^3} \rho_r$$

Gravity The dimensions of the gravity acceleration are $m.s^{-2}$, or LT^{-2} , so that

$$\mathbf{g}_m = \frac{L_m}{T_m^2} = \frac{K_L}{K_T^2} \frac{L_r}{T_r^2} = \frac{K_L}{K_T^2} \mathbf{g}_r$$

Body force The dimensions of body forces are those of $\rho\mathbf{g}$, i.e. $kg.m^{-2}.s^{-2}$, or $ML^{-2}T^{-2}$, so that

$$(\rho\mathbf{g})_m = \frac{M_m}{L_m^2 T_m^2} = \frac{K_M}{K_L^2 K_T^2} \frac{M_r}{L_r^2 T_r^2} = \frac{K_M}{K_L^2 K_T^2} (\rho\mathbf{g})_r$$

Thermal expansion The dimension of thermal expansion is $^{\circ}C^{-1}$, so that

$$\alpha_m = \frac{1}{\theta_m} = \frac{1}{K_\theta} \frac{1}{\theta_r} = \frac{1}{K_\theta} \alpha_r$$

Thermal conductivity The dimensions of the thermal conductivity are $W/m/K$, i.e. $kg.m.s^{-3}.K^{-1}$, so that

$$k_m = \frac{M_m L_m}{T_m^3 \theta_m} = \frac{K_M K_L}{K_T^3 K_\theta} \frac{M_r L_r}{T_r^3 \theta_r} = \frac{K_M K_L}{K_T^3 K_\theta} k_r$$

Thermal diffusivity The dimensions of thermal diffusivity are m^2/s so that

$$\kappa_m = \frac{L_m^2}{T_m} = \frac{K_L^2}{K_T} \frac{L_r^2}{T_r} = \frac{K_L^2}{K_T} \kappa_r$$

Heat capacity The dimensions are $J.kg^{-1}.K^{-1}$, i.e. $m^2.s^{-2}.K^{-1}$ so that

$$(c_p)_m = \frac{L_m^2}{T_m^2 \theta_m} = \frac{K_L^2}{K_T^2 K_\theta} \frac{L_r^2}{T_r^2 \theta_r} = \frac{K_L^2}{K_T^2 K_\theta} (c_p)_r$$

Radiogenic heat production The dimensions are $W.m^{-3}$, i.e. $kg.m^{-1}.s^{-3}$ so that

$$H_m = \frac{M_m}{L_m T_m^3} = \frac{K_M}{K_L K_T^3} \frac{M_r}{L_r T_r^3} = \frac{K_M}{K_L K_T^3} H_r$$

A constant The dimensions of A are $Pa^{-n}.s^{-1}$, or $M^{-n}L^nT^{(2n-1)}$ so that

$$A_m = \frac{L_m^n T_m^{2n+1}}{M_m^n} = \frac{K_L^n K_T^{2n-1}}{K_M^n} A_r$$

Activation energy The dimensions of Q are $J.mol^{-1}$, but the dimensions of $\tilde{Q} = Q/R$ are K , so that

$$\tilde{Q}_m = \theta_m = K_\theta \tilde{Q}_r$$

Activation volume The dimensions of V are $m^3.mol^{-1}$, but the dimensions of $\tilde{V} = V/R$ are $m^3.^\circ C.J^{-1}$, so that

$$\tilde{V}_m = \frac{K_L K_T^2 K_\theta}{K_M} \tilde{V}_r$$

The following coefficients are used in the code too :

$$K_\mu = \frac{\mu_m}{\mu_r} = \frac{K_M}{K_L K_T}$$

$$K_v = \frac{v_m}{v_r} = \frac{K_L}{K_T}$$

$$K_{stress} = \frac{K_M}{K_L K_T^2}$$

$$K_\rho = \frac{\rho_m}{\rho_r} = \frac{K_M}{K_L^3}$$

9.38 Spectral methods

Trubitsyn *et al.* (2008) [1286]

Chapter 10

Geodynamics GEO3-1313 syllabus (Utrecht University)

chapter9.tex

This course has officially be retired and is no more given by me with this content in the UU Batchelor program.

What follows was written by Arie van den Berg and was/is used as the syllabus for the 3rd year geodynamics course at Utrecht University. It is reproduced with Arie's permission and has been slightly modified by me.

10.1 Introduction

The internal constitution of the Earth has been investigated systematically from the nineteenth century on. With the advent of seismological instrumentation for the registration of tele-seismic events, by the end of that century, the main tool for obtaining direct information about distribution of the material properties controlling seismic wave propagation became available. Before this, mainly global properties could be determined from gravity and magnetic field observations, astronomical data and indications about the heatflow from the Earth's interior. As a result of the early seismological investigations the main internal structure of the Earth was revealed within the first few decades of the twentieth century with the discovery of the earth's core in 1906 by Oldham and Gutenberg (1912) and the solid inner core in 1936 by Lehmann.

From the radial distribution of the seismic velocity profile, obtained by processing the tables of traveltime versus epicentral distance, Williamson and Adams (1923) [1363] made a first estimate of the density profile for a compressible homogeneous mantle model, consistent with the total mass of the Earth and obtained at the same time strong indication for a high density core, compositionally distinct from the mantle. They concluded that "It is therefore impossible to explain the high density of the Earth on the basis of compression alone. The dense interior cannot consist of ordinary rocks compressed to a small volume; we must therefore fall back on the only reasonable alternative, namely, the presence of a heavier material, presumably some metal, which, to judge from its abundance in the Earth's crust, in meteorites and in the Sun, is probably iron."

Bullen¹ (1975) [168] further refined the analysis and showed the assumption of a homogeneous mantle to be inconsistent with the known moment of inertia of the Earth. In the 1940s and 1950s he introduced a global division of the Earth in concentric shells, labelled A through G, ranging from the Earth's crust (A), bounded by the moho discontinuity, to the inner core (G). Region C between,

¹Keith Edward Bullen (29 June 1906 - 23 September 1976) was a New Zealand-born mathematician and geophysicist. He is noted for his seismological interpretation of the deep structure of the Earth's mantle and core.

roughly 400km and 900km, characterized by rapid increase of the seismic velocities, was identified by Bullen as a transition region between the upper mantle region B and a homogeneous lower mantle, region D. The deduced inhomogeneity of the mantle was projected by Bullen in this C region. E through G were used to label subdivisions of the core. Region E indicated the liquid, adiabatic outer core, F a transition region between inner and outer core and G the solid inner core. Birch (1952) [89] published improved equations of state, based on finite-strain theory, thereby giving a more firm physical basis to interpretation of available data in terms of a compressible medium.

In the second half of the twentieth century the resolution and accuracy of the models were further improved using continuously improved seismological observations and a growing data set. It also became possible to obtain independent information about the radial density distribution from spectral analysis of radial eigen-vibrations of the Earth after very large earthquakes. This development resulted in the publication of the Preliminary Reference Earth Model (PREM) by Dziewonski and Anderson (1981) [357] which still serves as a global reference.

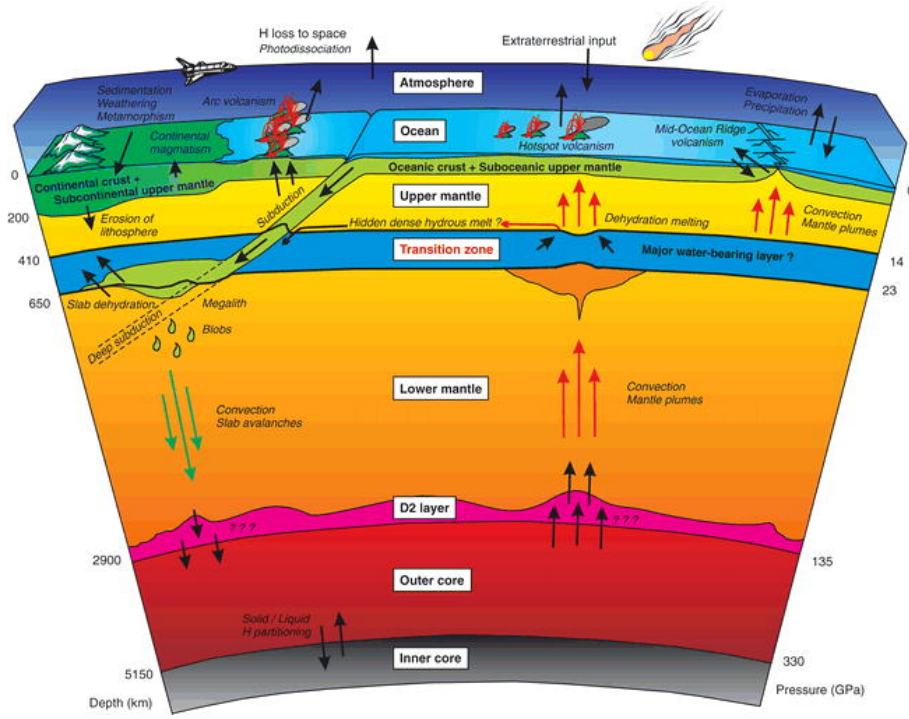
The improved seismological models indicated that the continuous rapid velocity increase in the transition zone (C) was actually a succession of several abrupt changes, confirming radial inhomogeneity in mineral phase and possibly in chemical composition of the mantle.

From geological and cosmochemical arguments a probable composition of the Earth had been derived consisting of a mantle with major element composition dominated by magnesium-iron silicates and an iron-nickel core with a small amount of lighter elements mixed in, most likely including mainly sulphur. In the 1960s this resulted in the definition of a so-called pyrolitic composition of the mantle by Ringwood which could explain the main mantle petrological observations regarding the complementary nature of basalts and ultra mafic mantle rocks found in ophiolites, kimberlites and mantle peridotite bodies (Ringwood, 1975 [1074]).

In experimental high-pressure and temperature work on the candidate mantle materials a series of phase transitions were found at pressure and temperature values relevant for the Earth's mantle which could be related to the seismic discontinuities revealed by the seismological data. From these the most prominent at approximately 410 and 660km depth were identified as the phase transition of the olivine component $(\text{Mg, Fe})_2\text{SiO}_4$ of the pyrolitic mantle to a denser wadsleyite crystal structure and, at 660km, a transition (dissociation) from a γ -spinel (known as ringwoodite) structure to a two-phase assemblage, post-spinel, i.e. magnesium-iron perovskite, $(\text{Mg, Fe})\text{SiO}_3$ and wüstite $(\text{Mg, Fe})\text{O}$.

It was also found that the 660 km boundary corresponds to an endothermic phase transition which would have implications for large scale circulation in the mantle, leading to long-standing speculations about the degree of layering in mantle convection (Christensen & Yuen (1985) [251], Albarede & van der Hilst (2002) [4]), <http://www.mantleplumes.org>.

A more recent development in this area is the discovery of a new phase transition of magnesium-perovskite to a denser form for pressure temperature conditions, approximately 125 GPa 2500 K, relevant for the D'' layer close to the core-mantle boundary (Lay *et al.* , 2005, van der Hilst *et al.* , 2007).



Effect of water on the phase relations in Earth's mantle and deep water cycle,

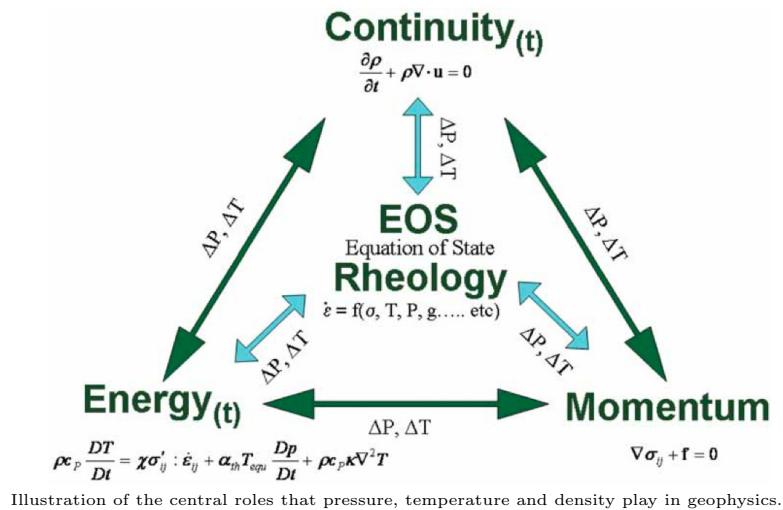
Litasov & Ohtani (2007) [793]

In the following sections the density distribution in the Earth's interior is treated in relation to the gravity field and internal pressure distribution of a self-gravitating compressible planet model and the link is shown with results from theoretical mineral physics and high pressure-temperature experimental data for mantle materials.

10.2 Global internal structure and temperature of the Earth

To understand the Earth's internal dynamics and evolution we need to know its internal structure and material properties. What do we know about Earth's global internal structure?

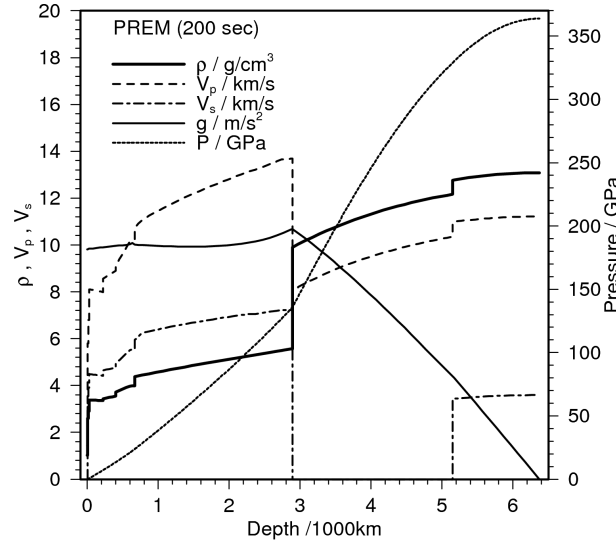
For a substance of given chemical composition, the material properties are determined by temperature and pressure. A full understanding of the Earth's internal dynamics therefore requires that we know the internal distribution of composition, temperature and pressure as illustrated in the following figure:



Taken from Regenauer-Lieb *et al.*, Phil. Mag., 2006 [1055].

The internal pressure distribution is directly linked with the Earth's own internal gravity field and density distribution because the local pressure gradient equals the local gravity acceleration times the density (see problem 6). In Section 10.4 density, gravity and pressure are treated together in a consistent way.

If the internal pressure distribution is known we can relate sharp transitions in the physical parameters as shown in the PREM model, illustrated in the following figure to phase transitions, solid-solid or solid-liquid, in the Earth's deep interior.



Radial (depth) distribution of density ρ , seismic velocities v_p and v_s , gravity acceleration g and pressure P in the PREM model (Dziewonski and Anderson (1981) [357]).

Phase transitions in ‘candidate’ materials for the Earth's interior are investigated under high pressure and temperature conditions in HPT laboratory experiments². Using theoretical mineral physics models the complete mineral phase diagram of mantle silicates can, in principle, be constructed from a limited set of experimental data (Stixrude & Lithgow-Bertelloni (2005) [1214, 1213] Jacobs & de Jong (2007) [629]). To constrain the possible candidate materials we also need to know the internal distribution of the Earth's chemical composition. Such composition models are derived from geological evidence and cosmochemical considerations.

Table I of (Dziewonski and Anderson (1981) [357]) gives an expression for the density as a function of the radius r , which I have turned into a python function:

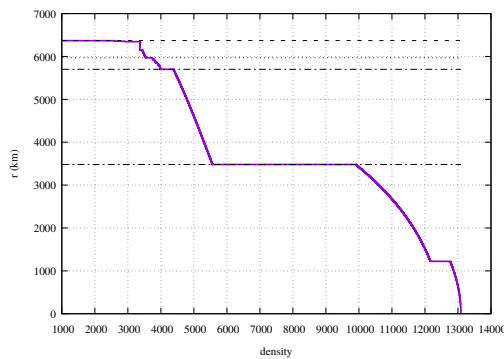
```
def prem_density(radius):
    x=radius/6371.e3
    if radius>6371.e3:
        densprem=0
    elif radius<=1221.5e3:
        densprem=13.0885-8.8381*x**2
    elif radius<=3480e3:
        densprem=12.5815-1.2638*x-3.6426*x**2-5.5281*x**3
    elif radius<=3630.e3:
        densprem=7.9565-6.4761*x+5.5283*x**2-3.0807*x**3
    elif radius<=5600.e3:
        densprem=7.9565-6.4761*x+5.5283*x**2-3.0807*x**3
    elif radius<=5701.e3:
        densprem=7.9565-6.4761*x+5.5283*x**2-3.0807*x**3
    elif radius<=5771.e3:
        densprem=5.3197-1.4836*x
    elif radius<=5971.e3:
```

²Deep Earth pressure and temperature conditions can be produced in a Diamond Anvil Cell (DAC), see http://en.wikipedia.org/wiki/Diamond_anvil_cell.

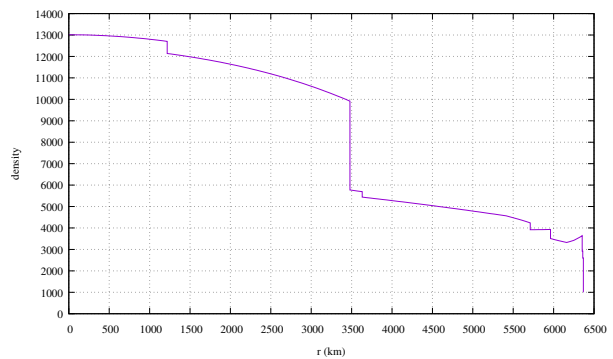
```

densprem=11.2494-8.0298*x
elif radius <=6151.e3:
    densprem=7.1089-3.8045*x
elif radius <=6291.e3:
    densprem=2.6910+0.6924*x
elif radius <=6346.e3:
    densprem=2.6910+0.6924*x
elif radius <=6356.e3:
    densprem=2.9
elif radius <=6368.e3:
    densprem=2.6
else:
    densprem=1.020
return densprem*1000

```



Top: PREM density as computed with the function above. Code available in `/images/prem/`;



Bottom: ak135 density from [696]. Data available in `/images/ak135/`

Early models of the Earth's density

The total Earth mass M_{\oplus} and average density $\langle \rho \rangle$ were not known before independent measurement of Newton's gravitational constant by Cavendish, (see Section 10.4). When the average density had been determined as approximately $5.5 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$ it became clear, from the lower density of surface rocks of around $2.7 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$, that the Earth's interior must consist of higher density material.

Besides the mass or average density the (average) moment of inertia I (defined in Section 10.3) provides a constraint on the radial distribution of density.

These two integral parameter values have been applied in several two-parameter models for the radial density distribution of the Earth. At the end of the nineteenth century Wiechert³ assumed that the compressibility of Earth materials would be negligible to first approximation and that Earth's high mean density was due to a dense, probably metallic, core. He assumed an iron core based on astronomical evidence of high iron content of the sun's outer layers (see also Section 10.5.4).

Wiechert considered in particular layered spherically symmetric models consisting of two uniform layers, core and mantle. Since the radius of the Earth's core had not yet been determined by seismology, Wiechert used the core radius R_c and density ρ_c as unknown parameters to be determined from the known data. Wiechert assumed the density of the mantle to be $\rho_m = 3.2 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$ and using known values for M and I he derived for the radius of the core $R_c/R = 0.779$ corresponding to a mantle depth of about 1400 km and a core density $\rho_c = 8.2 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$. This model is investigated in problem 3.

Later, after $R_c/R = 0.545$ had been determined using seismic data, Jeffreys substituted the known value of the core radius and derived for the mantle and core densities $\rho_c = 12.6 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$

³Emil Johann Wiechert (26 December 1861 - 19 March 1928) was a German physicist and geophysicist who made many contributions to both fields, including presenting the first verifiable model of a layered structure of the Earth and being among the first to discover the electron.

and $\rho_m = 4.14 \cdot 10^3 \text{ kg} \cdot \text{m}^{-3}$ (Bullen (1975) [168]). This model is investigated in problem 4.

The (radially averaged) density distribution in the Earth remains a topic of research [695].

10.3 The moment of inertia of a spherically symmetric density distribution

The moment of inertia I of a point mass of mass m , with respect to a given rotation axis is defined as $I = md^2$ where d is the distance from the point mass to the axis. This quantity relates the angular velocity ω , about the rotation axis, to the angular momentum J , of the point mass, in $J = I\omega$. This is an analogous relation as the one between the linear momentum p and the linear velocity v , $p = mv$. For an extended mass distribution in a volume V , a moment of inertia tensor, I_{ij} , relating the angular momentum vector \mathbf{J} to the rotation vector $\mathbf{\Omega}$ can be defined as $J_i = I_{ij}\Omega_j$, where the summation convention for repeated indices is implied. This tensor is described by a 3×3 matrix defined by volume integration over point masses in the volume. Here we only consider spherically symmetric mass distributions where the moment tensor is isotropic, $I_{ij} = I\delta_{ij}$, with scalar coefficient I .⁴ In simple terms, the moment of inertia is the same for any rotation axis through the centre of the spherically symmetric body.

The moment of inertia I can be determined from Earth's global gravity field and the precession rate of the rotation axis determined from astronomical data, see Bullen, *The Earth's density*, 1975. The principal moments of inertia can also be calculated with the hydrostatic equilibrium figure of the Earth [795].

The scalar moment of inertia is defined as a volume integral over point masses,

$$I = \int_V \rho(\vec{r}) d(\vec{r})^2 dV. \quad (10.1)$$

where $d(\vec{r})$ is the distance from point \vec{r} to the rotation axis.

For a *spherically symmetric* body of finite volume, it is often expressed in terms of the total mass M , the outer radius R and a prefactor f as,

$$I = fMR^2 \quad (10.2)$$

| Body | Value | Source |
|---------|---------------------|--------|
| Earth | 0.3307 | [3] |
| Mars | 0.3662 ± 0.0017 | [4] |
| Mercury | 0.346 ± 0.014 | [5] |
| Moon | 0.3929 ± 0.0009 | [6] |
| Venus | unknown | |

Taken from Wikipedia⁵

We have seen that the planetary mass and surface density were used to constrain models for the interior density distribution. These models are further constrained by the planets moment of inertia

⁴ δ_{ij} is the Kronecker delta, i.e. $\delta_{ij} = 1$ for $i = j$ and zero otherwise.

⁵https://en.wikipedia.org/wiki/Moment_of_inertia_factor

I that can be determined from (satellite) geodetic and astronomical observations. For Earth the following values for the total mass and moment of inertia prefactor have been found,

$$\begin{aligned} M &= 5.97 \cdot 10^{24} \text{kg} \\ I &= 0.3307 M R^2 \text{kg m}^2 \end{aligned}$$

where $R = 6371 \text{km}$ is the mean radius. The observed moment of inertia prefactor $f = 0.3307$ is smaller than the value 0.4 for a homogeneous sphere (see problem 2), another indication of mass concentration towards the earth's centre.

problem: 1. *Derive the following expression for the moment of inertia of a spherically symmetric Earth model with outer radius R ,*

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr \quad (10.3)$$

Hint: use the symmetry and compute $I = \frac{1}{3}(I_x + I_y + I_z)$, where I_x is the moment of inertia with respect to a rotation axis coinciding with the x -axis.

problem: 2. *Derive from Eq. (10.3) the value of the prefactor f of the moment of inertia for a uniform sphere. answer: $f = 2/5$.*

In general the moment of inertia prefactor f is an indicator of the degree of mass concentration towards the centre of a spherically symmetric mass distribution. Endmembers of mass concentration are a) a concentrated central point mass and b) all mass concentrated on a spherical surface of zero thickness.

Verify that the moment of inertia of the point mass endmember equals zero and that for the prefactor for a spherical shell of vanishing thickness we have $f = \frac{2}{3}$.

Add bit of theory for delta function in sph coords

Wiechert's two-layer model with a distinct core is constrained by the moment of inertia prefactor f , the mantle radius R and density ρ_m and the total mass M or, equivalently, the mean density $\langle \rho \rangle$. Expressions for the core radius R_c and density ρ_c can be formulated for this model as specified in the following exercise (Bullen, 1975).

for 2025: add figure for pb 2,3. Add plot of $\rho(r)$ for pb 3.

problem: 3. Derive a 2-parameter model for the earth's 1-D radial density distribution $\rho(r)$ consisting of two uniform layers (core and mantle) of radius R_c and R respectively and with contrasting uniform densities ρ_c and ρ_m for core and mantle respectively. Assume ρ_m to be known, leaving ρ_c and R_c as unknown parameters that can be determined from the known moment of inertia prefactor f and the average density $\langle \rho \rangle$.

Compute the total mass M

Compute the average density and arrive at:

$$\langle \rho \rangle = \frac{3}{R^3} \int_0^R \rho(r) r^2 dr \quad (10.4)$$

Use $I = fMR^2$ and the total mass to arrive at:

$$fR^5 \langle \rho \rangle = 2 \int_0^R \rho(r) r^4 dr \quad (10.5)$$

Derive the following expressions for R_c and ρ_c ,

$$\frac{R_c}{R} = \left(\frac{\frac{5}{2} f \frac{\langle \rho \rangle}{\rho_m} - 1}{\frac{\langle \rho \rangle}{\rho_m} - 1} \right)^{1/2}, \quad \rho_c = \rho_m \left\{ 1 + \left(\frac{R}{R_c} \right)^3 \left(\frac{\langle \rho \rangle}{\rho_m} - 1 \right) \right\} \quad (10.6)$$

In Bullen's two-layer model the core radius is assumed to be known from seismology. For this model the mantle and core densities can be expressed in the known parameters in the following problem.

problem: 4. Assume the core radius R_c to be a known parameter in the following. Derive a 2-parameter model for the earth's 1-D radial density distribution $\rho(r)$ consisting of two uniform layers (core and mantle), with a core and mantle radius R_c and R and different uniform densities ρ_m and ρ_c for mantle and core. Express the parameters ρ_m and ρ_c in terms of the mass and moment of inertia.

Hint: compute M first, then I , as a function of all other parameters. Establish a relationship of the form $(M, I)^T = A \cdot (\rho_c, \rho_m)^T$ where A is a 2×2 matrix.

Solution: in matrix-vector format,

$$\begin{pmatrix} \rho_c \\ \rho_m \end{pmatrix} = \frac{4\pi}{3\Delta} \begin{pmatrix} \frac{2}{5}(R^5 - R_c^5) & -(R^3 - R_c^3) \\ -\frac{2}{5}R_c^5 & R_c^3 \end{pmatrix} \begin{pmatrix} M \\ I \end{pmatrix} \quad (10.7)$$

where the determinant $\Delta = \frac{32\pi^2}{45} (R_c^3(R^5 - R_c^5) - R_c^5(R^3 - R_c^3))$.

Carry out live demo of python code. Code is in images/geodynamics

problem: 5. SKIP THIS PROBLEM. The numerical value of the interim expressions in (10.7) exceeds the magnitude of single precision real type variables in computer programs, that are limited to approximately $1.7 \cdot 10^{38}$. A work around for this problem may be to use double precision real variables that have a higher maximum magnitude of about 10^{308} .

An alternative solution is to switch to using non-dimensional parameters, denoted by primes, in the following way: define $R'_c = R_c/R$, $M_0 = 4/3 \cdot \pi R^3 \rho_0$ and $M = M_0 \cdot M/M_0 = M_0 \cdot M'$, $\rho_c = \rho_0 \rho'_c$, $\rho_m = \rho_0 \rho'_m$ and express the moment of inertia in the reference density ρ_0 and outer radius as, $I = f M R^2 = f 4/3 \cdot \pi R^5 \rho_0$. With these definitions rewrite (10.7) into the non-dimensional form,

$$\begin{pmatrix} \rho'_c \\ \rho'_m \end{pmatrix} = \frac{16\pi^2}{9\Delta'} \begin{pmatrix} \frac{2}{5}(1 - R'^5_c) & -(1 - R'^3_c) \\ -\frac{2}{5}R'^5_c & R'^3_c \end{pmatrix} \begin{pmatrix} M' \\ f \end{pmatrix} \quad (10.8)$$

where the determinant $\Delta' = \frac{32\pi^2}{45} (R'^3_c(1 - R'^5_c) - R'^5_c(1 - R'^3_c))$.

10.4 Density, gravity and pressure in the Earth

In the Earth's mantle major solid state phase transitions occur in the silicate material which constitutes the planetary mantle outside the metallic iron/nickel core. These phase transitions are induced by the increase in the static pressure from a 1 bar (10^5Pa) atmospheric value at the Earth's surface to $136 \cdot 10^9 \text{Pa}$ at the core mantle boundary at a depth of approximately 2900km. Phase transitions in the Earth's interior are associated with changes in the elastic wave velocities that can be deduced from seismological observations. In high pressure experiments, phase transitions in candidate mantle silicates can be studied and correlated with the seismological data to constrain the mineralogy and pressure/temperature distribution in the mantle. Knowledge of the internal material constitution of the Earth, such as the mineral phase, is a requirement for understanding the main geodynamical processes that determine Earth's evolution.

Density and pressure inside the Earth are linked with self-gravitation. This means that the hydrostatic or lithostatic pressure is a direct result of the gravity field generated by the Earth's own mass distribution. The lithostatic pressure can be expressed as the weight of a column of unit cross-sectional area extending from zero depth, at the Earth's surface, to the depth z of the evaluation point,

$$P(z) = \int_0^z \rho(z')g(z')dz' \quad (10.9)$$

where ρ is the mass density and g is the magnitude of the gravitational acceleration.

The gravity field defining g is generated by the Earth's own density distribution. Weak periodic gravity 'perturbations' are generated by celestial bodies, expressed in the external tides, both ocean tides and solid earth tides. The main tides are generated by the Earth's moon and by the Sun.

In the following section expressions for the gravity field in terms of the density distribution are given, based on Newton's law of gravitation.

In the description of the density distribution we will first neglect the role of self-compression and consider a number of one-dimensional (1-D), spherically symmetric, parameterized density distributions. Self-compression and compressibility are then treated in section 10.5.2. Self-compression and finite compressibility result in a continuous increase of density with pressure in agreement with several geophysical observations.

problem: 6. Derive the expression (10.9) (where the depth z is not to be confused with a cartesian coordinate) for the lithostatic pressure in a spherically symmetric planet from the elastostatic equation for a static medium,

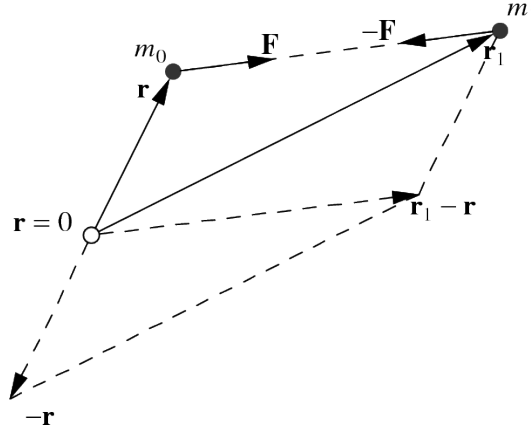
$$\partial_j \sigma_{ij} + \rho g_i = 0 \quad \Leftrightarrow \quad \vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} = \vec{0} \quad (10.10)$$

Hint: Assume hydrostatic conditions where the stress tensor can be written as $\sigma_{ij} = -P\delta_{ij}$, with δ_{ij} the Kronecker delta, and derive from equation (10.10) for the pressure gradient, $\vec{\nabla}P = \rho\vec{g}$.

Gravity field of a mass distribution

Newton formulated the attraction force acting on a point mass m_0 , located in a point with position vector $\vec{r} = (x, y, z)$, with x, y, z the cartesian coordinates, from a second point mass m_1 located at $\vec{r}_1 = (x_1, y_1, z_1)$, illustrated in the following figure as,

$$\vec{F}(\vec{r}) = \frac{\mathcal{G}m_0m_1}{|\vec{r}_1 - \vec{r}|^2} \vec{e}_{\vec{r}\vec{r}_1} \quad (10.11)$$



Vector diagram of the gravitational forces acting on the two point masses m_0, m_1 in vector locations \vec{r} and \vec{r}_1 respectively. From the expression for the gravity field (10.11) it follows that the forces on both masses are of equal magnitude and in opposite direction.

Where $\vec{e}_{\vec{r}\vec{r}_1}$ is the unit vector in \vec{r} pointing towards \vec{r}_1 and $\vec{F}(\vec{r}_1) = -\vec{F}(\vec{r})$. \vec{r}, \vec{r}_1 are the position vectors of the two point masses and $|\vec{r}_1 - \vec{r}| = \sqrt{(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2}$ is the distance between the points \vec{r} and \vec{r}_1 . \mathcal{G} is the gravitational constant $\mathcal{G} \simeq 6.67 \times 10^{-11} \text{N m}^2 \text{kg}^{-2}$, m_0, m_1 the mass of the respective pointmasses.

This gravitation effect is usually specified as a gravitation force per unit mass or acceleration vector \vec{g} ,

$$\vec{g}(\vec{r}) = \frac{\mathcal{G}m_1}{|\vec{r}_1 - \vec{r}|^2} \vec{e}_{\vec{r}\vec{r}_1} \quad (10.12)$$

It can be verified by inspection that the acceleration vector field can be written as the gradient of a scalar potential field $U(\vec{r})$ (i.e. the potential energy per unit mass) with

$$\vec{g} = -\vec{\nabla}U = \left(-\frac{\partial U}{\partial x}, -\frac{\partial U}{\partial y}, -\frac{\partial U}{\partial z}\right),$$

in Cartesian coordinates (see problem 9), and

$$U(\vec{r}) = -\frac{\mathcal{G}m_1}{|\vec{r}_1 - \vec{r}|} \quad (10.13)$$

The gravity acceleration and corresponding potential field are additive such that the total force or potential of a collection of N point masses is obtained by summation over individual point contributions,

$$\vec{g}(\vec{r}) = \sum_j^N \frac{\mathcal{G}m_j}{|\vec{r}_j - \vec{r}|^2} \vec{e}_{\vec{r}\vec{r}_j}, \quad U(\vec{r}) = -\sum_j^N \frac{\mathcal{G}m_j}{|\vec{r}_j - \vec{r}|} \quad (10.14)$$

With this definition and sign convention the potential field of a point source in the origin is represented by a potential well ($U(\vec{r}) < 0$). This is known as Coulomb's law and the equivalent form for a continuous mass distribution of density ρ (mass per unit volume) contained in a volume V is,

$$\vec{g}(\vec{r}) = \int_V \frac{\mathcal{G}\rho(\vec{r}')}{|\vec{r}' - \vec{r}|^2} \vec{e}_{\vec{r}\vec{r}'} dV(\vec{r}'), \quad U(\vec{r}) = -\int_V \frac{\mathcal{G}\rho(\vec{r}')}{|\vec{r}' - \vec{r}|} dV(\vec{r}') \quad (10.15)$$

Besides the integral expression for the gravity field defined in (10.15) there is also the differential form using the second order partial differential equations of Laplace and Poisson. It can be shown by verification (see hereafter) that U in (10.15) satisfies Poisson's equation,

$$\vec{\nabla}^2 U = 4\pi\mathcal{G}\rho \quad (10.16)$$

which reduces to Laplace's equation $\vec{\nabla}^2 U = 0$ outside the mass distribution in V (where $\rho = 0$).

To show that U in (10.15) satisfies Poisson's equation integrate the normal component of the acceleration field over an arbitrary closed surface S enclosing V and change the order of integration for the volume and surface integral.

$$\int_S \vec{\nabla} U(\vec{r}) \cdot \vec{n} dA(\vec{r}) = - \int_V \mathcal{G}\rho(\vec{r}') \left\{ \int_S \vec{\nabla} \left(\frac{1}{|\vec{r}' - \vec{r}|} \right) \cdot \vec{n} dA(\vec{r}) \right\} dV(\vec{r}') \quad (10.17)$$

The surface integral on the right is independent of the choice of the surface S as long as it contains \vec{r}' . We therefore replace this surface by a sphere of radius R centered at \vec{r}' and find for the surface integral the value -4π .

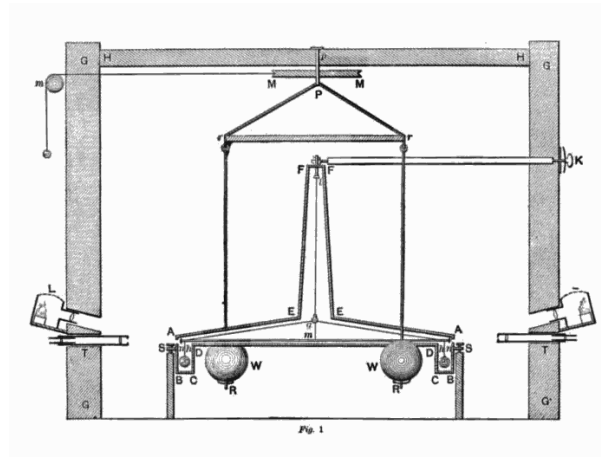
Next we apply the Gauss divergence theorem to the left hand surface integral to obtain,

$$\int_V \vec{\nabla}^2 U dV = \int_V 4\pi\mathcal{G}\rho dV \quad (10.18)$$

Note that the surface has been contracted on the volume V to obtain (10.18). Since the surface and enclosed volume are arbitrary we obtain the Poisson equation,

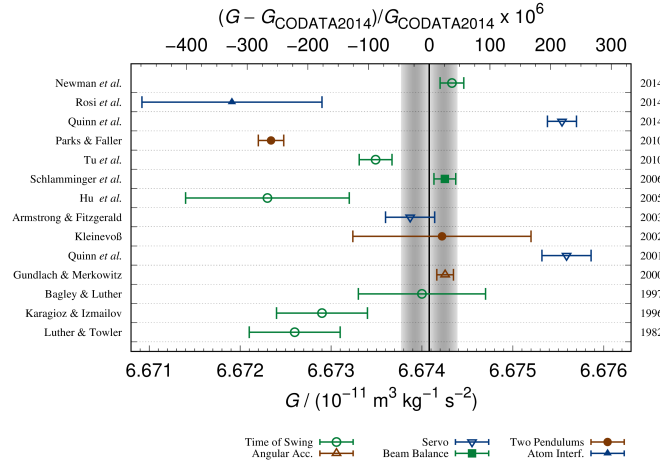
$$\vec{\nabla}^2 U = 4\pi\mathcal{G}\rho \quad (10.19)$$

In Newton's time the numerical value of \mathcal{G} had not been determined yet. As a result it was not possible to determine the mass of the Earth M_\oplus by measuring the gravitation force of the Earth on a known 'test mass'. This way only the value of GM_\oplus could be determined. Only with the experiment named after Cavendish (1798) ⁶ it became possible to measure \mathcal{G} directly, in a torsion balance experiment, by determining the gravitational attraction of two closely spaced test masses shown here:



Since Cavendish many experiments have been conducted in order to determine the value of this constant:

⁶http://en.wikipedia.org/wiki/Cavendish_experiment



This chart compares the results from a dozen experiments measuring \mathcal{G} . The vertical stripe represents the most recent recommended value for G (black line) with its error bar (gray). Far to the right are the two outlying BIPM measurements, in blue. Taken from <https://www.nist.gov/image/glabel2016plotfromstephanpng>

The recommended value is $\mathcal{G} = 6.67430(15) \cdot 10^{-11} \text{m}^3 \text{kg}^{-1} \text{s}^{-2}$, see for instance <https://physics.nist.gov/cgi-bin/cuu/Value?bg>.

Multipole expansion

The idea behind the multipole expansion is simple: the denominator in Eq. (10.13) can be rewritten in such a way that it can be expanded as a Taylor series. We have

$$\frac{1}{|\vec{r} - \vec{r}'|} = \frac{1}{\sqrt{(\vec{r} - \vec{r}') \cdot (\vec{r} - \vec{r}')}} = \frac{1}{\sqrt{\vec{r} \cdot \vec{r} - 2\vec{r} \cdot \vec{r}' + \vec{r}' \cdot \vec{r}'}} = \frac{1}{r^2} \frac{1}{\sqrt{1 - 2\frac{\vec{r} \cdot \vec{r}'}{r^2} + \left(\frac{r'}{r}\right)^2}}$$

The potential can be expanded in a series of Legendre polynomials.

$$(1 - 2XZ + Z^2)^{-1/2} = \sum_{n=0}^{\infty} Z^n P_n(X)$$

valid for $|X| \leq 1$ and $|Z| \leq q$. The coefficients P_n are the Legendre polynomials of degree n .
FINISH!

Gravitational potential energy

For two pairwise interacting point particles, the gravitational potential energy \mathcal{U} is given by

$$\mathcal{U} = \frac{-\mathcal{G}Mm}{R}$$

where M and m are the masses of the two particles, R is the distance between them. Close to the Earth's surface, the gravitational field is approximately constant, and the gravitational potential energy of an object reduces to

$$\mathcal{U} = mgh$$

where m is the object's mass, $g = \mathcal{G}M_E/R_E^2$ is the gravity of Earth, and h is the height of the object's center of mass above a chosen reference level.

Let us talk units

The SI units for (gravity) acceleration are m s^{-2} . However in the context of gravity, we will rarely encounter these.

The Gal is the commonly used unit in gravimetry:

$$0.01\text{m s}^{-2} = 1\text{Gal}$$

and often measurements are given in mGal or μGal .

As such, the acceleration due to Earth's gravity at its surface is 976 to 983 Gal, the variation being due mainly to differences in latitude and elevation.

Gravity Force Inside a Spherical Shell

In classical mechanics, the *shell theorem* gives gravitational simplifications that can be applied to objects inside or outside a *spherically symmetrical* body.

Isaac Newton proved the shell theorem and stated that:

1. A spherically symmetric body affects external objects gravitationally as though all of its mass were concentrated at a point at its centre.
2. If the body is a spherically symmetric shell (i.e., a hollow ball), no net gravitational force is exerted by the shell on any object inside, regardless of the object's location within the shell.



These two propositions are not easy to prove. The second one is very important: it states that if I stand mid-mantle at a radius of, say, 5000km, the 1371km-thick shell of rock above me does not contribute to the force of gravity that I am feeling. Only the rocks below my feet contribute to this force. At this location we can write

$$\frac{\mathcal{G}mM(r)}{r^2} = ma$$

where $M(r)$ is the mass inside a sphere of radius r . The mass m of my body cancels out, and we obtain

$$\frac{\mathcal{G}M(r)}{r^2} = a$$

The acceleration in this context is often called g and it clearly depends on r so that if density is constant, $M(r) = \frac{4\pi}{3}r^3\rho_0$ and then

$$g(r) = \mathcal{G}\frac{4\pi}{3}r\rho_0$$

It also follows that the gravity acceleration in the center of the planet ($r = 0$) must be zero and the gravity acceleration increases linearly with the distance to the center. If now the density is not constant (but radially symmetric, i.e. $\rho = \rho(r)$) then

$$g(r) = \mathcal{G} \frac{4\pi}{r^2} \int_0^r \rho(r') r'^2 dr'$$

Remember that this is true because of the spherical symmetry!

Nonuniqueness

Given a body with mass M at a distance r from me, the gravitational acceleration that I feel is

$$a = \frac{\mathcal{G}M}{r^2}$$

If the mass is now twice as far (distance $2r$) then

$$a' = \frac{\mathcal{G}M}{(2r)^2} = \frac{1}{4} \frac{\mathcal{G}M}{r^2} = \frac{1}{4}a$$

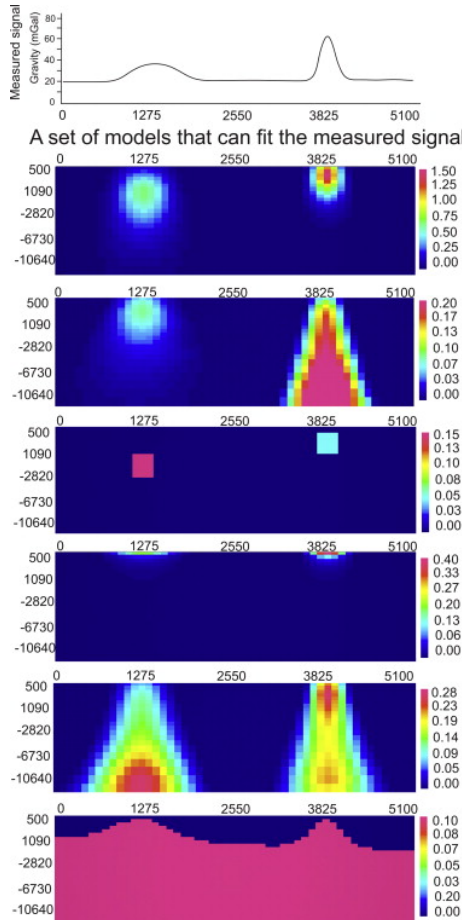
Because of the inverse square of the distance the acceleration is four times as small.

However, if I now 'make' the mass of the body four times as large and twice as far,

$$a'' = \frac{\mathcal{G}(4M)}{(2r)^2} = \frac{\mathcal{G}M}{r^2} = a$$

There lies a very important fact: There is an inescapable trade-off between distance and mass.

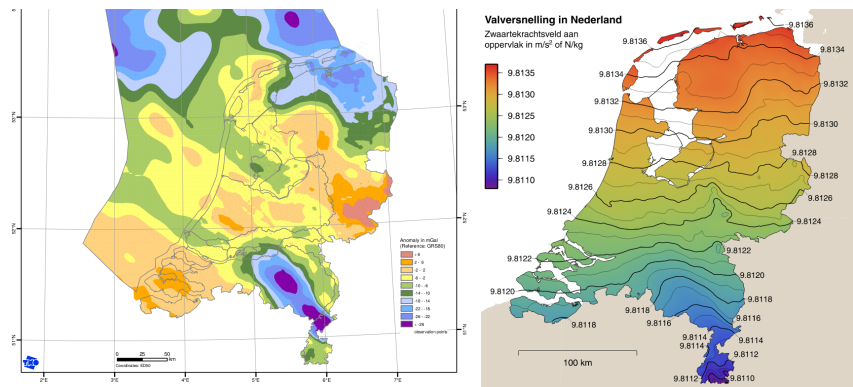
If gravity is measured at a single point in space nothing certain can be said about what lies below: the object generating the gravity anomaly could be 'close' and not so massive, or 'far' and really massive, both situations potentially leading to the same measurement.



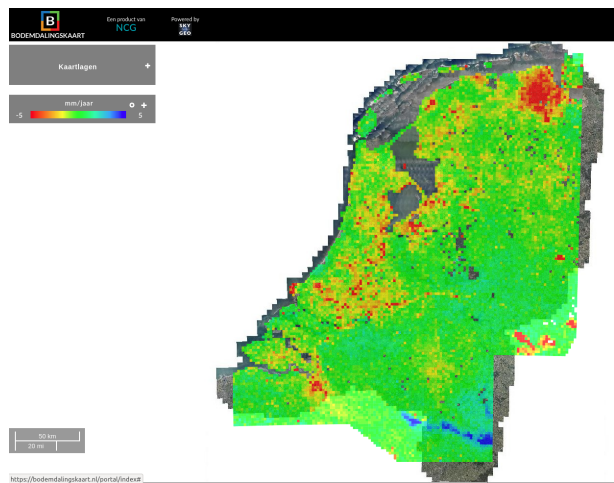
Taken from Meijde, Pail, Bingham, and Floberghagen [861] (2015).

The Netherlands

As explained in Crombaghs, Min, and Hees [287] (2002), in The Netherlands gravity values increase from south to north with about 1 milligal per kilometer. Smallest values occur in Limburg (981,100 milligal), while the largest values occur in Groningen (981,350 milligal). Local variations are limited to 1 milligal over some kilometers.



Left: Taken from <https://www.nlog.nl/en/gravity-and-magnetic-field>. The size of the Bouguer anomaly at a particular location is a measure of the mass deficit or mass excess in the underlying rocks. A mass deficit exists where the stratigraphic succession is composed of relatively light rocks; this yields a negative Bouguer anomaly. A mass excess exists where the stratigraphic succession is composed of relatively heavy rocks, this yields a positive anomaly. Right: Taken from https://upload.wikimedia.org/wikipedia/commons/c/ce/Valversnelling_in_Nederland.svg.



Taken from <https://bodemdalingkaart.nl/portal/index>. Is the continuous sinking of certain parts of the Netherlands visible in the satellite gravity rate measurements ?

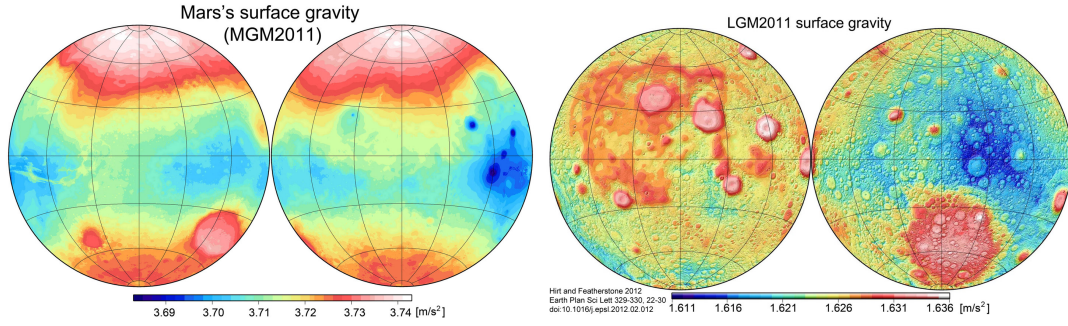
write about Anomalies, Additivity, Inversion

A few problems more to solve

problem: 7. Verify that the familiar surface value of the Earth’s gravity acceleration $g_0 = 9.8 \text{ m/s}^2$ corresponds to the value of a point mass at the Earth’s centre with the same mass as the Earth (see Table).

| | Radius km | Mass kg | Density kg/m ³ |
|---------|-------------------|----------------------|------------------------------|
| Earth | 6371 | $5.97 \cdot 10^{24}$ | 5.515×10^3 |
| Moon | 1738 | $7.34 \cdot 10^{22}$ | 3.34×10^3 |
| Mars | 3394 | $6.42 \cdot 10^{23}$ | 3.93×10^3 |
| Jupiter | 71492 | $1.9 \cdot 10^{27}$ | 1.326×10^3 |
| Sun | $6.96 \cdot 10^5$ | $1.99 \cdot 10^{30}$ | - |

Radius-mass parameters of Earth moon and planets.



Left: Mars gravity, taken from Hirt *et al.* (2011) [576]; Right: Moon gravity, taken from https://en.wikipedia.org/wiki/Gravitation_of_the_Moon

problem: 8. The PREM profile suggests that the magnitude of the gravity acceleration is approximately constant throughout the Earth's mantle. Assume an approximate uniform value of g in the Earth's mantle, equal to the surface value $g_0 \sim 9.8 \text{ m/s}^2$ and use an approximate average mantle density $\rho_m \sim 4.5 \times 10^3 \text{ kg/m}^3$ to obtain from Eq. (10.9) an approximation of the static pressure at the core mantle boundary at a depth of 2891 km.

problem: 9. Verify the consistency of the expression for the gravity acceleration and potential of a point mass in (10.12) and (10.13), i.e. prove from these expressions by explicit calculation of the gradient vector from the scalar potential field that $\vec{g} = -\vec{\nabla}U$.

Hint: specify the potential in (10.13) in Cartesian coordinates (i.e. write explicitly $|\vec{r}_1 - \vec{r}|$) and differentiate the result with respect to the coordinates x, y, z . What is the derivative of $(f(x))^\alpha$ with respect to x ? After a few steps you should then arrive at (10.12).

problem: 10. Apply the Poisson equation (10.16) to obtain the gravity field of a point-mass distribution with mass M , described by a Dirac delta function, $\rho(\vec{r}) = M\delta(\vec{r} - \vec{r}_0)$. Where the following property holds for the delta function,

$$\int_V \delta(\vec{r} - \vec{r}_0) dV = \begin{cases} 1, & \vec{r}_0 \in V \\ 0, & \vec{r}_0 \notin V \end{cases} \quad \text{or, more general} \quad \int_V f(\vec{r}) \delta(\vec{r} - \vec{r}_0) dV = \begin{cases} f(\vec{r}_0), & \vec{r}_0 \in V \\ 0, & \vec{r}_0 \notin V \end{cases} \quad (10.20)$$

Hint: integrate (10.16) over a spherical volume, centered at \vec{r}_0 and apply the Gauss divergence theorem: for a vector field $\vec{A} = (A_1, A_2, A_3)$ with divergence $\vec{\nabla} \cdot \vec{A} = \frac{\partial A_1}{\partial x} + \frac{\partial A_2}{\partial y} + \frac{\partial A_3}{\partial z}$

$$\int_V \vec{\nabla} \cdot \vec{A} dV = \int_{\partial V} \vec{A} \cdot \vec{n} dS \quad (10.21)$$

where ∂V is the closed boundary surface of V .

problem: 11. Check the dimensional units in (10.16) and verify that the gravitational potential has the dimension of energy per unit mass. This is in agreement with the identification of the gravity potential with the potential (gravitational) energy of a unit mass in the gravity field. ^a

^aThe local potential field value $U(\mathbf{r}_1)$ equals the negative of the (gravitational) potential energy $W(\mathbf{r}_1)$ of a unit point mass positioned at \mathbf{r}_1 . It can be shown that the change in potential energy ΔW that results from moving a unit mass from \mathbf{r}_1 to \mathbf{r}_2 follows directly from the potential field values $U(\mathbf{r}_1)$, $U(\mathbf{r}_2)$ and is independent of the path taken between \mathbf{r}_1 and \mathbf{r}_2 . This property defines a so called conservative field U .

To derive this result we compute the potential energy difference as the path (line) integral of the work done by the gravity force field on a unit mass and apply the gradient property $\vec{g} = -\vec{\nabla}U$. The work done by moving a unit point mass from a location \mathbf{r}_1 to \mathbf{r}_2 is defined by the line integral,

$$\Delta W = \int_{\mathbf{r}_1}^{\mathbf{r}_2} \mathbf{F} \cdot d\mathbf{r} = \int_{\mathbf{r}_1}^{\mathbf{r}_2} \mathbf{g} \cdot d\mathbf{r} = \int_{\mathbf{r}_1}^{\mathbf{r}_2} -\vec{\nabla}U \cdot d\mathbf{r} = \int_{U(\mathbf{r}_1)}^{U(\mathbf{r}_2)} -dU = -(U(\mathbf{r}_2) - U(\mathbf{r}_1)) = -\Delta U \quad (10.22)$$

Here the following gradient property has been used, relating the gradient vector to the differential of the scalar potential field,

$$dU = \frac{\partial U}{\partial x}dx + \frac{\partial U}{\partial y}dy + \frac{\partial U}{\partial z}dz = \vec{\nabla}U \cdot d\vec{r} \quad (10.23)$$

The gravitational potential field can thus be defined in terms of the work done by the gravity field to move a unit mass from infinity to the evaluation point.

$$W(\mathbf{r}_1) = \int_{\mathbf{r}_\infty}^{\mathbf{r}_1} \mathbf{g} \cdot d\mathbf{r} = \int_{\mathbf{r}_\infty}^{\mathbf{r}_1} -\vec{\nabla}U \cdot d\mathbf{r} = \int_{U(\mathbf{r}_\infty)}^{U(\mathbf{r}_1)} -dU = -U(\mathbf{r}_1) + U(\mathbf{r}_\infty) = -U(\mathbf{r}_1) \quad (10.24)$$

Where $U(\mathbf{r}_\infty) = 0$ has been used.

The above can be applied in the determination of the escape velocity from the surface of a planet. This is the minimum launch velocity to escape from the planet's gravity field. For a spherically symmetric planet the external gravity potential is given by (10.33). Moving an object from the surface, the gravity potential changes by $\Delta U = U(r) - U(R) = \mathcal{G}M(-\frac{1}{r} + \frac{1}{R})$. Applying an energy conservation argument we require the change in total (potential plus kinetic) energy per unit mass to be: $\Delta E = \Delta U + \Delta K = 0$. With $\Delta K = -v_{ex}^2/2$ we get $v_{esc} = \sqrt{2\mathcal{G}M/R}$.

problem: 12. Compute the surface escape velocities for different celestial bodies using the parameters given in the Table above.

To watch:

- Tossing Satellites into Orbit with SpinLaunch: Is that Possible? (3min)
<https://www.youtube.com/watch?v=-BCeanUiKwM>
- Can We Throw Satellites to Space? (42min)
<https://www.youtube.com/watch?v=yrc632oilWo>

problem: 13. The potential energy of a self-gravitating planet in its own gravity field is defined in terms of the volume density ρU as,

$$E = - \int_V \rho U dV \quad (10.25)$$

Derive the following expression for the potential energy of a spherically symmetric, uniform density model, using the expression for the internal gravity potential defined in (10.32)

$$E = \frac{8\pi}{5} \mathcal{G} \rho_0 M R^2 \quad (10.26)$$

Compute the potential energy value E , assuming a density $\rho_0 = 5.5 \cdot 10^3 \text{kg/m}^3$ and planetary radius $R = 6371 \text{km}$.

answer: $4.4 \cdot 10^{32} \text{J}$

The gravitational energy considered above plays an important role in major compositional differentiation processes that occurred in the early Earth and are still occurring today.

- A so called ‘core catastrophe’ occurred when the iron/nickel core of the Earth differentiated from the silicate mantle in the first few million years after the formation of the Earth in the early solar system. This event has probably freed enough potential energy to melt the mantle completely, resulting in a global magma ocean ⁷.
- Crystallization of the solid inner core from the liquid outer core, as a result of core cooling, is accompanied by compositional differentiation. The liquid outer core contains a lighter fraction, possibly sulfur, which stays behind in the liquid during freezing of the inner core. The enriched residual liquid near the inner core boundary is less dense than the average liquid of the outer core and this results in a gravitationally unstable layering that induces ‘chemically driven’ convective flow in the outer core. The potential energy released in this chemical convection is probably an important energy source in powering the geodynamo that generates the Earth’s present day magnetic field.

10.5 The gravitational potential for spherical problems

Starting from the Poisson equation,

$$\Delta U = 4\pi \mathcal{G} \rho$$

and using Gauss’ theorem (noting that $\Delta U = \vec{\nabla} \cdot \vec{\nabla} U$):

$$\int_V \Delta U dV = \int_V \vec{\nabla} \cdot \vec{\nabla} U dV = \int_{\Gamma} \vec{\nabla} U \cdot \vec{n} dS = \int_V 4\pi \mathcal{G} \rho dV = 4\pi \mathcal{G} \int_V \rho dV$$

where \vec{n} is the outward pointing normal vector.

A uniform sphere of mass M and radius a (and therefore density $\rho = M/(4\pi a^3/3)$) has the potential

$$U(r) = \begin{cases} -2\pi \mathcal{G} \rho (a^2 - r^2/3) & r \leq a \\ -\mathcal{G} M/r & r \geq a \end{cases} \quad (10.27)$$

Outside the sphere the potential is Keplerian, while inside it has the form of a parabola; both the potential and its derivative are continuous at the surface of the sphere.

A sphere with density profile

$$\rho(r) = \rho_0 (r/r_0)^{-2}$$

⁷https://en.wikipedia.org/wiki/Iron_catastrophe

has the potential

$$U(r) = 4\pi\mathcal{G}\rho_0 r_0^2 \ln(r/r_0) \quad (10.28)$$

problem: 14. Verify (10.27) and (10.28). Sketch the density field and the resulting gravity field and potential.

problem: 15. Assume a spherically symmetric non-rotating Earth in hydrostatic equilibrium. In spherical coordinates the divergence of a vector field $\vec{a}(r)$, which only depends on the radius r is

$$\vec{\nabla} \cdot \vec{a} = \frac{1}{r^2} \frac{d}{dr}(r^2 a_r)$$

a) Prove that the acceleration of gravity at radius r only depends on the mass contained in the sphere of radius R_i . Hint: start from $\vec{\nabla} \cdot \vec{g}$.

b) Assume that the mass of the Earth's core is M_c . Assume a linear density profile for the crust and mantle and determine the acceleration of gravity as a function of the radius in the mantle.

As it turns out, pairs of functions related by Poisson's equation provide convenient building-blocks for galaxy models. Three such functions often used in the literature are listed here; all describe models characterized by a total mass M and a length scale a :

- Plummer (1905) [325]

$$\rho(r) = \frac{3M}{4\pi a^3} \left(1 + \frac{r^2}{a^2}\right)^{-5/2} \quad U(r) = -\frac{\mathcal{G}M}{\sqrt{r^2 + a^2}}$$

- Hernquist (1990)

$$\rho(r) = \frac{M}{2\pi} \frac{a}{r(r+a)^3} \quad U(r) = -\frac{\mathcal{G}M}{r+a}$$

- Jaffe (1983)

$$\rho(r) = \frac{M}{4\pi} \frac{a}{r^2(r+a)^2} \quad U(r) = \frac{\mathcal{G}M}{a} \ln \frac{a}{r+a}$$

(VERIFY?)

10.5.1 The gravity and pressure field for parameterized density models with self-gravitation

In the following problems a number of simple density distributions are investigated that will serve as a reference for models more constrained by geophysical observations to be introduced in later sections. The gravity field can be determined by solving the governing Poisson equation (10.16) using suitable boundary conditions. For the special case of spherically symmetric mass distributions simple 1-D integral expressions can be used to derive the corresponding radial pressure distribution.

problem: 16. The internal and external gravity field for a simple model of a planet can be derived by solving the Poisson equation (10.16), and applying appropriate boundary conditions to the general solution. Consider a spherically symmetric planet of radius R and uniform density ρ_0 .

1. Derive expressions for the gravity potential field U and the gravity force field $g = |\mathbf{g}|$ inside and outside the planet.

Hints: Solve Poisson's equation in spherical coordinates for the interior ($r \leq R$) and exterior domain $r \geq R$ separately. The separate solutions for the interior $U_{\text{int}}, g_{\text{int}}$ and exterior $U_{\text{ext}}, g_{\text{ext}}$ domain each contain two integration constants which can be determined by applying the following boundary conditions,

$$\lim_{r \rightarrow \infty} U_{\text{ext}}(r) = 0, \quad \lim_{r \rightarrow 0} g_{\text{int}}(r) < \infty \quad (10.29)$$

Continuity of the gravity acceleration g at the surface $r = R$,

$$g_{\text{int}}(R) = g_{\text{ext}}(R) \quad (10.30)$$

Continuity of the gravity potential U at the surface $r = R$,

$$U_{\text{int}}(R) = U_{\text{ext}}(R) \quad (10.31)$$

Answers

$$g_{\text{int}} = \frac{4\pi}{3} \mathcal{G} \rho_0 r, \quad U_{\text{int}} = \frac{2\pi}{3} \mathcal{G} \rho_0 r^2 - \frac{3}{2} \frac{\mathcal{G} M}{R} \quad (10.32)$$

where $M = \frac{4\pi}{3} R^3 \rho_0$ is the planet mass and \mathcal{G} is the gravitational constant.

$$g_{\text{ext}} = \frac{\mathcal{G} M}{r^2}, \quad U_{\text{ext}} = -\frac{\mathcal{G} M}{r} \quad (10.33)$$

2. Verify that the external gravity force field is identical to the field of a concentrated point mass at $r = 0$.
3. Derive an expression for the radial distribution of the pressure in the planetary interior and compute the central pressure for a case with $\rho_0 = 5.5 \cdot 10^3 \text{ kg m}^{-3}$ and $R = 6.371 \times 10^6 \text{ m}$.

Solution: $P(r) = \frac{2\pi}{3} \rho_0^2 \mathcal{G} (R^2 - r^2)$

The gravity field of a spherically symmetric density distribution is identical to the field of an equivalent point-mass. (see problem 16 for the spatial case of a uniform density distribution). This can be formulated as follows,

$$g(r) = \frac{Gm(r)}{r^2}, \quad (10.34)$$

with

$$m(r) = \int_{V(r)} \rho dV = \int_0^r \rho(r') 4\pi r'^2 dr' \quad (10.35)$$

Here $m(r)$ is the mass inside a sphere of radius r and $g(r)$ is the corresponding magnitude of the gravity acceleration. For the corresponding gravity potential this implies, with $\int_r^\infty \frac{dU}{dr'} dr' = U(\infty) - U(r) = -U(r)$,

$$U(r) = - \int_r^\infty \frac{dU}{dr'} dr' = \int_r^\infty g_r(r') dr' = \int_r^\infty -g(r') dr' = - \int_r^\infty \frac{Gm(r')}{r'^2} dr' \quad (10.36)$$

where the radial vector component g_r has been expressed in the vector length g as $g_r = \mathbf{g} \cdot \mathbf{e}_r = -g$.

To derive (10.34), the potential field at the radial coordinate r can be split in contributions originating from an internal- and external density distribution $U(r) = U_i(r) + U_e(r)$. With corresponding pairs, $U_i \leftrightarrow \rho_i$, and $U_e \leftrightarrow \rho_e$, where $\rho_e(r') = 0$, $r' \leq r$, and $\rho_e(r') = \rho(r')$, $r' > r$. This follows from the linearity of the governing Poisson equation.

The field generated by the internal mass distribution is obtained by integrating the corresponding Poisson equation in spherical coordinates,

$$\frac{1}{r'^2} \frac{d}{dr'} r'^2 \frac{dU_i}{dr'} = 4\pi G \rho_i \quad (10.37)$$

$$\int_0^r \frac{d}{dr'} \left(r'^2 \frac{dU_i}{dr'} \right) dr' = \int_0^r 4\pi G \rho_i r'^2 dr' \quad (10.38)$$

The left hand side becomes simply $r^2 \frac{dU_i}{dr}$ so the radial component of the gravity acceleration becomes

$$g_r(r) = -\frac{dU_i}{dr} = -\frac{1}{r^2} \int_0^r 4\pi G \rho_i r'^2 dr' \quad (10.39)$$

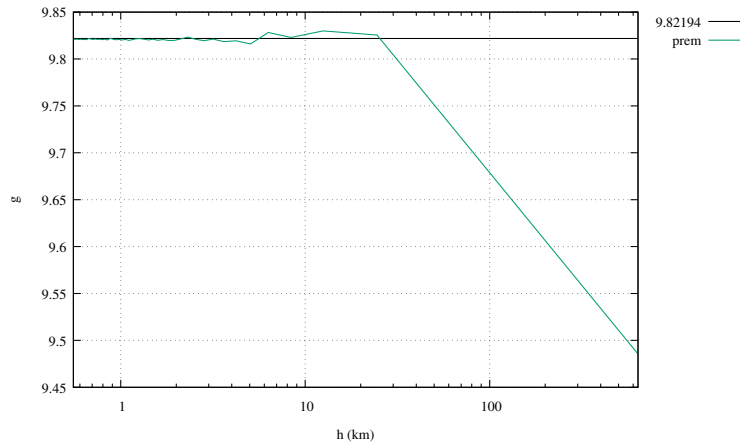
Furthermore the acceleration field g_e from the external mass distribution ρ_e for internal evaluation points $r' < r$ is zero. The corresponding gravity potential U_e is uniform, which follows from the relevant Poisson equation, in spherical coordinates for a spherically symmetric mass distribution,

$$\frac{1}{r'^2} \frac{d}{dr'} r'^2 \frac{dU_e}{dr'} = 4\pi G \rho_e = 0 \rightarrow r'^2 \frac{dU_e}{dr'} = A \rightarrow g_e(r') = -\frac{dU_e}{dr'} = -\frac{A}{r'^2} \quad (10.40)$$

A non-singular field requires $A = 0$, $g_e(r') = 0$, $r' \leq 0$ and,

$$\frac{dU_e}{dr'} = 0 \rightarrow U_e(r') = B, \quad r' \leq r \quad (10.41)$$

Looking back at the PREM model, it is a radial model $\rho(r)$ so can compute g_r at the surface of the Earth. The code is available in /images/prem/ and uses a somewhat naive numerical quadrature of Eq. (10.39):



Radial gravity component measured at the surface of the Earth as a function of the cell size used for the integration.

For very high resolutions we do recover the analytical value as computed in Section 18.1.

problem: 17. Verify that 10.36, applied to the special case of a homogeneous sphere of density ρ_0 , lead to the same expression for the internal and external potential and acceleration field as given in problem 16.

For a two-parameter spherically symmetric planet model consisting of a uniform core and mantle with radius R_c and R_m and contrasting densities ρ_c and ρ_m , the gravity field can also be determined by solving the Poisson equation for the particular density distribution and determination of the integration constants from the boundary conditions. However in this case the formula (10.34) are more convenient to obtain expressions for the gravity field.

problem: 18. *Derive expressions for the gravity acceleration and internal pressure distribution for the two-parameter model*

$$\rho(r) = \begin{cases} \rho_c, & r < R_c \\ \rho_m, & R_c < r \leq R \\ \rho_e = 0, & r > R \end{cases}, \quad g(r) = \begin{cases} g_c, & r < R_c \\ g_m, & R_c < r \leq R \\ g_e, & r > R \end{cases}, \quad P(r) = \begin{cases} P_c, & r < R_c \\ P_m, & r \geq R_c \end{cases} \quad (10.42)$$

using (10.34) and (10.9). See also (10.48).

Answer:

$$g_c(r) = \frac{4\pi}{3}G\rho_c r, \quad g_m(r) = \frac{G}{r^2} \left\{ \frac{4\pi}{3}\rho_m (r^3 - R_c^3) + M_c \right\}, \quad g_e(r) = \frac{G}{r^2} (M_m + M_c) \quad (10.43)$$

$$M_c = \frac{4\pi}{3}R_c^3\rho_c, \quad M_m = \frac{4\pi}{3}\rho_m (R^3 - R_c^3) \quad (10.44)$$

$$P_c(r) = P_m(R_c) + \frac{2\pi}{3}G\rho_c^2 (R_c^2 - r^2) \quad (10.45)$$

$$P_m(r) = \frac{2\pi}{3}G\rho_m^2 \left\{ R_m^2 - r^2 + 2 \left(\frac{\rho_c}{\rho_m} - 1 \right) R_c^3 \left(\frac{1}{r} - \frac{1}{R_m} \right) \right\} \quad (10.46)$$

10.5.2 The pressure effect on density

In the previous sections we considered the gravity field of a given mass distribution. For self-gravitating planets of sufficient size **the local density depends on the pressure, through selfcompression** i.e. the compression of the material caused by the planets own gravity field. As we have seen in previous sections the lithostatic pressure depends on the gravity field and the density distribution. It follows that the determination of the density, gravity and pressure are coupled problems that must be solved simultaneously and can not be solved separately. Here we will consider the solution of such coupled problems.

From observations of the average density of surface rocks of some $2.7 \cdot 10^3 \text{ kg/m}^3$ and the known mean density of the Earth $5.5 \cdot 10^3 \text{ kg/m}^3$, it follows that the surface density is less than half the mean Earth value. **The difference between both density values suggests a density increase in the interior which could be related either to different composition at depth**, for example corresponding to a dense metallic core, **and/or the effect of selfcompression in an otherwise homogeneous planet**. Solid state phase transitions of mantle material due to increasing pressure can also explain part of the high mean density value.

From the nineteenth century on, models of the internal density distribution of the earth have been investigated. These models have in common that the radial density distribution is parameterized in a simple way with a small number of parameters, typically two, which are then adjusted to the known data such as the surface density and the Earth's total mass or moment of inertia.

In the following the relation between density, gravity and pressure in a self-gravitating planet will be investigated in a more self consistent way.

For a spherically symmetric density distribution the corresponding magnitude of the gravity accel-

eration vector is given by (10.39),

$$g(r) = |\mathbf{g}(\mathbf{r})| = |g_r(r)| = \frac{4\pi G}{r^2} \int_0^r \rho(r') r'^2 dr' = \frac{Gm(r)}{r^2} \quad (10.47)$$

where $m(r)$ is the mass of a sphere of radius r and $\rho(r)$ is the corresponding radial density profile.

problem: 19. Use (10.47) to show that it is not possible to derive a unique radial mass distribution of a spherically symmetric planet from the observed surface value of the gravity field alone. This can be verified by showing that multiple density profiles exist that produce the same surface gravity. To illustrate this sketch a schematic internal radial profile of the gravity acceleration in a comparison of two spherically symmetric planets of identical mass M and radius R . The first one is a homogeneous planet with density ρ_0 and the second one is a differentiated planet with a uniform high density core $\rho_c = \rho_0 + \delta\rho$ and less dense mantle $\rho_m = \rho_0 - \delta\rho$. Verify that these assumptions correspond to this special case with volume fraction of the core $\phi_c = 1/2$.

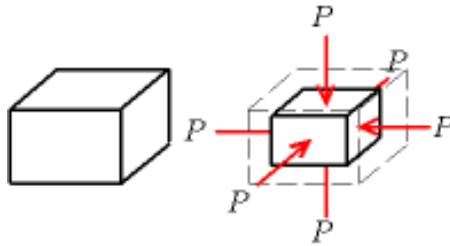
From the above results the lithostatic pressure distribution can be obtained by substitution for the gravity acceleration and integrating the pressure gradient $dP/dr = -\rho g$. Assuming a zero pressure value at the surface this results in,

$$P(r) = \int_r^R \rho(r') g(r') dr' = 4\pi G \int_r^R \rho(r') \left\{ \frac{1}{r'^2} \int_0^{r'} \rho(r'') r''^2 dr'' \right\} dr' \quad (10.48)$$

The pressure in the Earth's interior reaches values over 350 GPa as shown in Fig. ???. For such high pressure values the effect of self-compression on the density is significant. In the following this effect is further explored.

The bulkmodulus An isotropic linear elastic solid can be described by two independent elasticity parameters, for instance the Lamé parameters λ and μ ⁸. The bulkmodulus can be expressed in the Lamé parameters as, $K = \lambda + \frac{2}{3}\mu$. The bulkmodulus K and the shearmodulus μ are the most commonly used parameters to specify the elastic parameters of Earth materials.

The bulk modulus K of a substance measures the substance's resistance to uniform compression. It is defined as the ratio of the infinitesimal pressure increase to the resulting relative decrease of the volume. Its SI unit is the Pascal, and its dimensional form is $M^1 L^{-1} T^{-2}$.



The incompressibility K , or bulkmodulus, is defined as,

$$\frac{1}{K} = \frac{1}{\rho} \frac{d\rho}{dP} \quad (10.49)$$

By substitution of $dP = -\rho g dr$ in (10.49) we derive a differential equation for the density profile of a compressible planet model,

$$\frac{1}{K} = \frac{-1}{\rho^2 g} \frac{d\rho}{dr} \Rightarrow \frac{d\rho}{dr} = -\frac{\rho^2 g}{K} \quad (10.50)$$

⁸https://en.wikipedia.org/wiki/Elastic_modulus

Parameterization of the bulkmodulus The radial density distribution for a selfcompressing planet can be obtained from (10.50) once the bulkmodulus K is known. We will first consider simple cases where K is either a uniform constant or it is parameterized in terms of the density.

problem: 20. Assume both K and g in (10.50) to be uniform in the mantle and derive the following density profile,

$$\rho(z) = \frac{\rho_0}{1 - \frac{\rho_0 g z}{K}} \quad (10.51)$$

where $z = R - r$ is the depth coordinate and $\rho_0 = \rho(0)$ is the surface density value.

- Compute the depth z_1 where the expression (10.51) becomes singular, i.e. $\rho \rightarrow \infty$, suggesting infinite compression of the material. To do this assume Earth(mantle)-like values of the incompressibility, $K = 400\text{GPa}$ (see Fig.??) and the surface density $\rho_0 = 3 \cdot 10^3 \text{ kg/m}^3$.
- Now consider a simplified model of a large rocky exoplanet of Earth-like composition with $M = 8M_\oplus$ and $R = 1.5R_\oplus$. Assume uniform gravity (adapted for the given M, R) and uniform incompressibility K . Do you now find the singular depth z_1 within the depth range of the planet? Comment on the assumption of a uniform gravity field in view of the models presented in section 10.5.1.

problem: 21. The result of problem 20 gives the density depth distribution for the model with constant properties. The resulting expression (10.51) also contains the uniform gravity acceleration. A more fundamental relation between density and pressure, not including gravity, can be derived for this model with constant material property K as an equation of state (EOS) for the density.

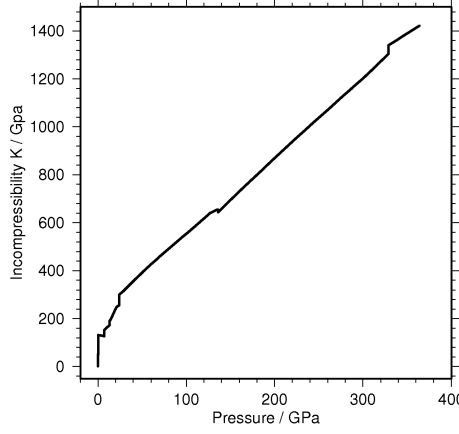
Derive from the definition of the bulkmodulus (10.49) the following logarithmic EOS for the density in terms of the static pressure,

$$P = \ln \left(\left(\frac{\rho}{\rho_0} \right)^K \right) \quad (10.52)$$

Show that the above EOS (10.52) can be inverted to obtain an explicit expression for density as a function of pressure.

The singular behavior in the density model of problem 20 is a result of the assumed uniform g and K in (10.50). While g is reasonably constant with depth in the mantle, as illustrated in Fig. ??, K is not. The incompressibility increases with increasing depth/pressure and as a result the compression remains finite for earth-like conditions. The incompressibility can be expressed in the density and the seismic wave velocities, $v_p = \sqrt{(\lambda + 2\mu)/\rho}$, $v_s = \sqrt{\mu/\rho}$. With $K = \lambda + \frac{2}{3}\mu$ this becomes $K = \rho(v_p^2 - 4/3v_s^2)$. A radial profile $K(P(r))$ can therefore be derived, from the seismic velocities determined from inversion of traveltime tables of longitudinal and shearwave seismic arrivals.

The $K(P(r))$ profile derived from the PREM model of Dziewonski and Anderson (1981) appears to be roughly linear as shown in the following figure:



Incompressibility profile derived from the PREM model.

A linear relation between bulkmodulus and pressure as suggested by this figure is also obtained using the following power law parameterization for the bulkmodulus in terms of the density $K(\rho)$.

$$K = C\rho^n \Rightarrow \ln(K) = \ln(C) + n \ln(\rho) \Rightarrow n = \frac{d \ln(K)}{d \ln(\rho)} = \frac{dK}{dP} = K'_0 \quad (10.53)$$

where C is a constant. The constant pressure derivative in this model implies a linear pressure relation $K(P) = K_0 + K'_0 P$. This appears to approximate the distribution of K in particular in the lower mantle as determined from seismological data in the PREM model. $K'_0 \approx 4$ for the magnesium-iron silicates (Mg, Fe)SiO₃ (perovskite) and dense oxides (Mg, Fe)O (wüstite), representative for the earth's deep mantle.

The Murnaghan e.o.s. An equation of state directly relating the density or specific volume, $V = 1/\rho$, to pressure can be derived from such an 'ansatz' of a linear pressure dependence $K = K_0 + K'_0 P$ as shown in the following,

$$\frac{1}{\rho} \frac{d\rho}{dP} = \frac{1}{K} \rightarrow \frac{1}{V} \frac{dV}{dP} = -\frac{1}{K} \rightarrow dP = -(K_0 + K'_0 P) \frac{1}{V} dV \quad (10.54)$$

$$\int_0^P \frac{dP'}{K_0 + K'_0 P'} = - \int_{V_0}^V \frac{1}{V'} dV' = \int_V^{V_0} \frac{1}{V'} dV' = \ln \left(\frac{V_0}{V} \right) \quad (10.55)$$

Substitution in the integral over pressure of $K_0 + K'_0 P' = x$, $dx = K'_0 dP'$ gives,

$$\int_{x_0=K_0}^{x_P=K_0+K'_0 P} \frac{1}{K'_0} \frac{dx}{x} = \frac{1}{K'_0} \ln \left(\frac{K_0 + K'_0 P}{K_0} \right) = \ln \left(\frac{V_0}{V} \right) \quad (10.56)$$

$$1 + \frac{K'_0 P}{K_0} = \left(\frac{V_0}{V} \right)^{K'_0} \rightarrow P = \frac{K_0}{K'_0} \left(\left(\frac{V_0}{V} \right)^{K'_0} - 1 \right) \quad (10.57)$$

This relation is known as the Murnaghan equation of state (EOS).

The Murnaghan equation of state is a relationship between the volume of a body and the pressure to which it is subjected. This is one of many state equations that have been used in earth sciences and shock physics to model the behavior of matter under conditions of high pressure. It owes its name to Francis D. Murnaghan who proposed it in 1944 to reflect material behavior under a pressure range as wide as possible to reflect an experimentally established fact: **the more a solid is compressed, the more difficult it is to compress further.**

The Murnaghan equation is derived, under certain assumptions, from the equations of continuum mechanics. It involves two adjustable parameters: the modulus of incompressibility K_0 and its

first derivative with respect to the pressure, K'_0 , both measured at ambient pressure. In general, these coefficients are determined by a regression on experimentally obtained values of volume V as a function of the pressure P . These experimental data can be obtained by X-ray diffraction or by shock tests. Regression can also be performed on the values of the energy as a function of the volume obtained from ab-initio and molecular dynamics calculations.

problem: 22. Derive an explicit expression for the pressure dependent density from the Murnaghan equation of state (10.57).

Answer:

$$\rho(P) = \rho_0 \left(\frac{K'_0 P}{K_0} + 1 \right)^{1/K'_0} \quad (10.58)$$

problem: 23. In problem 20 we have seen that a simple model with uniform incompressibility and gravity $K = K_0$ and $g = g_0$ leads to physically impossible solutions. In a refined version of this model, applied to the Earth's mantle, $g = g_0$ is maintained (compare Fig.??), and K is parameterized using the powerlaw relation (10.53).

Derive the following density profile for the model corresponding to (10.53).

$$\rho(r) = \rho_0 \left(1 + (n-1) \frac{\rho_0 g_0 z}{K_0} \right)^{\frac{1}{n-1}} \quad (10.59)$$

where $z = R - r$ is the depth coordinate and the 0 subscript refers to zero pressure conditions.

Note that the singularity for $\rho_0 g_0 z / K_0 = 1$ in problem 20 is absent in this model.

A more widely used and more accurate EOS for a higher pressure range is the equation derived by Birch (1952) from a consideration of elastic strain energy, known as the Birch-Murnaghan EOS (Poirier, 2000).

In other cases than the special simplified cases discussed above, in particular in problems 20 and 23, the gravity acceleration varies also with depth. Also more accurate equations of state may be necessary for very high pressure, encountered in the deep interior of large (exo)planets, that result in large compression. Such models can be formulated in a more general way by the following coupled set of equations for pressure, gravity and density.

$$\frac{dP}{dr} = -\rho g \quad (10.60)$$

$$g(r) = \frac{Gm(r)}{r^2} \quad (10.61)$$

$$F(\rho, P, T) = 0 \quad (10.62)$$

where the radial mass distribution $m(r)$ is defined as in (??). A model based on (10.60), (10.61), and (10.62) can be constructed for the internal structure (density, gravity, pressure) of a planet of given mass M and composition, i.e. with given parameters of the EOS (10.62) such as ρ_0, K_0, K'_0 in the Murnaghan EOS (10.57). Consider the application of such a model to a planet for which only the planet mass M is known.⁹ Assume a homogeneous terrestrial (rocky) planet without a distinct metallic core. Assuming an earth-mantle like composition, representative values of the EOS parameters can be used, to solve the coupled model equations in the following iterative scheme.

1. Define a grid along the radial coordinate $r_i, i = 1, \dots, N, r_1 = 0$. This grid defines a subdivision of the interior in $N - 1$ concentric layers and must be chosen large enough, i.e. $r_N > R$.

⁹Such models can be applied to exoplanets that are recently being discovered https://en.wikipedia.org/wiki/Methods_of_detecting_exoplanets. For some of these planets, detected from radial velocity variations of the star, only the planet mass M is known.

2. Choose an initial estimate of the central pressure $P^{(1)}(0)$.
3. In a loop over the internal layers, starting upward from the centre, first compute the pressure decrement over the layer from (10.60). This is then used to obtain the pressure at the next grid point and corresponding density from the EOS (10.62). From the computed density the corresponding mass distribution $m(r_i)$ and gravity $g(r_i)$ (10.61) follow.
4. The layer iteration in the previous item is stopped when a zero pressure value has been reached. The radial level reached this way now defines the next approximation of the planetary radius $R^{(j)}$ and $M^{(j)} = m(R^{(j)})$ is a new approximation of the planet mass M .
5. From the total mass defect $\Delta M^{(j)} = M^{(j)} - M$ a correction to the central pressure is computed as $\Delta P^{(j)}$, (problem 24). In the next iteration the radial integration is repeated from item 3 with an updated central pressure $P^{(j+1)}(0) = P^{(j)}(0) + \Delta P^{(j)}$ and this iterative procedure is repeated until convergence is reached, i.e. until $|\Delta M^{(j)}|/M$ drops below a specified tolerance value.

problem: 24. A correction for the central pressure in item 4 can be estimated by distributing the mass defect $\Delta M^{(j)}$ over a spherical shell of thickness $\Delta R^{(j)}$, positioned at the surface, and computing an approximate pressure $\Delta P^{(j)}$ at the bottom of this shell. Derive the following expression for the thickness of this spherical shell,

$$\frac{\Delta R^{(j)}}{R^{(j)}} = \left(\frac{\Delta M^{(j)}}{M^{*(j)}} + 1 \right)^{1/3} - 1 \quad (10.63)$$

Where $M^{*(j)} = \frac{4\pi}{3} \rho(R^{(j)}) R^{(j)3}$.

The correction for the central pressure is then defined as, $\Delta P^{(j)} = \rho(R^{(j)}) g(R^{(j)}) \Delta R^{(j)}$.

10.5.3 Adiabatic density distribution

In the previous section density models were based on assumptions about the parameterization of the bulkmodulus K . The density model of Williamson and Adams (1923), (Hemley, 2006) does not depend on a parameterized K . Instead it is defined in terms of the seismic wave velocities v_p and v_s that can be determined from inversion of seismological traveltime data as $K/\rho = v_p^2 - 4/3 v_s^2$.

The W-A model can be derived from thermodynamic principles for a homogeneous self-compressing layer which is in an adiabatic state. The bulkmodulus applied in this model is expressed in the seismic wave velocities which in turn depend on the elasticity parameters and the density. The elastic deformation process in seismic wave propagation occurs on a relatively short time scale (seconds-minutes) compared to the characteristic time scale of conductive heat transport in solids (see ??). Therefore (diffusive) heat exchange can be neglected and adiabatic conditions apply in seismic wave propagation. This implies that the elasticity parameters determined from seismic data, including the bulkmodulus K pertain to adiabatic conditions (see also Appendix ??).

Other processes such as convective mantle flow that occur on a much longer time scale may take place under more general (non-adiabatic) conditions.

In section ?? on the thermal state of the Earth it is shown that adiabatic conditions hold for the interior of a fluid layer when heat transport is dominated by advection and heat diffusion by conduction/radiation plays a minor role. Assuming the Earth's mantle to be in a state of vigorous thermal convection it also follows that the average temperature profile, the geotherm, corresponds to an adiabatic distribution.

In general the density differential can be written as,

$$d\rho = \left(\frac{\partial \rho}{\partial P} \right)_S dP + \left(\frac{\partial \rho}{\partial S} \right)_P dS \quad (10.64)$$

where the differential of the entropy S is dropped in case of adiabatic conditions and the pressure derivative is written in terms of the adiabatic bulkmodulus K_S defined in (10.49), $1/K_S = (\partial \rho / \partial P)_S / \rho$.

problem: 25. Derive the Williamson-Adams equation for a homogeneous adiabatic layer from the density differential (10.64) and assumption of isentropic (adiabatic) conditions with $dS \equiv 0$,

$$\frac{d\rho}{dr} = -\frac{\rho^2 g}{K_S} \quad (10.65)$$

The density solution of the W-A equation can be expressed in terms of the seismic parameter $\Phi = K_S / \rho$ which in turn can be obtained from seismic velocity models: $\Phi = v_p^2 - \frac{4}{3}v_s^2$ for P and S waves. $\sqrt{\Phi} = \sqrt{K_S / \rho}$ is known as the bulkvelocity. For a given bulkvelocity profile, obtained from seismic observations, the W-A density profile is derived from (10.65) as,

$$\ln \left(\frac{\rho(r)}{\rho(R)} \right) = \int_r^R \Phi^{-1}(r') g(r') dr' \quad (10.66)$$

problem: 26. Derive (10.66) by integration of the W-A equation (10.65).

In (10.66) the gravity acceleration g depends on the density distribution $\rho(r)$ in the lefthand side. Therefore the density profile can not be simply obtained from a seismologically determined $\Phi(r)$ profile and a single evaluation of the integral in (10.66). The expression represents an integral equation that can be solved iteratively as specified in problem 27.

problem: 27. Assume that a seismic parameter profile for the mantle $\Phi(r)$, obtained from seismic travel times, is available. Investigate how (10.66) can be used to compute a sequence of mantle density profiles $\rho^{(j)}(r)$, $j = 1, 2, \dots$ in an iterative procedure, by successive substitution. How would you define a starting profile $\rho^{(1)}(r)$ for this iterative procedure?
Hint: Substitute the density profile for iteration number j in the gravity acceleration in the right-hand side of (10.66) for the computation of an updated profile $j + 1$. This is an example of a general solution strategy for non-linear problems known as ‘successive substitution’ or Picard iteration.

Williamson and Adams (1923) [1363] used the iterative scheme in problem 27 to test the hypothesis that the mass concentration towards the Earth’s centre is completely explained by compression of a homogeneous self-gravitating sphere. They showed that integrating (10.66) from a surface value of $3.3 \cdot 10^3 \text{ kg/m}^3$ results in unrealistically high density values for depths greater than the core-mantle boundary. This way they concluded that an inhomogeneous earth with a dense, compositionally distinct core, probably iron-nickle, was required by the observations. The necessary multiple integrals in the evaluation of (10.66) had to be computed by means of graphical approximation methods in 1923, several decades before the advent of electronic computers.

In a later analysis Bullen (1936) showed that the assumption of a homogeneous selfcompressing mantle described by the W-A equation, and a chemically distinct dense core, leads to unrealistically high values of the moment of inertia for the core $I_c = f M_c R_c^2$, with a prefactor value $f \sim 0.57$ greater than the value of a core with uniform density, 0.4. Since this would imply a density decrease towards the centre Bullen concluded that the applicability of the W-A model for the whole mantle

can not be maintained and that instead a distinct mantle transition layer, labeled C-layer, must be included between the upper and lower mantle proper, related to transitions in mineral phase and/or composition (Bullen, 1975).

problem: 28.

1. Derive the following equation for the temperature distribution of a W-A layer (see Appendix ??),

$$\frac{dT}{dr} = -\frac{\alpha g}{c_P} T \quad (10.67)$$

where α and c_P are the thermal expansion coefficient and the specific heat at constant pressure.

Hint: Use the differential for the entropy,

$$dS = \left(\frac{\partial S}{\partial T} \right)_P dT + \left(\frac{\partial S}{\partial P} \right)_T dP \quad (10.68)$$

and the thermodynamic relations: $(\partial S / \partial T)_P = c_P / T$ and $(\partial S / \partial P)_T = -\alpha / \rho$.

2. Derive the expression for the temperature profile for an adiabatic layer, sometimes referred to as the ‘adiabat’, by solving equation (10.67),

$$T(r) = T(R) \exp \left(\int_r^R \frac{\alpha g}{c_P} dr' \right) \quad (10.69)$$

The temperature extrapolated to the surface, $T_P = T(R)$ is known as the potential temperature of the layer. The quantity $H_T = (\alpha g / c_P)^{-1}$ is known as the thermal scale height of the layer.

3. Derive an expression from (10.69) for the special case with a constant value of the scale height parameter.

The W-A equation for the density of an adiabatic layer can be generalized introducing the Bullen parameter η which is used as a measure of the departure of the actual density/temperature profile from an adiabat. This is done by writing,

$$\eta(r) = -\frac{\Phi}{\rho g} \frac{d\rho}{dr} \quad (10.70)$$

where $\eta(r)$ has been substituted for the constant value ($\equiv 1$) in the W-A equation.

Current density models

The concept of an adiabatic layer was essential when no independent determinations for the density distribution were available and the W-A equation was used to compute $\rho(r)$ for given values of the seismic parameter $\Phi(r)$ determined from seismological observations (Bullen, 1975).

During the 1970s a radial density distribution has been obtained for the Earth from inversion of seismological observations, incorporating spectral analysis of the Earth’s eigenvibrations, under the constraints of the given values for M and I . This, together with seismic velocities determined from bodywave traveltimes and surfacewave dispersion, has resulted in the Preliminary Reference Earth Model (PREM), (Dziewonski and Anderson, 1981 [357]).

Since $\rho(r)$ can be determined from analysis of the earth’s normal modes (radial eigenvibrations) the ‘adiabaticity’ of the mantle is no longer assumed.

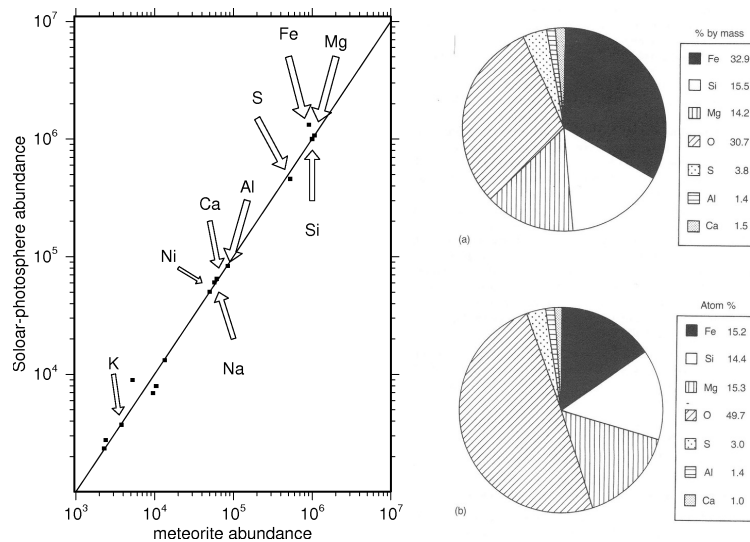
The degree of ‘adiabaticity’ is used in numerical modelling experiments as a diagnostic for the dynamic state - where a high degree of adiabaticity indicates vigorous thermal convection and predominantly convective heat transport (van den Berg and Yuen, 1998, [1305] Matyska and Yuen, 2000, [843] Bunge *et al.*, 2001).

Usually the outcome of such experiments shows that the upper and lower mantle separately are approximately adiabatic - away from boundary layers where conductive transport dominates. In recent years models of the deep lower mantle have become popular where a compositionally distinct dense layer occupies the bottom 30% (roughly) of the lower mantle (Kellogg *et al.* (1999) [693], Albaredo and van der Hilst (2002) [4]).

10.5.4 Earth’s chemical composition

For a complete description of the Earth’s interior we need to know its chemical composition, temperature and pressure. In section 10.4 the pressure is expressed in the density distribution and the related internal gravity field. Once the internal pressure distribution is known, sharp transitions or discontinuities in the material properties, like the seismic velocities v_p , v_s and the density in the PREM model, can be identified with mineral phase transitions and as such they can be related to the mineral (P, T) phase diagram of candidate mantle silicate materials in order to estimate the temperature in the Earth’s interior. Such phase diagrams are determined from experimental (HPT) and theoretical work in mineral physics.

What do we know about Earth’s bulk chemical composition? Candidate mantle materials have been defined based on cosmochemical and petrological considerations. Models of the chemical composition of the Earth are commonly based on the hypothesis that the planet was formed in a multi-stage accretion process from material that condensed from the original solar nebula approximately 4.6 billion years ago at the time of formation of the solar system. The chemical composition of chondritic meteorites, in particular the carbonaceous chondrites (CI type) (McBride and Gilmour (2003) [849]) show a strong correlation with the composition of the outer layer of the sun (photosphere), determined from spectral analysis of the solar light, as illustrated in the following figure:



Left: Element abundance (normalized with $Si = 10^6$), of the solar shallow photosphere compared to chondritic meteorites (Anders & Grevesse (1989) [18]).

Right: amounts of Earth’s major elements assuming a chondritic composition (Brown & Musset (1993) [155]).

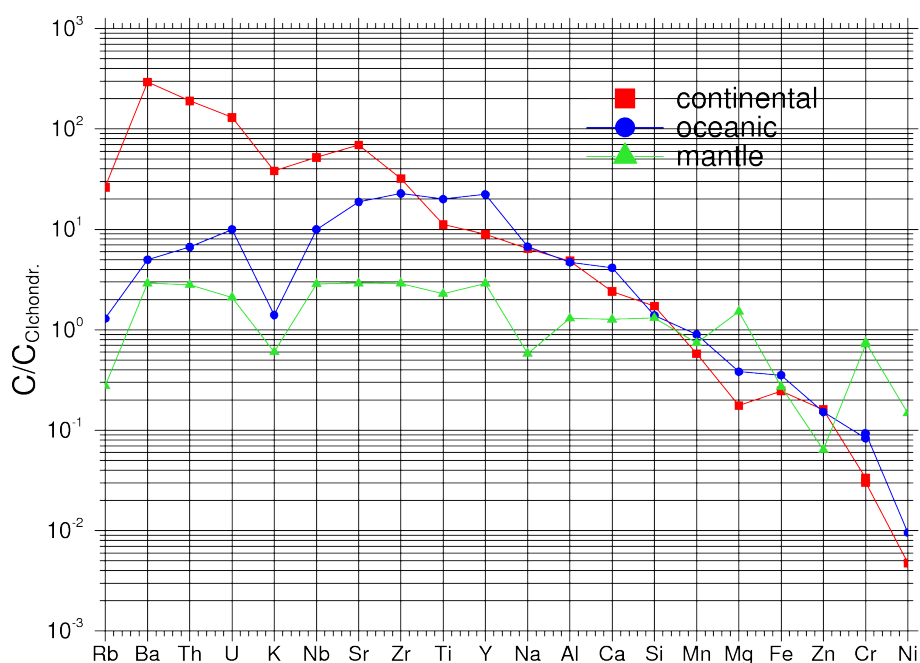
The solar-chondritic data in the lefthand frame show that Mg, Fe and Si are by far the most abundant (non-volatile) elements. According to the chondritic Earth hypothesis a similar abundance can be expected for the bulk-earth. This is illustrated in the righthand pie diagrams. Note the large proportion of oxygen, bound in oxides. In most crust-mantle rocks S is less abundant than Al or Ca. This is usually explained by assuming that S is relatively volatile and also ‘siderophile’, meaning

that a significant fraction may have ended up in the iron-nickel core during an early core-mantle differentiation.

The chondritic meteorites are thought to be representative of the undifferentiated material condensed from the solar nebula.

Around 1960 a model chemical composition for the bulk of the Earth's mantle, coined pyrolite, was introduced by Ringwood (see (Ringwood, 1975) and original references therein). This is still used as a reference model. The pyrolitic composition is associated with the main upper mantle rock type peridotite that is brought to the Earth's surface in small fragments included in volcanic rocks (xenoliths) and also in larger, kilometer sized, fragments in so called peridotite bodies (Spengler *et al.*, 2006). The pyrolitic composition of the upper mantle rocks is also strongly correlated with the composition of chondritic meteorites, in agreement with the hypothesis of a chondritic origin of the Earth.

Mantle peridotites are found with different degrees of depletion (mass fraction lost) by partial melting. More depleted material is denoted as harzburgite and the relatively undepleted peridotite is known as lherzolite. During progressive partial melting the mineral composition of the residual rock material, a mineral assemblage consisting of olivine, pyroxene and garnet, shifts towards the olivine composition. The olivine enriched harzburgitic residue appears to be the chemical complement of the basaltic melt product, with respect to the original lherzolitic mantle source rock. This depletion relation, between oceanic and continental crust on the one hand and peridotitic mantle rock on the other, is reflected in the element abundance of crust and mantle rocks, illustrated in the following figure:



Chemical abundance of crustal and mantle rocks, normalized with respect to CI chondritic values. Data from (McBride and Gilmour, 2003 [849]).

This figure shows abundance ratio's relative to the CI-chondritic composition. The curve for mantle rock appears to be relatively close to the chondritic composition, whereas the crustal material is enriched with respect to the mantle in most elements shown.

A notable exception to this crustal enrichment is found for magnesium which appears to be enriched in average mantle peridotite. This is in agreement with the previous observation that the olivine/pyroxene content ratio of the residual increases with the degree of partial melting. Magnesium content increases with the olivine (forsterite Mg_2SiO_4)/pyroxene MgSiO_3 ratio.

An other observation that can be made from the figure above is the apparent depletion of the siderophile elements Fe and Ni, both in crust and mantle material, with respect to the chondritic composition. This is usually explained by the formation of a liquid Fe, Ni rich metal core of the

Earth during the first few million years of the accretion process, in the early solar system. During this event the molten liquid metal would have differentiated from the silicate mantle, leaving the mantle depleted in siderophile (iron loving) elements.

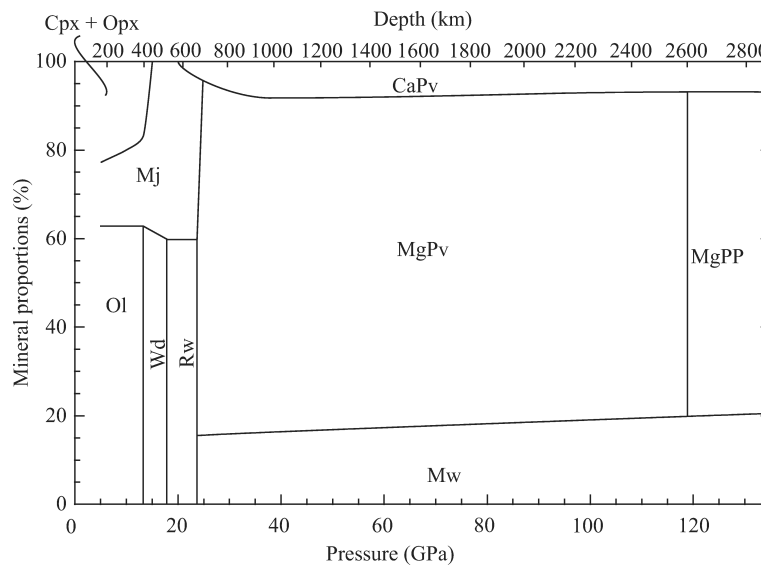
Core formation is also sometimes used as an explanation of the apparent K (potassium) depletion of both mantle and crust with respect to chondrites. In this explanation K is dissolved in liquid iron in significant quantity at high pressure and temperature (Rama Murthy *et al.*, 2003). An alternative explanation for the Earth's K depletion is an escape of K due to significant volatilization during the planetary accretion process.

problem: 29. From Figure ?? it can be concluded that the Earth's mantle and crust lost roughly 2/3 of its original iron content corresponding to a chondritic composition. Verify how this iron-depletion of crust and mantle could be explained by differentiation of the Earth's mostly-iron core. Use the following data in your argument: a) The mass fraction of the core $X_c = M_c/M_\oplus = 0.315$. b) The Fe mass fraction $X_{mFe} \sim 10\%$ of the pyrolitic mantle, c) The mass fraction of lighter elements in the core - (S, Si, O) amounts to about 20%. d) The Fe mass fraction of the bulk Earth $X_{\oplus Fe} \sim 33\%$ (Fig. ??)

10.5.5 Phase transitions as anchor points of the geotherm

Major phase boundaries in the Earth's mantle and core have been identified with sharp transitions in the seismic wave velocities and the density distribution of the PREM model.

The depth distribution of the mineral composition for a pyrolitic mantle model is shown in the following figure:



Pressure/depth distribution of mineral assemblage for a pyrolitic mantle model. Cpx: clinopyroxene, Opx: orthopyroxene, Mj: majorite garnet, Ol: olivine, Wd: wadsleyite, Rw: ringwoodite, CaPv: $CaSiO_3$ perovskite MgPv: $MgSiO_3$ -rich perovskite, MgPP: $MgSiO_3$ -rich post-perovskite, Mw: magnesiowüstite.

(From: (Hirose, 2007))

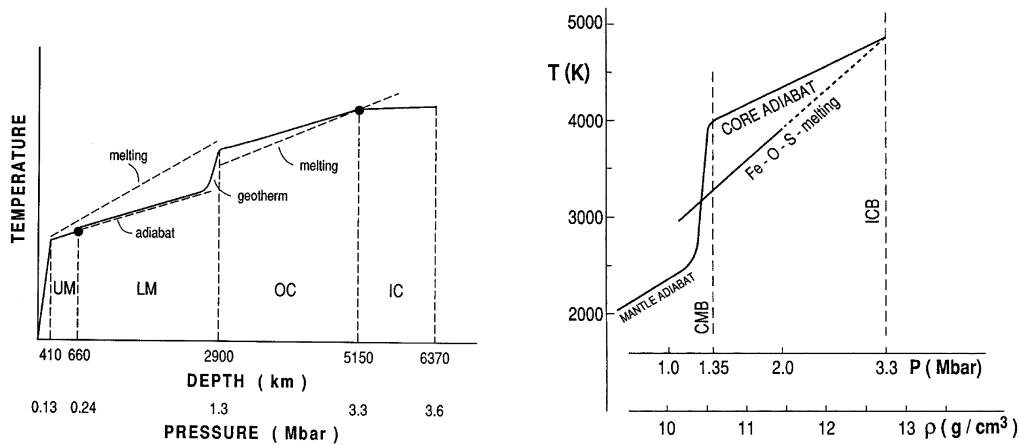
This figure clearly illustrates the different mineral composition of the upper and lower mantle regions separated by the major phase boundary near 660 km depth (~ 24 GPa), where the ringwoodite polymorph of olivine, $(Mg, Fe)_2SiO_4$, transforms (dissociates) into a mineral assemblage of perovskite, $(Mg, Fe)SiO_3$ and magnesiowüstite, $(Mg, Fe)O$.

For a given mantle composition, for instance for a pyrolitic mantle, the pressure-temperature mineral phase diagram can be determined for the relevant P, T range of the Earth's mantle by HPT experiments and mineral physics theory. A sharp transition at a pressure P_t in the PREM model can then be located at the corresponding pressure in the phase diagram by the intersection of the P_t isobar with the diagram phase boundaries. The (possibly multiple) intersection points

define the corresponding transition temperature T_t . The pressure-temperature point located in the phase diagram defines an ‘anchor point’ that constrains the geotherm. In this procedure the phase transition is used as a mantle/core thermometer.

This way several (P, T) ‘anchor points’ of the geotherm have been determined, related to the solid state phase transition near 660 km depth and the solid/liquid inner/outer core boundary at 1220 km from the Earth’s centre.

The following figure from Boehler (1996) [106]) illustrates the determination of anchor points of the geotherm at the phase boundary near 660 km depth ($P_{660} = 24\text{GPa}, T_{660} = 1900 \pm 100\text{ K}$) and at the boundary between the outer and inner core at 5150 km depth, ($P_{ICB} = 330\text{GPa}, T_{ICB} = 4850 \pm 200\text{ K}$).



Schematic radial temperature distribution in the mantle and core, constrained by major phase transitions (Boehler, 1996), (UM-upper mantle, LM lower mantle, OC outer core, IC inner core). The temperature of the upper/lower mantle boundary is constrained by the γ -spinel to postspinel phase transition at 660 km depth. The temperature at the inner/outer core boundary at 5150 km depth (radius 1220 km) is constrained by the melting temperature of the hypothetical core ‘Fe-O-S’ alloy. The right hand frame shows a schematic core temperature distribution (geotherm) labeled ‘CORE ADIABAT’ in the liquid outer core versus pressure and the melting curve (liquidus) of the core ‘Fe-O-S’ alloy. (CMB core-mantle boundary, ICB inner core boundary). The ICB is determined by the intersection of the liquidus and the geotherm. During core cooling the ICB moves outward as the inner core grows by crystallisation.

Starting from these anchor points the temperature is then extrapolated from both sides to the core mantle boundary at 2900 km depth. For this temperature extrapolation assumptions have to be made about the dominant heat transport mechanism and in this case it is assumed that heat transport operates mainly through thermal convection. This will be further investigated in later sections dealing with heat transport in the Earth’s mantle.

problem: 30. Estimate the temperature near the bottom of the mantle by adiabatic extrapolation of the temperature $T_{660} \sim 1900\text{K}$ of the phase transition near 660 km depth, to the depth of the core mantle boundary, using the general expression for the adiabat in a homogeneous layer.
Hints: apply the result of problem 28 and assume uniform values of the ‘scale height parameter’ $H_T = (\alpha g/c_p)^{-1}$, with $\alpha = 2 \cdot 10^{-5}\text{K}^{-1}$, $g = 10\text{ms}^{-2}$, $c_p = 1250\text{Jkg}^{-1}\text{K}^{-1}$. Further: approximate the adiabat by a linear depth function, in agreement with the schematic diagram of Boehler (1996) - see figure above -to obtain a uniform adiabatic temperature gradient.

The ‘head’ of the extrapolated outer core adiabat is at a temperature of approximately 4000 K and the ‘foot’ of the lower mantle adiabat at approximately 2700 K. This result indicates a large temperature contrast of about 1300 K across the CMB.

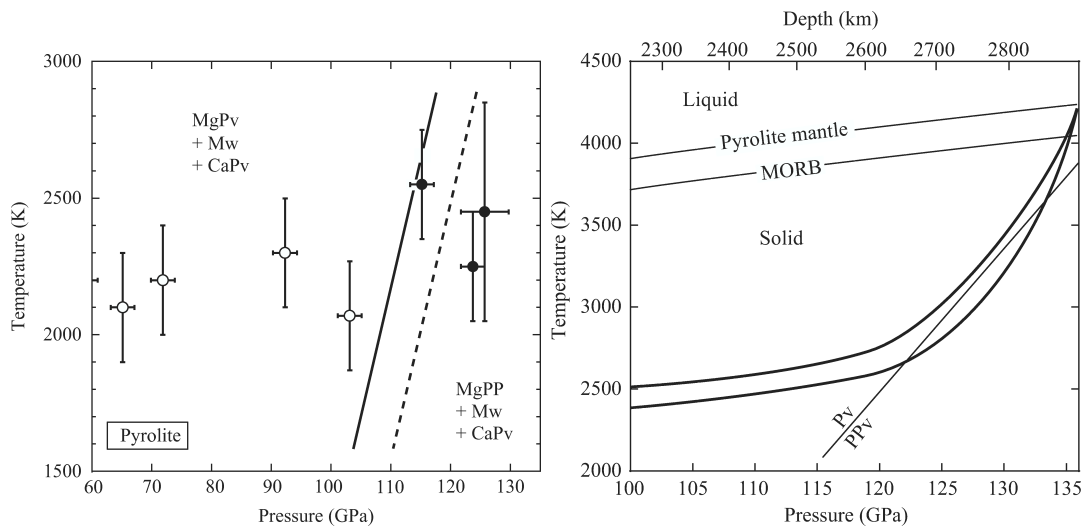
How can such a large contrast be explained physically? As we will see later, this can be explained by interpreting the CMB as a boundary between two separately convecting fluid layers, each with a thermal boundary layer where the main heat transport mechanism shifts from convection in the interior of the fluid layers, to conduction near the boundary interface, where vertical convective transport vanishes with the flow velocity component normal to the boundary. Separately convecting

layers are in agreement with the large density contrast across the CMB where the density almost doubles, as illustrated in the PREM profile. The resulting strong temperature contrast across the CMB is consistent with a lower mantle in a state of vigorous thermal convection.

problem: 31. *Explain why we can not turn this argument around and conclude from these indications for a strong temperature contrast at CMB that the mantle convects vigorously.*

Hint: Check Appendix ?? for the assumptions made for an adiabatic geotherm in the lower mantle.

More recent developments, providing independent information, shed new light on the temperature distribution in the bottom layer of the lower mantle. A previously unknown mantle phase transition has been identified, in the main constituent magnesium-perovskite, to a ($\sim 1.5\%$) denser phase (post-perovskite) both in experimental HPT and theoretical (mineral physics) work at temperatures and pressure conditions corresponding to a region in the lowermost mantle close to the core-mantle boundary. This is illustrated in the figure hereafter showing experimental data points delineating the phase boundary.



Left: phase relations near the bottom of the mantle for pyrolitic material (Hirose, 2007). The solid- and dashed line correspond to different pressure calibration of the HPT experiments. The Clapeyron slope of the phase boundary is assumed 11.5 MPa/K . CaPv: CaSiO_3 perovskite MgPv: MgSiO_3 -rich perovskite, MgPP: MgSiO_3 -rich post-perovskite, Mw: magnesio-wüstite. Right: schematic temperature profiles in the lower mantle in relation to the perovskite (PV) to postperovskite (PPV) phase transition and the melting curve for pyrolitic mantle material and subducted basaltic crust (MORB) (Hirose *et al.* (2007) [573]).

This phase transition has a high valued positive slope of the phase boundary (Clapeyron parameter) $dP_t/dT \sim 10 \text{ MPaK}^{-1}$. The intercept of the phase boundary with the core mantle boundary at $\sim 136 \text{ GPa}$ appears to be at a temperature several hundred Kelvin below the temperature of the liquid metal outer core as illustrated in the right part of the above figure. As a consequence the geotherm may intersect the phase boundary at multiple depth's, depending on the local mantle temperature, a phenomenon known as 'double crossing' (Hernlund *et al.* ., 2005). When a double crossing of the geotherm occurs, a thin layer exists directly bordering the core, where perovskite is the stable phase while on top of this bottom PV layer, a postperovskite layer exists with a variable thickness of up to several hundred kilometers.

A further implication of the phase diagram illustrated in is that the PPV layer will be absent in hot regions where the geotherm is completely above the PV-PPV phase boundary. This post-perovskite phase boundary has also been associated with the top of the D'' layer at variable height $\sim 100 - 300 \text{ km}$ above the CMB (Lay *et al.* ., 2005).

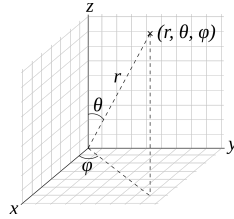
These seismological interpretations of the postperovskite phase boundary have been based on limited resolution methods applying 1-D radial velocity models. In a more recent development, techniques related to seismic wave migration methods, used in the oil and gas exploration industry,

are applied to delineate reflecting interfaces in 2-D and 3-D models in seismic stratigraphy of the CMB region (van der Hilst *et al.* , 2007). This way the spatial resolution has been brought down to about 20 km, allowing mapping of detailed structures in the lowermost mantle. An important target of these high resolution seismic methods is the bottom interface of a postperovskite layer, associated with the ‘double crossing’, where mantle material transforms back from postperovskite into perovskite due to the steep increase in temperature in the bottom thermal boundary layer, illustrated in the figure above, related to the temperature contrast across the CMB.

In a similar way as for the spinel-postspinel phase transition the temperature at the seismic interfaces can then be estimated from the given depth(pressure) and the experimentally determined parameters of the postperovskite phase transition. This way a mantle adiabatic geotherm and boundary layer structure (error function) have been estimated with a CMB temperature $T_{cmb} \sim 4000$ K (van der Hilst *et al.* , 2007). The ‘foot’ of the adiabatic mantle geotherm derived from this lies at a temperature of approximately 2500 K. Both the estimated CMB temperature and the foot of the adiabat seem to confirm independent earlier findings based on adiabatic temperature extrapolation over large depth ranges (Boehler, 1996 [106]).

The temperature contrast of about 1500 K across the core-mantle boundary resulting from these interpretations identify the bottom of the mantle as a thermal boundary layer, characteristic of a vigorously convecting layer where the boundary interface has a fixed or slowly varying temperature, as we will see in the section on heat transport in the mantle. As such these results from mineral physics and seismology have produced new evidence for strong mantle convective flow near the core-mantle boundary.

10.6 geostationary orbit



In spherical coordinates¹⁰ the position, velocity and acceleration of a point are given by

$$\begin{aligned}\vec{r} &= r\vec{e}_r \\ \vec{v} &= \dot{r}\vec{e}_r + r\dot{\theta}\vec{e}_\theta + r\dot{\phi}\sin\theta\vec{e}_\phi \\ \vec{a} &= (\ddot{r} - r\dot{\theta}^2 - r\dot{\phi}^2\sin^2\theta)\vec{e}_r \\ &\quad + (r\ddot{\theta} + 2\dot{r}\dot{\theta} - r\dot{\phi}^2\sin\theta\cos\theta)\vec{e}_\theta \\ &\quad + (r\ddot{\phi}\sin\theta + 2\dot{r}\dot{\phi}\sin\theta + 2r\dot{\theta}\dot{\phi}\cos\theta)\vec{e}_\phi\end{aligned}$$

For an orbit at constant height ($r = R$, $\dot{r} = 0$), and constant angular velocities (i.e. $\ddot{\theta} = 0$ and $\ddot{\phi} = 0$) we arrive at

$$\vec{v} = R\dot{\theta}\vec{e}_\theta + R\dot{\phi}\sin\theta\vec{e}_\phi \quad (10.71)$$

$$\vec{a} = (-R\dot{\theta}^2 - R\dot{\phi}^2\sin^2\theta)\vec{e}_r + (-R\dot{\phi}^2\sin\theta\cos\theta)\vec{e}_\theta + (2R\dot{\theta}\dot{\phi}\cos\theta)\vec{e}_\phi \quad (10.72)$$

If the orbit is the equatorial plane, we have $\theta = \pi/2$, $\sin\theta = 1$, $\cos\theta = 0$ (and of course $\dot{\theta} = 0$) so now

$$\vec{v} = R\dot{\phi}\vec{e}_\phi \quad (10.73)$$

$$\vec{a} = -R\dot{\phi}^2\vec{e}_r \quad (10.74)$$

¹⁰https://en.wikipedia.org/wiki/Spherical_coordinate_system

The acceleration is the so-called centripetal¹¹ acceleration. We coin $\dot{\phi} = \omega$ is the constant angular velocity.

“ A centripetal force (from Latin centrum, ”center” and petere, ”to seek”) is a force that makes a body follow a curved path. The direction of the centripetal force is always orthogonal to the motion of the body and towards the fixed point of the instantaneous center of curvature of the path.

One common example involving centripetal force is the case in which a body moves with uniform speed along a circular path. The centripetal force is directed at right angles to the motion and also along the radius towards the centre of the circular path. The mathematical description was derived in 1659 by the Dutch physicist Christiaan Huygens.

In the case of an object that is swinging around on the end of a rope in a horizontal plane, the centripetal force on the object is supplied by the tension of the rope. The rope example is an example involving a ’pull’ force.

Newton’s idea of a centripetal force corresponds to what is nowadays referred to as a central force. When a satellite is in orbit around a planet, gravity is considered to be a centripetal force even though in the case of eccentric orbits, the gravitational force is directed towards the focus, and not towards the instantaneous center of curvature.”¹²

“ A geostationary orbit, also referred to as a geosynchronous equatorial orbit, is a circular geosynchronous orbit 35,786 km in altitude above Earth’s equator, 42,164 km in radius from Earth’s center, and following the direction of Earth’s rotation. An object in such an orbit has an orbital period equal to Earth’s rotational period, one sidereal day, and so to ground observers it appears motionless, in a fixed position in the sky. ”¹³

The centripetal force of an orbiting body of mass m is then

$$\vec{F}_c = m\vec{a} = -mR\dot{\phi}^2 \vec{e}_r$$

The gravitational force is

$$\vec{F}_g = -\frac{\mathcal{G}Mm}{R^2}\vec{e}_r$$

From Newton’s second law of motion (sum of forces = mass * acceleration) we can write

$$m\vec{a} = -mR\dot{\phi}^2 \vec{e}_r = -\frac{\mathcal{G}Mm}{R^2}\vec{e}_r$$

or,

$$R\dot{\phi}^2 = \frac{\mathcal{G}M}{R^2}$$

We have $\dot{\phi} = \omega = \frac{v}{R}$ so

$$R\left(\frac{v}{R}\right)^2 = \frac{\mathcal{G}M}{R^2}$$

$$v^2 = \frac{\mathcal{G}M}{R}$$

The velocity is given by $2\pi R/T$ where T is the desired period so now

$$\left(\frac{2\pi R}{T}\right)^2 = \frac{\mathcal{G}M}{R}$$

where T is the orbital period (i.e. one sidereal day), and is equal to 86164.09054 s. In the end

$$R = \left(\frac{\mathcal{G}MT^2}{4\pi^2}\right)^{1/3}$$

¹¹Moving or tending to move towards a centre, as opposed to centrifugal: moving or tending to move away from a centre.

¹²https://en.wikipedia.org/wiki/Centripetal_force

¹³https://en.wikipedia.org/wiki/Geostationary_orbit

The resulting orbital radius is 42,164 kilometres. Subtracting the Earth's equatorial radius, 6,378 kilometres, gives the altitude of 35,786 kilometres.

The orbital speed is calculated by multiplying the angular speed by the orbital radius:

$$v = \omega R \simeq 3074.6m/s$$

10.7 Programming exercises - February 2024

gravity_exercises.tex

Background

We have seen that the calculation of the gravity vector and/or the gravity potential for a mass distribution in 3D space is of the form

$$\xi(\vec{r}) = \mathcal{G} \int_V f(\vec{r}, \vec{r}') \rho(\vec{r}') d\vec{r}'$$

where ξ is either g_x , g_y , g_z or U and f is a function of the coordinates \vec{r} and \vec{r}' .

Let us now assume that the body under consideration can be subdivided into N_e smaller blocks/elements. By virtue of the linearity of the integral, we have

$$\xi(\vec{r}) = \mathcal{G} \sum_{e=1}^{N_e} \int_{V_e} f(\vec{r}, \vec{r}') \rho(\vec{r}') d\vec{r}'$$

We can further assume that inside each element the density is constant so that

$$\xi(\vec{r}) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e \int_{V_e} f(\vec{r}, \vec{r}') d\vec{r}'$$

We will now make a strong assumption which is only valid when elements are (very) small: we will assume that we can replace $f(\vec{r}, \vec{r}')$ by $f(\vec{r}, \vec{r}_e)$ where \vec{r}_e is the location of the 'center' of the element. We then get:

$$\xi(\vec{r}) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e f(\vec{r}, \vec{r}_e) \int_{V_e} d\vec{r}'$$

And finally the integral term is simply the volume of the element V_e :

$$\xi(\vec{r}) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e f(\vec{r}, \vec{r}_e) V_e$$

In the end, assuming that the body of interest can be split into many small elements of constant density, the gravity fields at a location $\vec{r} = (x, y, z)$ can be computed as follows:

$$g_x(x, y, z) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e V_e \frac{x - x_e}{|\vec{r} - \vec{r}_e|^3} \quad (10.75)$$

$$g_y(x, y, z) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e V_e \frac{y - y_e}{|\vec{r} - \vec{r}_e|^3} \quad (10.76)$$

$$g_z(x, y, z) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e V_e \frac{z - z_e}{|\vec{r} - \vec{r}_e|^3} \quad (10.77)$$

$$U(x, y, z) = -\mathcal{G} \sum_{e=1}^{N_e} \rho_e V_e \frac{1}{|\vec{r} - \vec{r}_e|} \quad (10.78)$$

where

$$|\vec{r} - \vec{r}_e| = \sqrt{(x - x_e)^2 + (y - y_e)^2 + (z - z_e)^2}$$

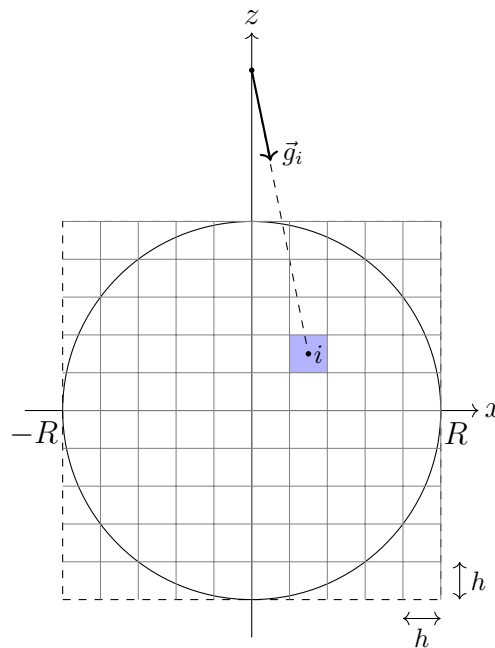
The following exercises are designed to test this approach which lends itself to numerical implementation. The basic idea is rather simple: generate a cloud of points in a regular manner such that we can assign them a corresponding volume and a density (and therefore a mass) when they are in the geometry of interest, and then use the formula above to compute the gravity vector and potential, and finally compare these values with the analytical solutions we derived for simple spherical bodies.

All quantities in the code(s) must be expressed in S.I. units, i.e. m, s, kg.

NO JUPYTER NOTEBOOK

Exercise 1: Full sphere

- (1A) We consider a domain of size $2R \times 2R \times 2R$ centered on the origin. It is partitioned in $N \times N \times N$ cells as shown in the following figure.



2D representation of the exercise. $N = 10$

Compute the total number of points NP as a function of N , the associated volume dV of a point (i.e. the volume of the cell the point is in) as a function of R and N and the size of a cell h as a function of R and N . Here is how your code should look like:

```
N=10
NP=
h=
dV=
```

- (1B-1) To get started, start in 1D. Assume that we consider a 1D cube, i.e. the segment $[-R, R]$ that is divided into N cells. What is h ? First, using the `linspace` function compute the x -coordinates of the N cell centers. Second, compute the same array \mathbf{x} using a single for loop. This second approach will be the building block of the next question.
- (1B-2) In the middle of each cell we place a point. Compute and store the coordinates of the N^3 points (use $R = 6371$ km). Please use arrays \mathbf{x} , \mathbf{y} and \mathbf{z} to store the coordinates. how your code should look like:

```

x = np.zeros(NP, dtype=np.float64) # x coordinates of all points
y = np.zeros(NP, dtype=np.float64) # y coordinates of all points
z = np.zeros(NP, dtype=np.float64) # z coordinates of all points
for i in range(?):
    for j in range(?):
        for k in range(?):
            x[?] = ?
            y[?] = ?
            z[?] = ?

```

Important: these arrays are N^3 long because they contain the coordinates of *all* points (cell centers).

Help: draw on paper a $3 \times 3 \times 3$ elements grid. Place the axis system on the plot and explicitly write the arrays for all 27 points. Example:

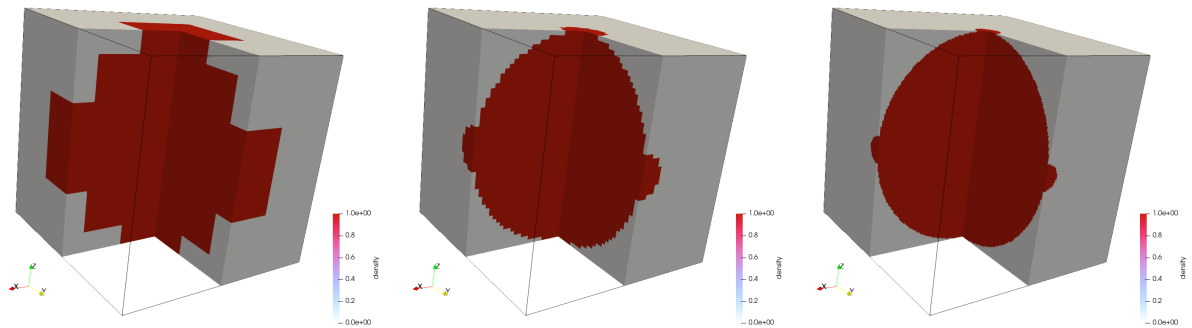
```

x[0]=..... y[0]=..... z[0]=.....
x[1]=..... y[1]=..... z[1]=.....
x[2]=..... y[2]=..... z[2]=.....
...
x[26]=.... y[26]=.... z[26]=.....

```

Once you have done so, run your code for $N = 3$, print the arrays and compare their content with what you have on paper (tip: for this test temporarily set $R = 3$).

- (1C) Assign a density $\rho_0 = 3000 \text{ kg m}^{-3}$ to points (cells) inside a sphere of radius R and zero otherwise, store these values in the `rho` array.



Example of density field ($\rho = 1$) for a 10^3 , 50^3 and 150^3 mesh.

- (1D) Fix $N = 10$. Compute the total mass of the sphere

$$M_s = \int_V \rho \, dV = \sum_e \rho_e V_e$$

and its volume¹⁴

$$V_s = \int_V dV = \sum_e V_e.$$

- (1E) Fix $N = 10$. Compute the moment of inertia of the sphere with respect to rotation axis z using Eq. (10.1).

¹⁴Note that the sums run over the cells which center lies inside the sphere.

- (1F) Repeat the last measurements (1D & 1E) with different values of $N \in (20, 30, 40, 50, \dots?)$. For both the mass and moment of inertia report on the relative error as a function of h .
- (1G) Compute the coordinates of 6 points situated at $z = 10^m$ meters *above the north pole* with $m = 0, 1, 2, 3, 4, 5$ and store these coordinates in arrays **xm**, **ym** **zm**.
- ```
xm = np.zeros(6, dtype=np.float64) # x coordinates of all points
ym = np.zeros(6, dtype=np.float64) # y coordinates of all points
zm = np.zeros(6, dtype=np.float64) # z coordinates of all points
```
- (1H) Fix  $N = 10$  for now. Compute the gravity potential  $U$  and acceleration vector components  $g_x, g_y, g_z$  at each of these 5 points using Eq. (10.14) (actually its discretised version, i.e. Eqs. (10.75), (10.76), (10.77) and (10.78).
- (1I) Plot the computed quantities as a function of  $z$  and plot on the same graphic the analytical values.
- (1J) Fix  $m = 4$ . Progressively increase  $N$  and record the absolute error on the gravity vector norm as a function of  $h$ . Plot this in log-log scale. Discuss.
- (1K) Use the `prem_density` function to assign the PREM [357] density to the points. Compute the mass of the planet with this new density distribution and compare it with the mass of the Earth. Compute the gravity at the surface. hint: use a large( $r$ )  $N$  for good results. How long are you willing to wait?
- (1L) Bonus: time how long it takes to compute  $U, g_x, g_y, g_z$  at a single location for  $N = 20$ . Report these times. Can you think of a way (and implement it) to arrive at the same results in less time?

## Exercise 2: Hollow sphere

This is based on the previous exercise.

- For points with radius  $r$  such that  $R/2 \leq r \leq R$  assign a density  $\rho_0 = 3000 \text{ kg m}^{-3}$  and zero otherwise.
- Compute the gravity potential and vector components on the  $x$ -axis between  $r = 0$  and  $r = 3R$  with steps of  $R/100$ .
- Plot the results and the analytical solution on the same plot as a function of  $r$ .
- Repeat the exercise with different values of  $N$ . Discuss.

## Exercise 3: Full sphere - revisited

...NOT for 2024...

We are now going to re-do the first exercise but this time we do not want any point outside of the sphere. We shall therefore use the spherical coordinates (see Section 2.3.4). We will use three for loops, one over  $r \in [0, R]$  values, one over  $\theta \in [0, \pi]$  values and one over  $\phi \in [-\pi, \pi]$  values. The number of points in each direction in this space is still  $N$  so that the total number of points is still  $N^3$ .

- Compute and store the coordinates of the points in the  $r, \theta, \phi$  space. Store these in arrays **r**, **theta**, **phi**.



- Use these coordinates to compute and store the Cartesian coordinates of these points.
- Plot this cloud of points in 3D. Discuss.
- Repeat the calculations of the first exercise.
- The cost of the calculation is the same as in exercise 1, but what about accuracy?

## Report

The report should contain results from exercises 1 and 2. I expect one pdf file per student (maximum 10 pages) and the corresponding python file(s), all delivered in a single zip file. Your report does not have to follow questions 1A to 1K in sequential order. Please send all files in a single email per **March 29th, 2024, 23:59**.

You should have the following guidelines in mind when writing your report:

- Layout: is the document visually pleasing? Is it well structured? are the student names and numbers present ?
- Is there a complete bibliography (if/when applicable)?
- Introduction: is the context clear? Are the methods presented?
- Figures: Are they properly numbered? captioned? all figures must be referenced in the text. Are they of good quality? are they readable? are all axis labelled?
- Text: Overall quality of the language. Are there still typos ? Do all sentences make sense?
- If results are wrong, was an attempt made to document/explain the (probable) source of the problem?
- Discussion: are the results properly discussed, analyzed? are potential problems, errors, limitations discussed?
- Conclusion: Are the report's findings summarized and when applicable generalized?

For concrete examples of what not to do, check Appendix [J](#).

## 10.8 Exam - February 2020

This is not the full 2020 exam, only the first two exercises.

### Exercise 1

Let us consider a spherically symmetric body of radius  $R$ .

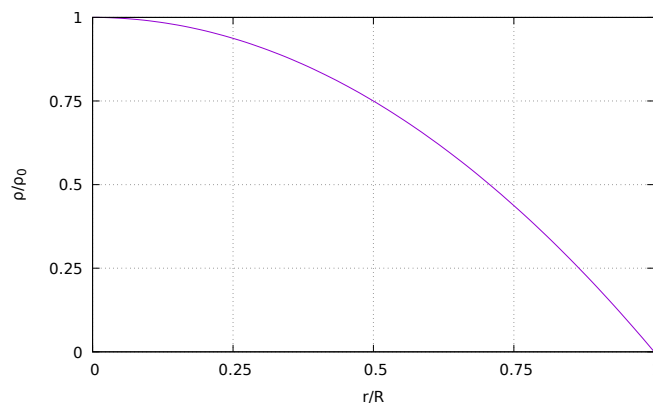
1. (1 pt) use symmetry considerations with regards to the 3 axis to arrive at the moment of inertia  $I$ :

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr \quad (10.79)$$

2. (1/2 pt) The density inside the sphere is given by

$$\rho(r) = \rho_0 \left[ a \left( \frac{r}{R} \right)^2 + b \left( \frac{r}{R} \right) + c \right]$$

Compute coefficients  $a, b, c$  such that the density matches the curve on the following figure:



3. (1 pt) Compute the total mass  $M$  of the planet using this expression for the density.
4. (1 pt) Compute the moment of inertia  $I$
5. (1/2 pt) Can  $I$  be written  $I = fMR^2$ ? if so, give  $f$ .
6. (1/2 pt) What are the dimensions of  $\rho$ ,  $I$  and  $M$ ?

### Exercise 1 - answer

Question 1 was treated in class. Looking at the figure, we see that  $\rho(r=0) = \rho_0$ , i.e.  $c = 1$ . Also, we see that  $\rho(r=R) = 0$ , i.e.

$$\rho_0 [a + b + 1] = 0$$

or,  $a + b = -1$ . Finally, we see that  $\rho(r=R/2) = 3\rho_0/4$ , so

$$\rho_0 \left[ \frac{a}{4} + \frac{b}{2} + 1 \right] = \frac{3}{4} \rho_0$$

or,  $a/4 + b/2 = -1/4$ . This leads to  $a = -1$  and  $b = 0$  so

$$\rho(r) = \rho_0 \left[ 1 - \left( \frac{r}{R} \right)^2 \right]$$

The total mass is given by

$$\begin{aligned}
M &= \int_V \rho(r) dV \\
&= 4\pi \int_0^R \rho(r) r^2 dr \\
&= 4\pi \int_0^R \rho_0 \left[ 1 - \left( \frac{r}{R} \right)^2 \right] r^2 dr \\
&= 4\pi \rho_0 \int_0^R \left[ r^2 - \frac{r^4}{R^2} \right] dr \\
&= \frac{8\pi}{15} \rho_0 R^3
\end{aligned} \tag{10.80}$$

The moment of inertia  $I$  is

$$\begin{aligned}
I &= \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr \\
&= \frac{16\pi}{105} \rho_0 R^5
\end{aligned} \tag{10.81}$$

Finally,

$$I = \frac{16\pi}{105} \rho_0 R^5 = \frac{2}{7} \left( \frac{8\pi}{15} \rho_0 R^3 \right) R^2$$

so  $f = 2/7$ .

## Exercise 2

Let us consider the same sphere as in the previous exercise, and the same density profile  $\rho(r)$ . The gravitational potential satisfies the Poisson equation:

$$\Delta U = 4\pi \mathcal{G} \rho(\vec{r}) \tag{10.82}$$

and we have the following relationship between the gravitational acceleration vector and the potential:  $\vec{g} = -\vec{\nabla} U$ .

- (1/2 pt) Write explicitly Eq.(10.82) for a point inside the sphere and a point outside the sphere.
- (1 pt) Compute  $g(r)$  and  $U(r)$  for a point inside the sphere as a function of  $r$ . Use  $\lim_{r \rightarrow 0} g(r) \neq \infty$  to get rid of an integration constant.
- (1 pt) Compute  $g(r)$  and  $U(r)$  for a point outside the sphere as a function of  $r$ . Use  $\lim_{r \rightarrow \infty} U(r) = 0$  to get rid of another integration constant.
- (1 pt) Use the continuity of  $g(r)$  and  $U(r)$  at  $r = R$  to compute the last two remaining integration constants.

## Exercise 2 - answer

Inside the sphere Eq.(10.82) is

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi \mathcal{G} \rho_0 \left[ a \left( \frac{r}{R} \right)^2 + b \left( \frac{r}{R} \right) + c \right]$$

We could use the values of  $a$ ,  $b$  and  $c$  obtained above but I here choose not to, in order to show that this exercise can be carried out independently from the first one. Then:

$$\begin{aligned}\frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) &= 4\pi \mathcal{G} \rho_0 \left[ a \frac{r^4}{R^2} + b \frac{r^3}{R} + cr^2 \right] \\ r^2 \frac{\partial U}{\partial r} &= 4\pi \mathcal{G} \rho_0 \left[ a \frac{r^5}{5R^2} + b \frac{r^4}{4R} + c \frac{r^3}{3} \right] + A \\ \frac{\partial U}{\partial r} &= 4\pi \mathcal{G} \rho_0 \left[ a \frac{r^3}{5R^2} + b \frac{r^2}{4R} + c \frac{r}{3} \right] + \frac{A}{r^2}\end{aligned}$$

so that

$$g(r) = -\frac{\partial U}{\partial r} = -4\pi \mathcal{G} \rho_0 \left[ a \frac{r^3}{5R^2} + b \frac{r^2}{4R} + c \frac{r}{3} \right] + \frac{A}{r^2}$$

We use  $\lim_{r \rightarrow 0} g(r) \neq \infty$  to arrive at  $A = 0$ . Finally

$$U_{in}(r) = 4\pi \mathcal{G} \rho_0 \left[ a \frac{r^4}{20R^2} + b \frac{r^3}{12R} + c \frac{r^2}{6} \right] + B$$

With  $a = -1$ ,  $b = 0$  and  $c = 1$ :

$$\begin{aligned}g_{in}(r) &= -4\pi \mathcal{G} \rho_0 \left[ -\frac{r^3}{5R^2} + \frac{r}{3} \right] \\ U_{in}(r) &= 4\pi \mathcal{G} \rho_0 \left[ -\frac{r^4}{20R^2} + \frac{r^2}{6} \right] + B\end{aligned}$$

Outside the sphere we must solve

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 0$$

$$\begin{aligned}r^2 \frac{\partial U}{\partial r} &= C \\ \frac{\partial U}{\partial r} &= \frac{C}{r^2} \\ U_{out}(r) &= -\frac{C}{r} + D\end{aligned}$$

We use  $\lim_{r \rightarrow \infty} U(r) = 0$  to arrive at  $D = 0$  so that

$$\begin{aligned}g_{out}(r) &= -\frac{C}{r^2} \\ U_{out}(r) &= -\frac{C}{r}\end{aligned}$$

Both fields should match at  $r = R$ :

$$\begin{aligned}g_{in}(r = R) &= g_{out}(r = R) \\ U_{in}(r = R) &= U_{out}(r = R)\end{aligned}$$

i.e.

$$-4\pi \mathcal{G} \rho_0 \left[ a \frac{R^3}{5R^2} + b \frac{R^2}{4R} + c \frac{R}{3} \right] = -\frac{C}{R^2}$$

so

$$C = 4\pi\mathcal{G}\rho_0 R^3 \left[ \frac{a}{5} + \frac{b}{4} + \frac{c}{3} \right]$$

Now,

$$4\pi\mathcal{G}\rho_0 \left[ a\frac{R^4}{20R^2} + b\frac{R^3}{12R} + c\frac{R^2}{6} \right] + B = -\frac{C}{R}$$

$$4\pi\mathcal{G}\rho_0 R^3 \left[ \frac{a}{20} + \frac{b}{12} + \frac{c}{6} \right] + BR = -4\pi\mathcal{G}\rho_0 R^3 \left[ \frac{a}{5} + \frac{b}{4} + \frac{c}{3} \right]$$

so

$$BR = -4\pi\mathcal{G}\rho_0 R^3 \left[ \frac{a}{5} + \frac{a}{20} + \frac{b}{4} + \frac{b}{12} + \frac{c}{3} + \frac{c}{6} \right]$$

$$B = -4\pi\mathcal{G}\rho_0 R^2 \left[ \frac{a}{4} + \frac{b}{3} + \frac{c}{2} \right]$$

With  $a = -1$ ,  $b = 0$  and  $c = 1$ :

$$B = -\pi\mathcal{G}\rho_0 R^2 \quad C = \frac{8\pi}{15}\mathcal{G}\rho_0 R^3$$

Finally

$$g_{in}(r) = -4\pi\mathcal{G}\rho_0 \left[ -\frac{r^3}{5R^2} + \frac{r}{3} \right]$$

$$U_{in}(r) = 4\pi\mathcal{G}\rho_0 \left[ -\frac{r^4}{20R^2} + \frac{r^2}{6} \right] - \pi\mathcal{G}\rho_0 R^2$$

$$g_{out}(r) = -\frac{8\pi}{15}\mathcal{G}\rho_0 R^3 \frac{1}{r^2}$$

$$U_{out}(r) = -\frac{8\pi}{15}\mathcal{G}\rho_0 R^3 \frac{1}{r}$$

Also, remembering that the mass of the planet is  $M = \frac{8\pi}{15}\rho_0 R^3$  so that

$$g_{out}(r) = -\frac{\mathcal{G}M}{r^2}$$

$$U_{out}(r) = -\frac{\mathcal{G}M}{r}$$

No surprise there...

## 10.9 Exam - March 2021

### Exercise 1

Let us consider a spherically symmetric body of radius  $R$ .

1. (1 pt) use symmetry considerations with regards to the 3 axis to arrive at the moment of inertia  $I$ :

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr \quad (10.83)$$

2. (1 pt) The density inside the sphere is given by

$$\rho(r) = a \frac{r}{R} + b$$

Compute the total mass  $M$  of the planet using this expression for the density.

3. (1 pt) Compute the moment of inertia  $I$  also as a function of  $a$  and  $b$ .
4. (1/2 pt) Can  $I$  be written  $I = fMR^2$ ? if so, give  $f$ .
5. (1/2 pt) What are the dimensions of  $a$ ,  $b$ ,  $\rho$ ,  $I$  and  $M$ ?
6. (1/2 pt) Set  $a = 0$  and  $b = \rho_0$  and look again at  $M$  and  $I$ . Conclude.

### Exercise 1 - answer

The first question is answered in Problem 1, see Section 10.3.

The mass of the planet is given by

$$\begin{aligned} M &= \iiint_V \rho(r) dV \\ &= \iiint_V \rho(r) r^2 \sin \theta dr d\theta d\phi \\ &= \iiint_V \left( a \frac{r}{R} + b \right) r^2 \sin \theta dr d\theta d\phi \\ &= 4\pi \int_0^R \left( a \frac{r}{R} + b \right) r^2 dr \\ &= 4\pi \left[ \frac{a}{R} \int_0^R r^3 dr + b \int_0^R r^2 dr \right] \\ &= 4\pi \left[ \frac{a}{R} \frac{1}{4} R^4 + b \frac{1}{3} R^3 \right] \\ &= 4\pi R^3 \left( \frac{a}{4} + \frac{b}{3} \right) \end{aligned} \quad (10.84)$$

The moment of inertia of the planet is given by

$$\begin{aligned}
 I &= \frac{8\pi}{3} \iiint_0^R \rho(r) r^4 dV \\
 &= \frac{8\pi}{3} \iiint_0^R \left( a \frac{r}{R} + b \right) r^4 dV \\
 &= \frac{8\pi}{3} \left( \frac{a}{6} + \frac{b}{5} \right) R^5 \\
 &= \underbrace{(4\pi R^3) \left( \frac{a}{4} + \frac{b}{3} \right)}_M \left( \frac{a}{4} + \frac{b}{3} \right)^{-1} R^2 \frac{2}{3} \left( \frac{a}{6} + \frac{b}{5} \right) \\
 &= \underbrace{\left( \frac{a}{4} + \frac{b}{3} \right)^{-1} \frac{2}{3} \left( \frac{a}{6} + \frac{b}{5} \right)}_f M R^2
 \end{aligned} \tag{10.85}$$

The dimensions of  $a$ ,  $b$ ,  $I$  and  $M$  are as follows:

$$[a] = [b] = ML^{-3} \quad [I] = ML^2 \quad [M] = M$$

When  $a = 0$  and  $b = \rho_0$  we find the standard mass of a constant density sphere  $M = \frac{4}{3}\pi R^3 \rho_0$  and  $f = 2/5$ .

## Exercise 2

Let us consider the same sphere as in the previous exercise, and the same density profile  $\rho(r)$ . The gravitational potential satisfies the Poisson equation:

$$\Delta U = 4\pi \mathcal{G} \rho(\vec{r}) \tag{10.86}$$

and we have the following relationship between the gravitational acceleration vector and the potential:  $\vec{g} = -\vec{\nabla} U$ .

- (1/2 pt) Write explicitly Eq.(10.86) for a point inside the sphere and a point outside the sphere.
- (1 pt) Compute  $g(r)$  and  $U(r)$  for a point inside the sphere as a function of  $r$ . Use  $\lim_{r \rightarrow 0} g(r) \neq \infty$  to get rid of an integration constant.
- (1 pt) Compute  $g(r)$  and  $U(r)$  for a point outside the sphere as a function of  $r$ . Use  $\lim_{r \rightarrow \infty} U(r) = 0$  to get rid of another integration constant.
- (1 pt) Use the continuity of  $g(r)$  and  $U(r)$  at  $r = R$  to compute the last remaining integration constants.
- (1/2 pt) Set  $a = 0$  and  $b = \rho_0$  and sketch the obtained fields.

## Exercise 2 - answer

outside  $\Delta U = 0$ . Inside  $\Delta U = 4\pi \mathcal{G} (a \frac{r}{R} + b)$

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi \mathcal{G} \left( a \frac{r}{R} + b \right)$$

$$\begin{aligned}
\Rightarrow \quad & \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi\mathcal{G} \left( a \frac{r^3}{R} + br^2 \right) \\
\Rightarrow \quad & r^2 \frac{\partial U}{\partial r} = 4\pi\mathcal{G} \left( a \frac{r^4}{4R} + b \frac{r^3}{3} \right) + A \\
\Rightarrow \quad & \frac{\partial U}{\partial r} = 4\pi\mathcal{G} \left( a \frac{r^2}{4R} + b \frac{r}{3} \right) + \frac{A}{r^2}
\end{aligned}$$

so that

$$g(r) = -\frac{\partial U}{\partial r} = -4\pi\mathcal{G} \left[ a \frac{r^2}{4R} + b \frac{r}{3} \right] + \frac{A}{r^2}$$

We use  $\lim_{r \rightarrow 0} g(r) \neq \infty$  to arrive at  $A = 0$ . Finally

$$g_{in}(r) = -4\pi\mathcal{G} \left( a \frac{r^2}{4R} + b \frac{r}{3} \right) \quad (10.87)$$

and

$$U_{in} = -\int g(r)dr = 4\pi\mathcal{G} \left( a \frac{r^3}{12R} + b \frac{r^2}{6} \right) + B$$

Outside the sphere we must solve

$$\begin{aligned}
& \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 0 \\
\Rightarrow \quad & r^2 \frac{\partial U}{\partial r} = C \\
\Rightarrow \quad & \frac{\partial U}{\partial r} = \frac{C}{r^2} \\
\Rightarrow \quad & U_{out}(r) = -\frac{C}{r} + D
\end{aligned}$$

We use  $\lim_{r \rightarrow \infty} U(r) = 0$  to arrive at  $D = 0$  so that

$$\begin{aligned}
g_{out}(r) &= -\frac{C}{r^2} \\
U_{out}(r) &= -\frac{C}{r}
\end{aligned}$$

Both fields should match at  $r = R$ :

$$\begin{aligned}
g_{in}(r = R) &= g_{out}(r = R) \\
U_{in}(r = R) &= U_{out}(r = R)
\end{aligned}$$

i.e.

$$\begin{aligned}
& -4\pi\mathcal{G} \left( a \frac{R^2}{4R} + b \frac{R}{3} \right) = -\frac{C}{R^2} \\
\Rightarrow \quad & C = 4\pi\mathcal{G}R^3 \left( \frac{a}{4} + \frac{b}{3} \right) = M\mathcal{G}
\end{aligned} \quad (10.88)$$

so unsurprisingly (!):

$$\begin{aligned}
g_{out}(r) &= -\frac{M\mathcal{G}}{r^2} \\
U_{out}(r) &= -\frac{M\mathcal{G}}{r}
\end{aligned}$$



and

$$4\pi\mathcal{G}\left(a\frac{R^3}{12R}+b\frac{R^2}{6}\right)+B=-\frac{M\mathcal{G}}{R}$$

$$4\pi\mathcal{G}R^2\left(a\frac{1}{12}+b\frac{1}{6}\right)+B=-\frac{M\mathcal{G}}{R}$$

$$B=-\frac{M\mathcal{G}}{R}-4\pi\mathcal{G}R^2\left(a\frac{1}{12}+b\frac{1}{6}\right)$$

so

$$U_{in}=4\pi\mathcal{G}\left(a\frac{r^3}{12R}+b\frac{r^2}{6}\right)-\frac{M\mathcal{G}}{R}-4\pi\mathcal{G}R^2\left(a\frac{1}{12}+b\frac{1}{6}\right)$$

If  $a=0$  and  $b=\rho_0$ , then

$$g_{in}(r)=-4\pi\mathcal{G}\rho_0\frac{r}{3}$$

and

$$U_{in}=4\pi\mathcal{G}\left(\rho_0\frac{r^2}{6}\right)-\frac{M\mathcal{G}}{R}-4\pi\mathcal{G}R^2\left(\rho_0\frac{1}{6}\right)=\frac{2\pi}{3}\mathcal{G}\rho_0r^2-\frac{3}{2}\frac{M\mathcal{G}}{R}$$

which is the solution to problem 14.

## 10.10 WORK in PROGRESS. DUH.

What follows on this page is an unfinished attempt to link spherical harmonics with my 2018 paper.

We start from the Poisson equation for the gravity potential:

$$\Delta U = 4\pi\mathcal{G}\rho(\vec{r}) \quad (10.89)$$

As a consequence, inside a domain where  $\rho = 0$ , the equation becomes  $\Delta U = 0$ .

Let us assume that the spherical coordinates are appropriate for the problem at hand, and that the potential can be decomposed as follows:

$$U(r, \theta, \phi) = U_r(r)U_{\perp}(\theta, \phi)$$

The full Laplacian operator in spherical coordinates is given by<sup>15</sup>:

$$\Delta U = \underbrace{\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right)}_{\Delta_r} + \underbrace{\frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left( \sin \theta \frac{\partial U}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2 U}{\partial \phi^2}}_{\Delta_{\perp}}$$

we then have:

$$(\Delta_r + \Delta_{\perp})(U_r U_{\perp}) = 0$$

i.e.,

$$U_{\perp} \Delta_r U_r + U_r \Delta_{\perp} U_{\perp} = 0$$

Assuming  $U_{\perp} = \sum_l \sum_m U_{lm} Y_{lm}$ , knowing that spherical harmonics functions verify

$$r^2 \Delta_{\perp} Y_l^m(\theta, \phi) = -l(l+1) Y_l^m(\theta, \phi)$$

and assuming for now that the problem at hand is 1st degree ( $l=1$ ), then

$$\Delta_{\perp} Y_l^m(\theta, \phi) = -\frac{2}{r^2} Y_l^m(\theta, \phi)$$

and then

$$\Delta_r U_r - U_r \frac{2}{r^2} = 0$$

make a link with my 2018 paper.

---

<sup>15</sup>[https://en.wikipedia.org/wiki/Laplace\\_operator](https://en.wikipedia.org/wiki/Laplace_operator)

In spherical coordinates, the Laplacian is given by

$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial}{\partial r} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial}{\partial \theta} \left( \sin \theta \frac{\partial}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2}{\partial \phi^2}$$

We wish to solve Laplace's equation  $\Delta T(r, \theta, \phi) = 0$  using the method of separation of variables:

$$T(r, \theta, \phi) = R(r)\Theta(\theta)\Phi(\phi)$$

We can insert this decomposition into the Laplace equation and multiply it by  $r^2/R\Theta\Phi$  to obtain

$$\frac{1}{R} \frac{d}{dr} \left( r^2 \frac{dR}{dr} \right) + \frac{1}{\Theta \sin \theta} \frac{d}{d\theta} \left( \sin \theta \frac{d\Theta}{d\theta} \right) + \frac{1}{\Phi \sin^2 \theta} \frac{d^2 \Phi}{d\phi^2} = 0$$

For reasons that will become clear later, the separation constant is taken to be  $-m^2$ :

$$\frac{1}{\Phi} \frac{d^2 \Phi}{d\phi^2} = -m^2 \quad (10.90)$$

$$-\frac{\sin \theta}{\Theta} \frac{d}{d\theta} \left( \sin \theta \frac{d\Theta}{d\theta} \right) - \frac{\sin^2 \theta}{R} \frac{d}{dr} \left( r^2 \frac{dR}{dr} \right) = -m^2 \quad (10.91)$$

The first equation yields

$$\Phi(\phi) = \begin{cases} e^{im\phi} \\ e^{-im\phi} \end{cases} \quad \text{for } m = 0, 1, 2, 3, \dots$$

Note that  $m$  must be an integer since  $\phi$  is a periodic variable and  $\Phi(\phi + 2\pi) = \Phi(\phi)$ . In the case of  $m = 0$ , the general solution is  $\Phi(\phi) = a\phi + b$ , but we must choose  $a = 0$  to be consistent with  $\Phi(\phi + 2\pi) = \Phi(\phi)$ . Hence in the case of  $m = 0$ , only one solution is allowed.

Eq. (10.91) can now be recast in the following form:

$$\frac{1}{R} \frac{d}{dr} \left( r^2 \frac{dR}{dr} \right) = -\frac{1}{\Theta \sin \theta} \frac{d}{d\theta} \left( \sin \theta \frac{d\Theta}{d\theta} \right) + \frac{m^2}{\sin^2 \theta} \quad (10.92)$$

where the separation variable at this step is denoted by  $l(l+1)$  for reasons that will shortly become clear. The resulting radial equation is

$$\frac{1}{R} \frac{d}{dr} \left( r^2 \frac{dR}{dr} \right) = l(l+1)$$

or,

$$r^2 \frac{d^2 R}{dr^2} + 2r \frac{dR}{dr} - l(l+1)R = 0$$

The solution is of the form  $R = r^s$ . To determine the exponent  $s$ , we insert this solution back into the above ODE. The end result is

$$s(s+1) = l(l+1) \quad \Rightarrow \quad s = l \text{ or } s = -l-1$$

or,

$$R(r) = \begin{cases} r^l \\ r^{-(l+1)} \end{cases}$$

Eq. (10.92) also yields:

$$\frac{1}{\sin \theta} \frac{d}{d\theta} \left( \sin \theta \frac{d\Theta}{d\theta} \right) + \left[ l(l+1) - \frac{m^2}{\sin^2 \theta} \right] \Theta = 0$$

One can then carry out the following change of variables  $x = \cos \theta$  and  $y = \Theta(\theta)$  so that the above equation reduces to:

$$(1 - x^2) \frac{d^2 y}{dx^2} - 2x \frac{dy}{dx} + \left[ l(l+1) - \frac{m^2}{\sin^2 \theta} \right] y = 0$$

This equation is the differential equation for associated Legendre polynomials<sup>16</sup>. We then have

$$y = P_l^m(x) \quad \text{for } l = 0, 1, 2, 3, \dots \quad \text{and } m = -l, -l+1, \dots, 0, \dots, l-1, l$$

and

$$P_l^m(x) = \frac{(-1)^m}{2^l l!} (1 - x^2)^{m/2} \frac{d^{l+m}}{dx^{l+m}} (x^2 - 1)^l$$

with  $m \geq 0$  and  $l \geq 0$ . The first few polynomials are

$$\begin{aligned} P_0^0(\cos \theta) &= 1 \\ P_1^{-1}(\cos \theta) &= \frac{1}{2} \sin \theta \\ P_1^0(\cos \theta) &= \cos \theta \\ P_1^{+1}(\cos \theta) &= -\sin \theta \\ P_2^{-2}(\cos \theta) &= \frac{1}{8} \sin^2 \theta \\ P_2^{-1}(\cos \theta) &= \frac{1}{2} \sin \theta \cos \theta \\ P_2^0(\cos \theta) &= \frac{1}{2} (3 \cos^2 \theta - 1) \\ P_2^{+1}(\cos \theta) &= -3 \sin \theta \cos \theta \\ P_2^{+2}(\cos \theta) &= 3 \sin^2 \theta \end{aligned}$$

In our case the differential equation for the associated Legendre polynomials, given above, depends on  $m^2$  and is therefore not sensitive to the sign of  $m$ . Consequently,  $P_l^m(x)$  and  $P_l^{-m}(x)$  must be equivalent solutions and hence proportional to each other, and one can show that

$$P_l^{-m}(\cos \theta) = (-1)^m \frac{(l-m)!}{(l+m)!} P_l^m(\cos \theta) \quad (10.93)$$

Combining all the results obtained above, we have found that the general solution to Laplace's equation is of the form

$$T(r, \theta, \phi) = \left\{ \begin{matrix} r^l \\ r^{-(l+1)} \end{matrix} \right\} P_l^m(\cos \theta) \left\{ \begin{matrix} e^{im\phi} \\ e^{-im\phi} \end{matrix} \right\}$$

where  $l = 0, 1, 2, 3, \dots$  and  $m = -l, -l+1, \dots, l-1, l$ .

When solving the Laplace's equation in spherical coordinates, it is traditional to introduce the spherical harmonics,  $Y_l^m(\theta, \phi)$ :

$$Y_l^m(\theta, \phi) = (-1)^m \sqrt{\frac{2l+1}{4\pi} \frac{(l-m)!}{(l+m)!}} P_l^m(\cos \theta) e^{im\phi} \quad \text{for } l = 0, 1, 2, 3, \dots \quad \text{and } m = -l, -l+1, \dots, l-1, l \quad (10.94)$$

<sup>16</sup>[https://en.wikipedia.org/wiki/Associated\\_Legendre\\_polynomials](https://en.wikipedia.org/wiki/Associated_Legendre_polynomials)

The phase factor  $(-1)^m$ , introduced originally by Condon and Shortley, is convenient for applications in quantum mechanics. Note that Eq. (10.93) implies that

$$Y_l^{-m}(\theta, \phi) = (-1)^m Y_l^m(\theta, \phi)^*$$

where the star means complex conjugation.

The normalization factor in Eq. (10.94) has been chosen such that the spherical harmonics are normalized to one. In particular, these functions are orthonormal and complete. The orthonormality relation is given by:

$$\int Y_l^m(\theta, \phi) Y_{l'}^{m'}(\theta, \phi) d\Omega = \delta_{ll'} \delta_{mm'}$$

where  $d\Omega = \sin \theta d\theta d\phi$  is the differential solid angle in spherical coordinates.

It is important to note that there are different normalisations for spherical harmonics. In this document we choose:

$$Y_l^m(\theta, \phi) = \sqrt{\frac{2l+1}{4\pi} \frac{(l-m)!}{(l+m)!}} P_l^m(\cos \theta) e^{im\phi} \quad \text{for } l = 0, 1, 2, 3, \dots \text{ and } m = -l, -l+1, \dots, l-1, l \quad (10.95)$$

which is Eq.(7.8.1) in the Schubert, Turcotte & Olson book [1140]. In this case the  $(-1)^m$  is inside the  $P_l^m$ . The first few spherical harmonics are shown below in the real representation (i.e. using  $\cos m\phi$  instead of  $e^{im\phi}$ )<sup>17 18</sup>:

$$\begin{aligned} Y_0^0(\theta, \phi) &= \sqrt{\frac{1}{4\pi}} \\ Y_1^{-1}(\theta, \phi) &= \sqrt{\frac{3}{8\pi}} \cos \phi \sin \theta \\ Y_1^0(\theta, \phi) &= \sqrt{\frac{3}{4\pi}} \cos \theta \\ Y_1^{+1}(\theta, \phi) &= -\sqrt{\frac{3}{8\pi}} \cos \phi \sin \theta \\ Y_2^{-2}(\theta, \phi) &= \sqrt{\frac{15}{32\pi}} \cos(2\phi) \sin^2 \theta \\ Y_2^{-1}(\theta, \phi) &= \sqrt{\frac{15}{8\pi}} \cos \phi \sin \theta \cos \theta \\ Y_2^0(\theta, \phi) &= \sqrt{\frac{5}{16\pi}} (3 \cos^2 \theta - 1) \\ Y_2^{+1}(\theta, \phi) &= -\sqrt{\frac{15}{8\pi}} \cos \phi \sin \theta \cos \theta \\ Y_2^{+2}(\theta, \phi) &= \sqrt{\frac{15}{32\pi}} \cos(2\phi) \sin^2 \theta \end{aligned}$$

replace those my complex ones !

Another normalisation is sometimes used:

$$Y_l^m(\theta, \phi) = \sqrt{(2l+1) \frac{(l-m)!}{(l+m)!}} P_l^m(\cos \theta) e^{im\phi} \quad \text{for } l = 0, 1, 2, 3, \dots \text{ and } m = -l, -l+1, \dots, l-1, l \quad (10.96)$$

<sup>17</sup>[https://en.wikipedia.org/wiki/Table\\_of\\_spherical\\_harmonics](https://en.wikipedia.org/wiki/Table_of_spherical_harmonics)

<sup>18</sup><https://mathworld.wolfram.com/SphericalHarmonic.html>


with

$$\frac{1}{4\pi} \int Y_l^m(\theta, \phi) Y_{l'}^{m'}(\theta, \phi) d\Omega = \delta_{ll'} \delta_{mm'}$$

**Remark.** In [1412] the authors use a normalized associated Legendre polynomial that is related to the associated Legendre polynomial  $P_l^m$  as:

$$p_{lm}(\theta, \phi) = \sqrt{\frac{2l+1}{2\pi(1+\delta_{m0})} \frac{(l-m)!}{(l+m)!}} P_l^m(\cos \theta)$$

Note the absence of the  $(-1)^m$  term and the presence of the kronecker delta in the denominator.

 Relevant Literature: SHTools: Tools for Working with Spherical Harmonics [1354]

## 10.11 Gravity benchmarks

gravity\_benchmarks.tex

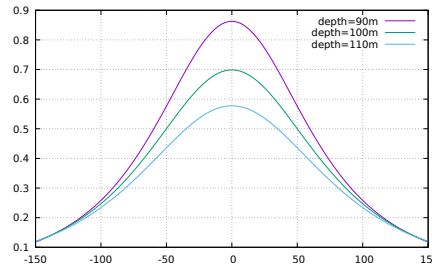
There are many analytical solutions for buried bodies of simple shape. Hereafter are the most common ones:

### 10.11.1 Buried sphere (3D)

To calculate the pull of gravity, we can use the fact that a sphere has the same gravitational pull as a point mass located at its centre. The distance between the measurement point and the center of the sphere is  $\sqrt{x^2 + d^2}$ , so

$$g_z = \frac{\mathcal{G} M_{\text{sphere}} d}{(x^2 + d^2)^{3/2}}$$

Let us take the following example: radius  $a=50\text{m}$ ,  $\Delta\rho = 2000$ , variable depth  $d=100\text{m}$



$g_z$  has its maximum value directly above the sphere at  $x = 0\text{m}$  and is given by

$$g_z^{\text{max}} = \frac{\mathcal{G} M_{\text{sphere}} d}{(d^2)^{3/2}} = \frac{\mathcal{G} M_{\text{sphere}}}{d^2}$$

We can then find the half width of the curve by finding  $x_{1/2}$  such that

$$\frac{\mathcal{G} M_{\text{sphere}} d}{(x_{1/2}^2 + d^2)^{3/2}} = \frac{g_z^{\text{max}}}{2} = \frac{\mathcal{G} M_{\text{sphere}}}{2d^2}$$

or, FINISH , derive  $x_{1/2}$

### 10.11.2 Buried horizontal cylinder (3D)

anticline can be approximated by a horizontal cylinder

$$g_z = \frac{2\mathcal{G}\pi a^2 d \Delta\rho}{x^2 + d^2}$$

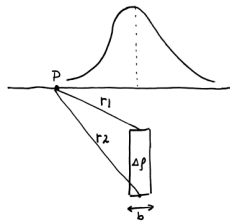
the maximum value of  $g_z$  is located directly above the axis of the cylinder

$g_{z\max}$  for a cylinder is larger than  $g_{z\max}$  for a sphere of the same radius.

Cannot distinguish a buried sphere from a cylinder with just a single profile. Need to collect gravity on a grid and make a map.

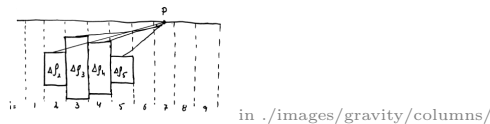
### 10.11.3 Buried column (2D)

$$g_z = 2\mathcal{G}\Delta\rho b \ln \frac{r_2}{r_1}$$



### 10.11.4 Buried columns (2D)

$$g_z = 2\mathcal{G} \sum_i \Delta\rho_i b_i \ln \frac{r_{2,i}}{r_{1,i}}$$



### 10.11.5 Uniform layer of rock

A layer of rock has an infinite extent, thickness  $\Delta z$  and a density  $\rho$ . The gravitational attraction of this slab at the point P at height  $z$  above the layer is

$$g_z = 2\pi\mathcal{G}\rho\Delta z$$

Note that  $g_z$  does not depend on the distance from the layer to the measurement point.

### 10.11.6 A constant density shell (Root *et al.* , 2021)

Results & raw data in `./images/benchmark_gravity/bench1`

The shell is defined between  $R_1$  and  $R_2$ . It contains a single material of density  $3300\text{kg m}^{-3}$ . The layer is centered around depth 100km. Gravity is measured 250km above the surface, i.e.  $r = 6621\text{km}$ . The thickness of the shell is 2, 5 or 10km.

The analytical gravity vector norm is given by

$$g = \frac{4\pi}{3}\rho\frac{R_2^3 - R_1^3}{r^2}\mathcal{G}$$

where we take  $\mathcal{G} = 6.67384 \cdot 10^{-11}$  (default in ASPECT ) or  $\tilde{\mathcal{G}} = 6.67428 \cdot 10^{-11}$  sometimes.

| shell thickness<br>(km) | volume<br>(m <sup>3</sup> ) | mass<br>(kg) | gravity using $\mathcal{G}$<br>(m s <sup>-2</sup> ) | gravity using $\tilde{\mathcal{G}}$<br>(m s <sup>-2</sup> ) |
|-------------------------|-----------------------------|--------------|-----------------------------------------------------|-------------------------------------------------------------|
| 2                       | 9.8835614e17                | 3.2615753e21 | 496.542034795                                       | 496.574771345                                               |
| 5                       | 2.4708905e17                | 8.1539385e21 | 1241.35514223                                       | 1241.43698361                                               |
| 10                      | 4.9417817e18                | 1.630788e22  | 2482.71067903                                       | 2482.87436182                                               |

In the ASPECT input file there are three main parameters which may influence the results:

- the radial resolution, controlled in the input file by: `set Number of slices = 1,2,3,4`
- the tangential/lateral resolution, controlled by: `set Initial lateral refinement = 3,4,5,6`
- the number of (additional) quadrature points, controlled by: `set Quadrature degree increase =0,1,...6`

We set here the default values at 1, 6 and 3 respectively.

|                                   | lat. res. 3    | lat. res. 4     | lat. res. 5      | lat. res. 6      | lat. res. 7       |
|-----------------------------------|----------------|-----------------|------------------|------------------|-------------------|
| nslice=1 (1 cells radial) # cells | 384            | 1,536           | 6,144            | 24,576           | 98,304            |
|                                   | $6 \times 64$  | $6 \times 256$  | $6 \times 1,024$ | $6 \times 4,096$ | $6 \times 16,384$ |
|                                   | $6 \times 8^2$ | $6 \times 16^2$ | $6 \times 32^2$  | $6 \times 64^2$  | $6 \times 128^2$  |
| nslice=2 (2 cells radial) # cells | 768            | 3,072           | 12,288           | 49,152           | 196,608           |
| nslice=3 (3 cells radial) # cells | 1,152          | 4,608           | 18,432           | 73,728           | 294,912           |
| nslice=4 (4 cells radial) # cells | 1,536          | 6,144           | 24,576           | 98,304           | 393,216           |
| average area (m2)                 | 1.328292e+12   | 3.320732e11     | 8.30183e10       | 2.075457e10      | 5.188644e9        |
| approx size (km)                  | 1152km         | 576             | 288km            | 144km            | 72km              |
| approx size (degree)              | 10.5           | 5.2             | 2.6              | 1.3              | 0.65              |

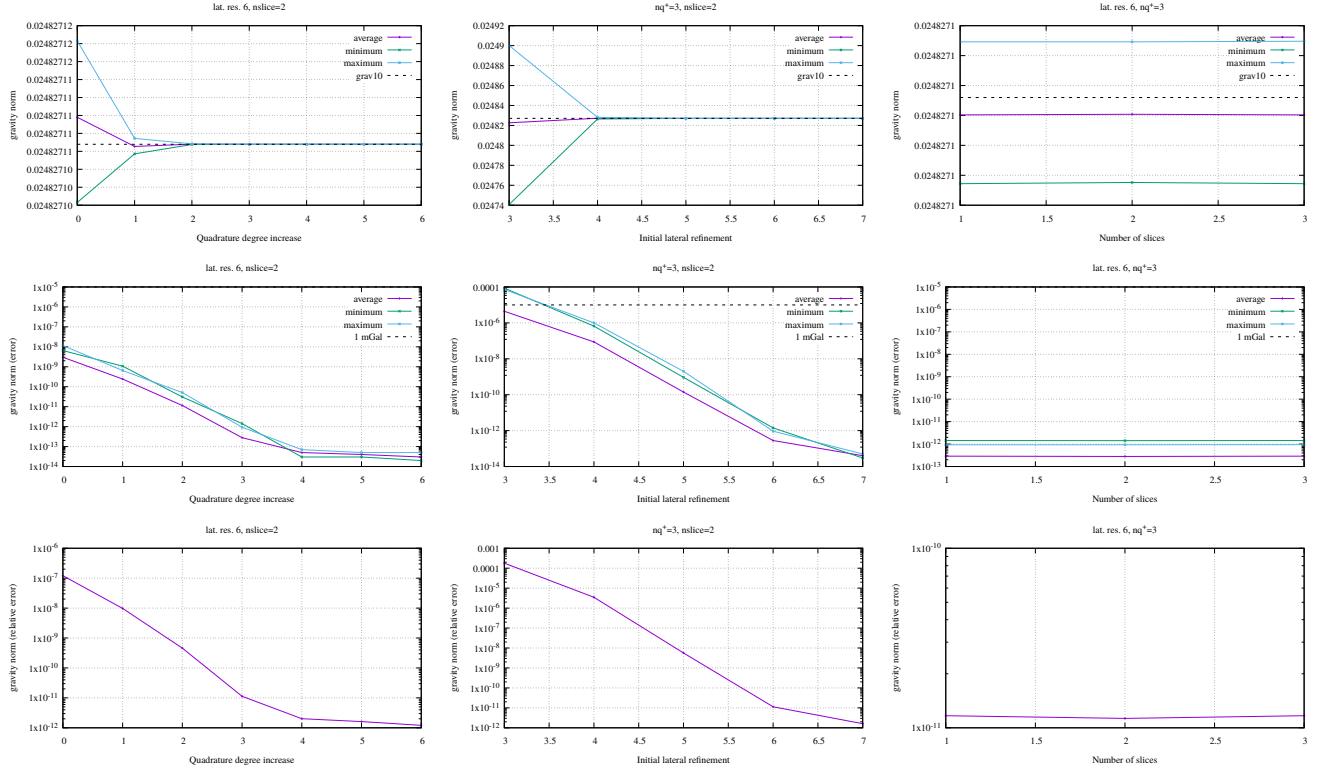
Earth has a surface of  $\mathcal{S} = 4\pi R^2 \simeq 5.1006447 \cdot 10^{14}\text{m}^2$ . An average degree resolution means that this surface would be tessellated in blocks of approximately  $2\pi R/360 \simeq 111\text{km}$  size. There would then be about 41,398 of such blocks. If a resolution of 2 degrees is required, then the blocks would be about 220km in size and there would be about 10,349 blocks.

Results obtained with ASPECT with  $\tilde{\mathcal{G}}$  are in the following table:

| Thickness (km)          | 2                                | 5                                   | 10                                  |
|-------------------------|----------------------------------|-------------------------------------|-------------------------------------|
| Shell formula (mGal)    | 496.574771345                    | 1241.43698361                       | 2482.87436182                       |
| $m = 4, \sim 5^\circ$   | 496.554320/496.602897/496.574854 | 1241.385829/1241.507337/1241.437190 | 2482.771870/2483.015344/2482.874775 |
| $m = 5, \sim 2.6^\circ$ | 496.574748/496.574819/496.574771 | 1241.436926/1241.437102/1241.436983 | 2482.874246/2482.874599/2482.874361 |
| $m = 6, \sim 1.3^\circ$ | 496.574771/496.574771/496.574771 | 1241.436984/1241.436984/1241.436984 | 2482.874362/2482.874362/2482.874362 |



## Results for a 10km thick shell with Aspect

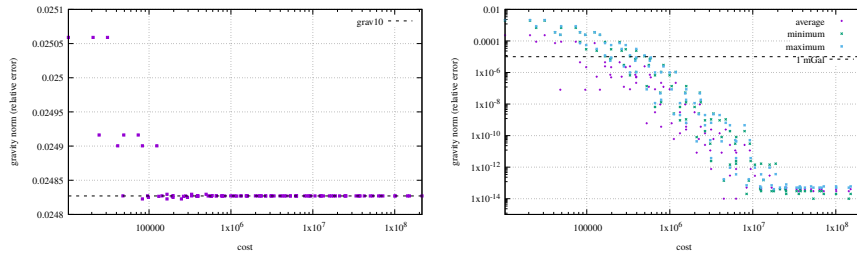


in ./images/benchmark\_gravity/bench1/

I then define the concept of 'cost'. In terms of computational cost, there is a tradeoff between resolution and number of quadrature points. The cost is then defined as

$$C = nel \times nq^3$$

where  $nel$  is the number of elements in the mesh and  $nq$  is the number of quadrature points per element and per dimension.



in ./images/benchmark\_gravity/bench1/

Preliminary conclusion:  $nq \geq 3$ ,  $nslice$  not so important here,  $lat\ res \geq 6$   
 TODO RUN 2 and 5 km shells !

### 10.11.7 The WINTERC mono-layer benchmark (Root *et al.* , 2021)

A single data file is used, `rho_56km_SH.W32.txt`, stored in `images/benchmark_gravity/bench2`. It contains density values for a single layer comprised between 56 and 80km depths, i.e. there is no radial variation of the density. Because of how ASPECT works, the density values need to be transformed into initial temperatures. Using the simple material model we have

$$\rho = \rho_0(1 - \alpha(T - T_0))$$

so

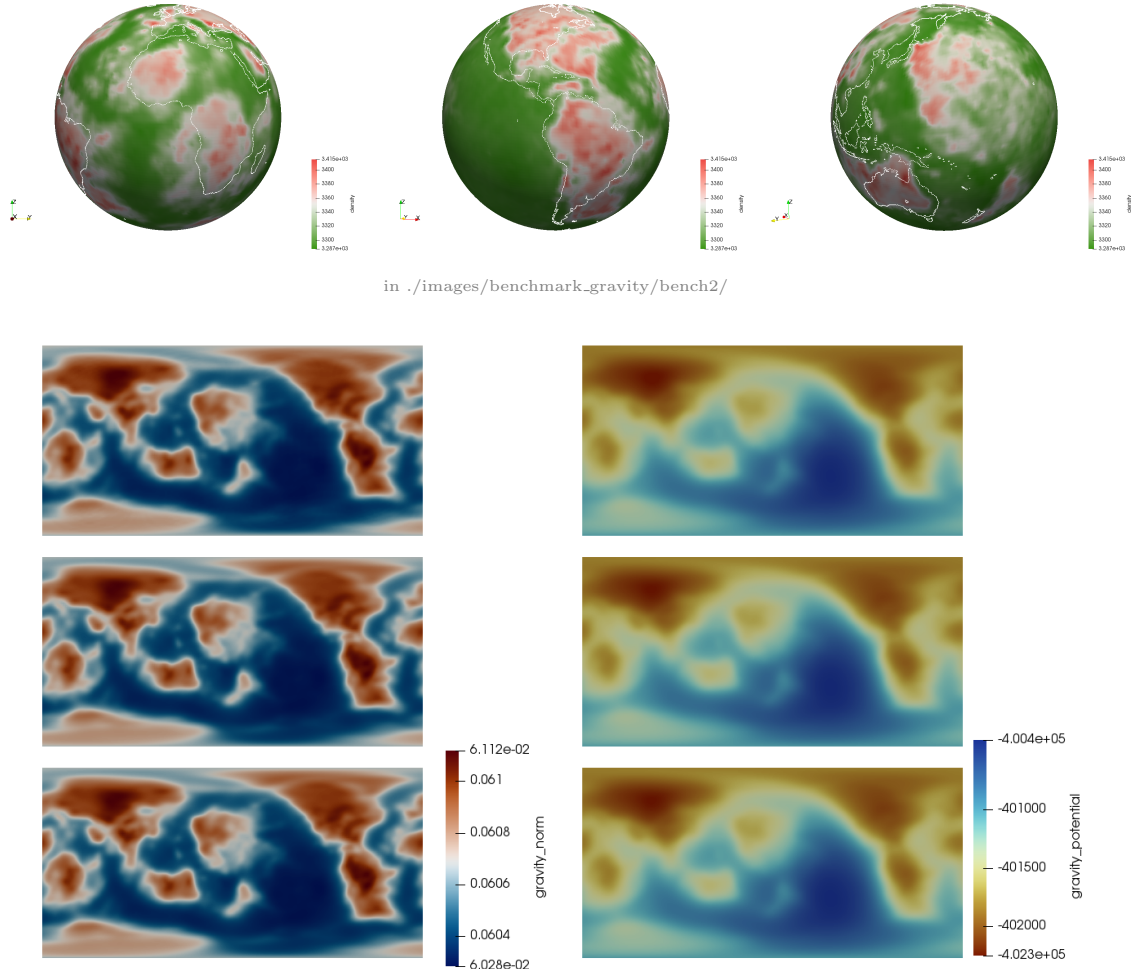
$$T = \frac{1}{\alpha} \left( 1 - \frac{\rho}{\rho_0} \right)$$

and we here take  $\alpha = 3 \cdot 10^{-5}$ ,  $T_0 = 0$  and  $\rho_0 = 3300$ , so that the temperatures range between -1198.9 and 141.5. These values make no sense, but all we want is that the densities  $\rho(T)$  generated by the material model are those of the original dataset.

Furthermore, the `rho_56km_SH.W32.txt` file contains 720x360 lines, i.e. the resolution is a half degree for longitude and latitude. These range from 0.25 to 359.75 and from -89.75 to 89.75 respectively. These must be transformed into spherical coordinates  $\phi \in [0, 2\pi]$  and  $\theta \in [0, \pi]$ .

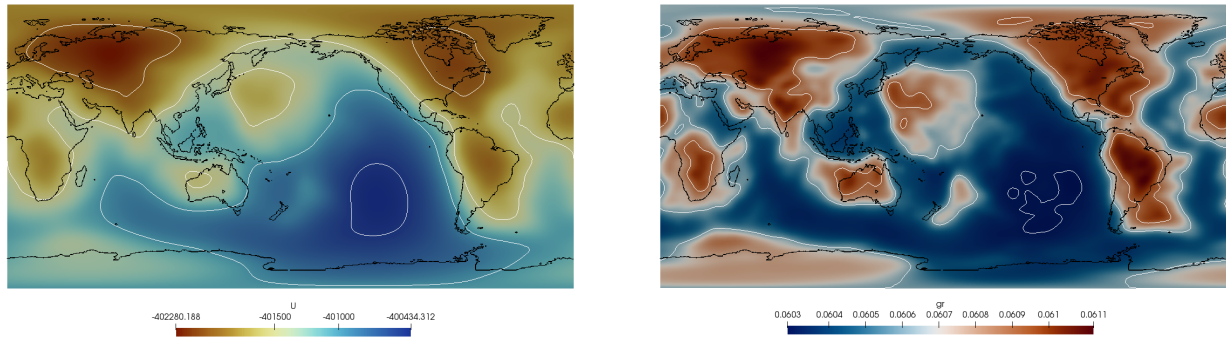
Also the original data file contains longitude, latitude and density values for a thick layer. The ascii data file which ASPECT can read requires radial values in increasing order as well, so for each combination  $\phi - \theta$  I generate two values, one at radius 6371-81km and one at 6371-55km so that the depth layer 56-80km fits in it. The data format of the ascii file is specified in `data/initial-temperature/ascii-data/test/shell_3d.txt` in ASPECT .

Stone 98 reads in `rho_56km_SH.W32.txt` and generates the `bench2.txt` file which is to be read by ASPECT . Note that the first line of this file is mandatory and reads: `# POINTS: 2 720 360`



Results obtained with ASPECT . Top to bottom, level 4, 5, 6. Measurements grid is 181x91 points.

in ./images/benchmark\_gravity/bench2/



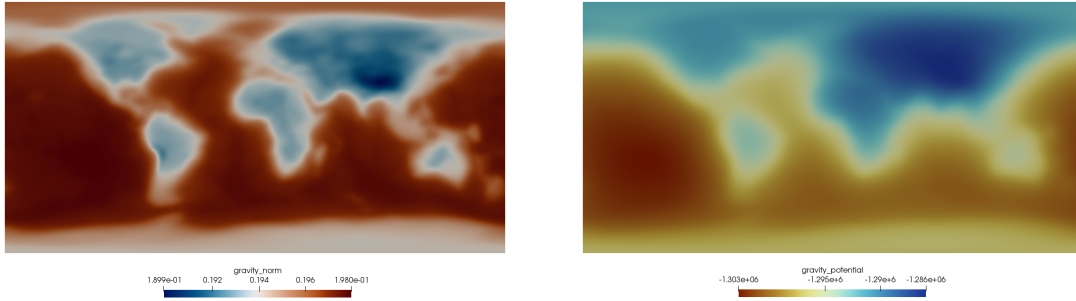
From Stone 98. Resolution of measurement grid is 181x91. It took about 19,100 seconds to run, averaging 1.16s per measurement point. Potential isocontours at -400.5e3, -401e3, -401.5e3, -402e3. Radial acceleration contours at 0.0603, 0.0606 and 0.0609.

in ./python\_codes/images/

| level | avrg density | total mass | number of cells | time  |
|-------|--------------|------------|-----------------|-------|
| 4     | 3323         | 3.981e+22  | 1,536           | 503s  |
| 5     | 3323         | 3.981e+22  | 6,144           | 2030s |
| 6     | 3323         | 3.981e+22  | 24.576          | 8190s |

### 10.11.8 Moho benchmark (Root *et al.* , 2021)

We consider an 80km thick shell with a density interface inside, using CRUST1.0 Moho for the boundary (upper dens = 2900kg/m<sup>3</sup>, lower dens = 3300 kg/m<sup>3</sup>). Stone 97 reads the CRUST1.0 file and transforms it into `bench3.ascii` in the right ASPECT format.



bench3.ascii file has 301 points in the radial direction, i.e. a 300m resolution. 181x91 gravity measurement points. Lateral refinement level is 6, 25 slices.

in `./images/benchmark_gravity/bench3/`

### 10.11.9 Gravity potential and gravity field of a two-layer spherically symmetric planet

Note: see also problem 16 in previous section.

Let us assume that the planet consists of two layers: the core (of density  $\rho_c$ ) and the mantle (of density  $\rho_m$ ). Its outer radius is  $R_2$  and the cmb is at  $R_1$ . We wish to compute the gravitational potential of this system and we start from the spherically symmetric Poisson equation:

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi \mathcal{G} \rho \quad (10.97)$$

We have three domains on the  $r$ -axis ( $a$  stands for air):

- $r \in [0, R_1]$ , density  $\rho_c$
- $r \in [R_1, R_2]$ , density  $\rho_m$
- $r \in [R_2, +\infty)$ , density  $\rho_a = 0$

We denote  $\hat{\rho} = 4\pi \mathcal{G} \rho$ . In the end we have to solve Eq. 10.97 in all three domains:

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = \hat{\rho}_c \quad r \in [0, R_1] \quad (10.98)$$

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = \hat{\rho}_m \quad r \in [R_1, R_2] \quad (10.99)$$

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 0 \quad r \in [R_2, +\infty) \quad (10.100)$$

The generic solution of Eq. (10.97) is obtained as follows:

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = \hat{\rho} \quad (10.101)$$

$$\Rightarrow \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = \hat{\rho} r^2 \quad (10.102)$$

$$\Rightarrow r^2 \frac{\partial U}{\partial r} = \frac{1}{3} \hat{\rho} r^3 + A \quad (10.103)$$

$$\Rightarrow \frac{\partial U}{\partial r} = \frac{1}{3} \hat{\rho} r + \frac{A}{r^2} = g(r) \quad (10.104)$$

$$\Rightarrow U(r) = \frac{1}{6} \hat{\rho} r^2 - \frac{A}{r} + B \quad (10.105)$$

So now we have the following set of equations:

$$g(r) = \frac{1}{3} \hat{\rho}_c r + \frac{A}{r^2} \quad U(r) = \frac{1}{6} \hat{\rho}_c r^2 - \frac{A}{r} + B \quad r \in [0, R_1] \quad (10.106)$$

$$g(r) = \frac{1}{3} \hat{\rho}_m r + \frac{C}{r^2} \quad U(r) = \frac{1}{6} \hat{\rho}_m r^2 - \frac{C}{r} + D \quad r \in [R_1, R_2] \quad (10.107)$$

$$g(r) = \frac{1}{3} \hat{\rho}_a r + \frac{E}{r^2} \quad U(r) = \frac{1}{6} \hat{\rho}_a r^2 - \frac{E}{r} + F \quad r \in [R_2, +\infty) \quad (10.108)$$

We know that  $\rho_a = 0$  and we additionally impose

$$\lim_{r \rightarrow 0} g(r) = 0 \quad \lim_{r \rightarrow +\infty} U(r) = 0$$

which automatically leads to  $A = F = 0$ .

In the end we have:

$$g(r) = \frac{1}{3}\hat{\rho}_c r \quad U(r) = \frac{1}{6}\hat{\rho}_c r^2 + B \quad r \in [0, R_1] \quad (10.109)$$

$$g(r) = \frac{1}{3}\hat{\rho}_m r + \frac{C}{r^2} \quad U(r) = \frac{1}{6}\hat{\rho}_m r^2 - \frac{C}{r} + D \quad r \in [R_1, R_2] \quad (10.110)$$

$$g(r) = \frac{E}{r^2} \quad U(r) = -\frac{E}{r} \quad r \in [R_2, +\infty) \quad (10.111)$$

We have now four unknowns  $B, C, D, E$ . Since the two fields  $g$  and  $U$  must be continuous at  $r = R_1$  and  $r = R_2$ , then we have four constraints and that will allow us to determine the four unknowns.

Let us start with the continuity at  $r = R_1$ :

$$g(r = R_1) = \frac{1}{3}\hat{\rho}_c R_1 = \frac{1}{3}\hat{\rho}_m R_1 + \frac{C}{R_1^2} \quad (10.112)$$

$$U(r = R_1) = \frac{1}{6}\hat{\rho}_c R_1^2 + B = \frac{1}{6}\hat{\rho}_m R_1^2 - \frac{C}{R_1} + D \quad (10.113)$$

The first equation yields

$$C = \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m)$$

which we plug in the second equation:

$$\frac{1}{6}\hat{\rho}_c R_1^2 + B = \frac{1}{6}\hat{\rho}_m R_1^2 - \frac{R_1^2}{3}(\hat{\rho}_c - \hat{\rho}_m) + D$$

$$\frac{1}{6}(\hat{\rho}_c - \hat{\rho}_m)R_1^2 + B = -\frac{R_1^2}{3}(\hat{\rho}_c - \hat{\rho}_m) + D$$

$$\left(\frac{1}{6} + \frac{1}{3}\right)(\hat{\rho}_c - \hat{\rho}_m)R_1^2 + B = D$$

$$\frac{1}{2}(\hat{\rho}_c - \hat{\rho}_m)R_1^2 + B = D$$

We cannot go any further so we turn to  $r = R_2$ :

$$g(r = R_2) = \frac{1}{3}\hat{\rho}_m R_2 + \frac{C}{R_2^2} = \frac{E}{R_2^2} \quad (10.114)$$

$$U(r = R_2) = \frac{1}{6}\hat{\rho}_m R_2^2 - \frac{C}{R_2} + D = -\frac{E}{R_2} \quad (10.115)$$

in which we insert the known value of  $C$ :

$$\frac{1}{3}\hat{\rho}_m R_2 + \frac{1}{R_2^2} \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) = \frac{E}{R_2^2} \quad (10.116)$$

$$\frac{1}{6}\hat{\rho}_m R_2^2 - \frac{1}{R_2} \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) + D = -\frac{E}{R_2} \quad (10.117)$$

Multiplying the first equation by  $R_2^2$  and the second one by  $R_2$ :

$$\frac{1}{3}\hat{\rho}_m R_2^3 + \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) = E \quad (10.118)$$

$$\frac{1}{6}\hat{\rho}_m R_2^3 - \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) + DR_2 = -E \quad (10.119)$$

The first equation then gives us  $E$ :

$$E = \frac{1}{3}\hat{\rho}_m(R_2^3 - R_1^3) + \frac{1}{3}R_1^3\hat{\rho}_c$$

So in the end, we have

$$C = \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) \quad (10.120)$$

$$E = \frac{1}{3}\hat{\rho}_m(R_2^3 - R_1^3) + \frac{1}{3}R_1^3\hat{\rho}_c \quad (10.121)$$

$$D = \frac{1}{R_2} \left( -E - \frac{1}{6}\hat{\rho}_m R_2^3 + \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) \right) \quad (10.122)$$

$$= \frac{1}{R_2} \left( -\frac{1}{3}\hat{\rho}_m(R_2^3 - R_1^3) - \frac{1}{3}R_1^3\hat{\rho}_c - \frac{1}{6}\hat{\rho}_m R_2^3 + \frac{R_1^3}{3}(\hat{\rho}_c - \hat{\rho}_m) \right) \quad (10.123)$$

$$= \frac{1}{6R_2} [(-2R_2^3 + 2R_1^3 - R_2^3 - 2R_1^3)\hat{\rho}_m + (-2R_1^3 + 2R_1^3)\hat{\rho}_c] \quad (10.124)$$

$$= \frac{1}{6R_2} (-3R_2^3\hat{\rho}_m) \quad (10.125)$$

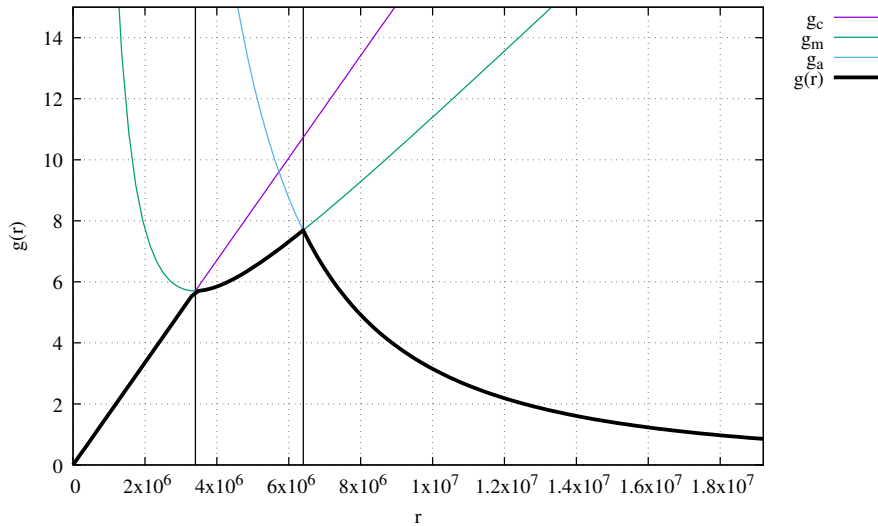
$$= -\frac{R_2^2}{2}\hat{\rho}_m \quad (10.126)$$

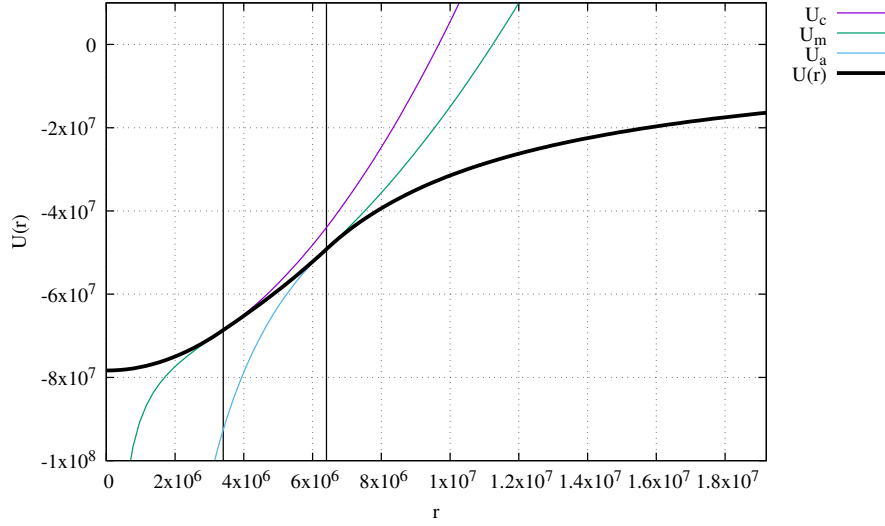
$$B = D - \frac{1}{2}(\hat{\rho}_c - \hat{\rho}_m)R_1^2 \quad (10.127)$$

$$= -\frac{R_2^2}{2}\hat{\rho}_m - \frac{1}{2}(\hat{\rho}_c - \hat{\rho}_m)R_1^2 \quad (10.128)$$

$$= \frac{1}{2}(R_1^2 - R_2^2)\hat{\rho}_m - \frac{1}{2}R_1^2\hat{\rho}_c \quad (10.129)$$

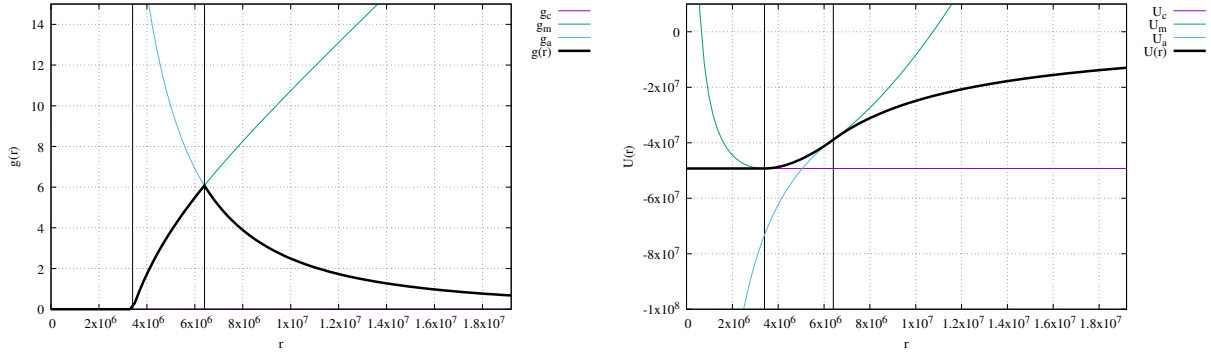
We set  $R_1 = 3400\text{km}$ ,  $R_2 = 6400\text{km}$ ,  $\rho_c = 6000\text{kg/m}^3$  and  $\rho_m = 4000\text{kg/m}^3$  and obtain the following fields:





We find that the fields fulfill all conditions and are continuous, as expected.

Also, setting  $\rho_c = 0$  (i.e. the planet is a hollow sphere), we recover the following fields:



These are similar to the results presented in Appendix A of Thieulot [1259].

We can also consider the case of a constant density planet, i.e.  $\rho_m = \rho_c = \rho_0 = 5000\text{kg/m}^3$ . In that case  $C = 0$ ,

$$E = \frac{1}{3}\hat{\rho}_0 R_2^3 = \frac{4\pi}{3}\mathcal{G}\rho_0 R_2^3 = M\mathcal{G}$$

$$D = -\frac{1}{2}R_2^2\hat{\rho}_0$$

$$B = -\frac{1}{2}R_2^2\hat{\rho}_0$$

so that

$$g(r) = \frac{4\pi}{3}\mathcal{G}\rho_0 r \quad U(r) = \frac{1}{6}\hat{\rho}_0 r^2 - \frac{1}{2}R_2^2\hat{\rho}_0 \quad r \in [0, R_1] \quad (10.130)$$

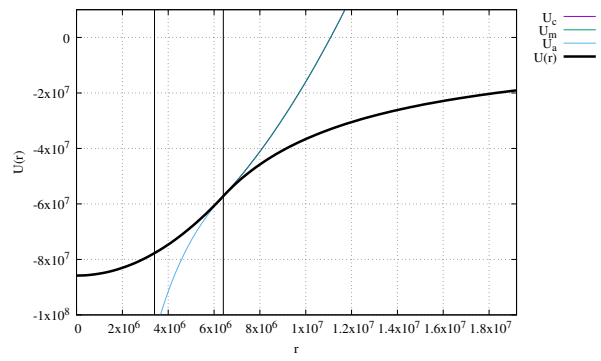
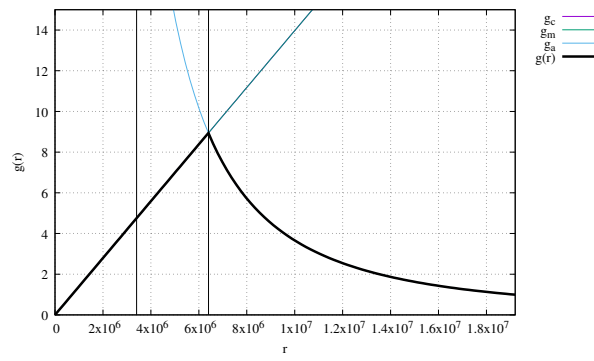
$$g(r) = \frac{4\pi}{3}\mathcal{G}\rho_0 r \quad U(r) = \frac{1}{6}\hat{\rho}_0 r^2 - \frac{1}{2}R_2^2\hat{\rho}_0 \quad r \in [R_1, R_2] \quad (10.131)$$

$$g(r) = \frac{\mathcal{G}M}{r^2} \quad U(r) = -\frac{\mathcal{G}M}{r} \quad r \in [R_2, +\infty) \quad (10.132)$$

Unsurprisingly we obtain the same expressions inside the core and the mantle and we recover the expressions of problem 14 (see previous section).

In the end the fields are as follows:





## 10.12 Gravity forward calculations in practice

$$\vec{g}(\vec{r}) = -\vec{\nabla}U(\vec{r})$$

with

$$U(\vec{r}) = \mathcal{G} \iiint \frac{\rho(\vec{r}')}{|\vec{r} - \vec{r}'|} dV$$

with (in a Cartesian coordinates system):

$$|\vec{r} - \vec{r}'| = \sqrt{(x - x')^2 + (y - y')^2 + (z - z')^2} \quad dV = dx' dy' dz'$$

We then have

$$\begin{aligned} g_x(x, y, z) &= -\mathcal{G} \iiint \frac{\partial}{\partial x} \frac{\rho(\vec{r}')}{|\vec{r} - \vec{r}'|} dx' dy' dz' \\ &= \mathcal{G} \iiint \frac{\rho(\vec{r}')(x - x')}{|\vec{r} - \vec{r}'|^3} dx' dy' dz' \\ g_y(x, y, z) &= -\mathcal{G} \iiint \frac{\partial}{\partial y} \frac{\rho(\vec{r}')}{|\vec{r} - \vec{r}'|} dx' dy' dz' \\ &= \mathcal{G} \iiint \frac{\rho(\vec{r}')(y - y')}{|\vec{r} - \vec{r}'|^3} dx' dy' dz' \\ g_z(x, y, z) &= -\mathcal{G} \iiint \frac{\partial}{\partial z} \frac{\rho(\vec{r}')}{|\vec{r} - \vec{r}'|} dx' dy' dz' \\ &= \mathcal{G} \iiint \frac{\rho(\vec{r}')(z - z')}{|\vec{r} - \vec{r}'|^3} dx' dy' dz' \end{aligned}$$

In order to compute the gravity tensor  $\mathbf{T} = \vec{\nabla}\vec{g}$ , also called gravitational gravity tensor [1090] or Marussi tensor, we need to compute  $\frac{\partial^2 U}{\partial \alpha \partial \beta}$  with  $\alpha, \beta = x, y, z$  and we obtain (See Arroyo *et al.* (2015) [30]):

$$T_{xx}(x, y, z) = -\mathcal{G} \iiint \frac{3(x - x')^2 - (\vec{r} - \vec{r}')^2}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.133)$$

$$T_{yy}(x, y, z) = -\mathcal{G} \iiint \frac{3(y - y')^2 - (\vec{r} - \vec{r}')^2}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.134)$$

$$T_{zz}(x, y, z) = -\mathcal{G} \iiint \frac{3(z - z')^2 - (\vec{r} - \vec{r}')^2}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.135)$$

$$T_{xy}(x, y, z) = -\mathcal{G} \iiint \frac{3(x - x')(y - y')}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.136)$$

$$T_{xz}(x, y, z) = -\mathcal{G} \iiint \frac{3(x - x')(z - z')}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.137)$$

$$T_{yz}(x, y, z) = -\mathcal{G} \iiint \frac{3(y - y')(z - z')}{|\vec{r} - \vec{r}'|^5} \rho(\vec{r}') dx' dy' dz' \quad (10.138)$$

Note that the trace satisfies the Laplace equation  $T_{xx} + T_{yy} + T_{zz} = \Delta U = 0$ .

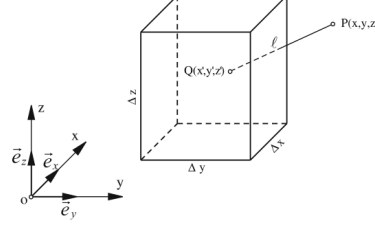
redo all calculations to be sure

Unless the geometry is conveniently chosen with lots of symmetry and the density is also very simple the above integrals cannot be computed analytically and one must resort to numerical integration based on a tessellation of the space, such as prisms or tesseroids.

## Prisms

The gravitational potential  $U$  of a right rectangular parallelepiped (prism) of homogeneous mass-density  $\rho_0$  is described by Newton's integral [556]

$$U(x, y, z) = \mathcal{G}\rho_0 \int_{z_1}^{z_2} \int_{y_1}^{y_2} \int_{x_1}^{x_2} \frac{1}{|\vec{r} - \vec{r}'|} dx' dy' dz' = \mathcal{G}\rho_0 \int_{z_1}^{z_2} \int_{y_1}^{y_2} \int_{x_1}^{x_2} \frac{1}{\sqrt{(x-x')^2 + (y-y')^2 + (z-z')^2}} dx' dy' dz'$$



Taken from [556].

The denominator is the distance between the computation point  $P(x, y, z)$  and the running integration point  $Q(x', y', z')$ . The coordinate axes have been assumed to be parallel to the edges of the prism, which extends between the coordinate surfaces related to the bounds  $x_1, x_2, y_1, y_2, z_1, z_2$ . It is well known that the integral can be solved analytically (Mader (1951) [822], see also Nagy *et al.* [922, 921]), resulting in the formula for the potential  $U(x, y, z)$ :

$$U(x, y, z) = \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( A + B + C - \frac{1}{2}D \right)$$

with

$$\begin{aligned} A &= (x - x_i)(y - y_j) \ln \left| \frac{z - z_k + r_{ijk}}{\sqrt{(x - x_i)^2 + (y - y_j)^2}} \right| \\ B &= (y - y_j)(z - z_k) \ln \left| \frac{x - x_i + r_{ijk}}{\sqrt{(y - y_j)^2 + (z - z_k)^2}} \right| \\ C &= (x - x_i)(z - z_k) \ln \left| \frac{y - y_j + r_{ijk}}{\sqrt{(z - z_k)^2 + (x - x_i)^2}} \right| \\ D &= (x - x_i)^2 \arctan \frac{(y - y_j)(z - z_k)}{(x - x_i)r_{ijk}} \\ &\quad + (y - y_j)^2 \arctan \frac{(z - z_k)(x - x_i)}{(y - y_j)r_{ijk}} \\ &\quad + (z - z_k)^2 \arctan \frac{(x - x_i)(y - y_j)}{(z - z_k)r_{ijk}} \end{aligned}$$

and

$$r_{ijk} = \sqrt{(x - x_i)^2 + (y - y_j)^2 + (z - z_k)^2}$$

**Remark.** The direct application of this equation will fail when the computation point  $P$  is situated on an edge or on a corner of the prism; the respective limit values have been derived by Nagy *et al.* [922, 921]

The gravity vector and tensor can then be computed [1004, 30, 279]:

$$\begin{aligned}
g_x(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \times \\
&\quad \left[ (y - y_j) \ln((z - z_k) + r_{ijk}) + (z - z_k) \ln((y - y_j) + r_{ijk}) - (x - x_i) \arctan \frac{(y - y_j)((z - z_k))}{(x - x_i)r_{ijk}} \right] \\
g_y(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \times \\
&\quad \left[ (z - z_k) \ln((x - x_i) + r_{ijk}) + (x - x_i) \ln((z - z_k) + r_{ijk}) - (y - y_j) \arctan \frac{(x - x_i)(z - z_k)}{(y - y_j)r_{ijk}} \right] \\
g_z(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \times \\
&\quad \left[ (y - y_j) \ln((x - x_i) + r_{ijk}) + (x - x_i) \ln((y - y_j) + r_{ijk}) - (z - z_k) \arctan \frac{(x - x_i)(y - y_j)}{(z - z_k)r_{ijk}} \right] \\
T_{xx}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \arctan \frac{(y - y_j)(z - z_k)}{(x - x_i)r_{ijk}} \\
T_{yy}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \arctan \frac{(x - x_i)(z - z_k)}{(y - y_j)r_{ijk}} \\
T_{zz}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \arctan \frac{(x - x_i)(y - y_j)}{(z - z_k)r_{ijk}} \\
T_{xy}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \ln((z - z_k) + r_{ijk}) \\
T_{xz}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \ln((y - y_j) + r_{ijk}) \\
T_{yz}(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \ln((x - x_i) + r_{ijk})
\end{aligned}$$

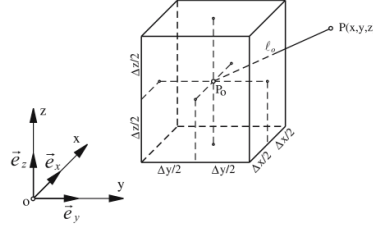
Note that Heck & Seitz [556] report that the logarithmic terms can be transformed in order to provide a better numerical stability and then

$$\begin{aligned}
g_x(x, y, z) &= \mathcal{G}\rho_0 \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left[ (y - y_j) \ln \left| \frac{(z - z_k) + r_{ijk}}{\sqrt{(x - x_i)^2 + (y - y_j)^2}} \right| \right. \\
&\quad \left. + (z - z_k) \ln \left| \frac{(y - y_j) + r_{ijk}}{\sqrt{(x - x_i)^2 + (z - z_k)^2}} \right| - (x - x_i) \arctan \frac{(y - y_j)(z - z_k)}{(x - x_i)r_{ijk}} \right] \quad (10.139)
\end{aligned}$$

The gravitational potential of the homogeneous rectangular prism, neglecting terms of order four and higher in  $x, y, z$ , is then given by MacMillan's (1930)<sup>19</sup> formula:

$$U(x, y, z) = \mathcal{G}\rho_0 \Delta_x \Delta_y \Delta_z \left[ \frac{1}{l_0} + \frac{3(x_0 - x)^2 - l_0^2}{24l_0^5} \Delta_x^2 + \frac{3(y - y_0)^2 - l_0^2}{24l_0^5} \Delta_y^2 + \frac{3(z - z_0)^2 - l_0^2}{24l_0^5} \Delta_z^2 + \mathcal{O}(\Delta^4) \right]$$

<sup>19</sup>MacMillan WD (1930) Theoretical Mechanics, vol 2: the Theory of the potential. McGraw-Hill, New York (reprinted by Dover Publications, New York 1958)



Taken from [556].

It is obvious that the zero-order approximation is identical with the potential of a point-mass at  $P_0$  when the total mass of the prism  $m = \rho_0 \Delta_x \Delta_y \Delta_z$  is concentrated at its geometrical centre  $P_0$ :

$$U(x, y, z) = \mathcal{G} \rho_0 \Delta_x \Delta_y \Delta_z \frac{1}{l_0}$$

It is also common [348] to look at:

- the differential curvature magnitude (*DCM*) which is also known as the horizontal directive tendency, computed by a combination of components of tensor  $T_{xx}$ ,  $T_{xy}$  and  $T_{yy}$ . It emphasizes greatly the effects of shallower sources (Saad, 2006);

$$DCM = \sqrt{(T_{xx} - T_{yy})^2 + 2T_{xy}^2}$$

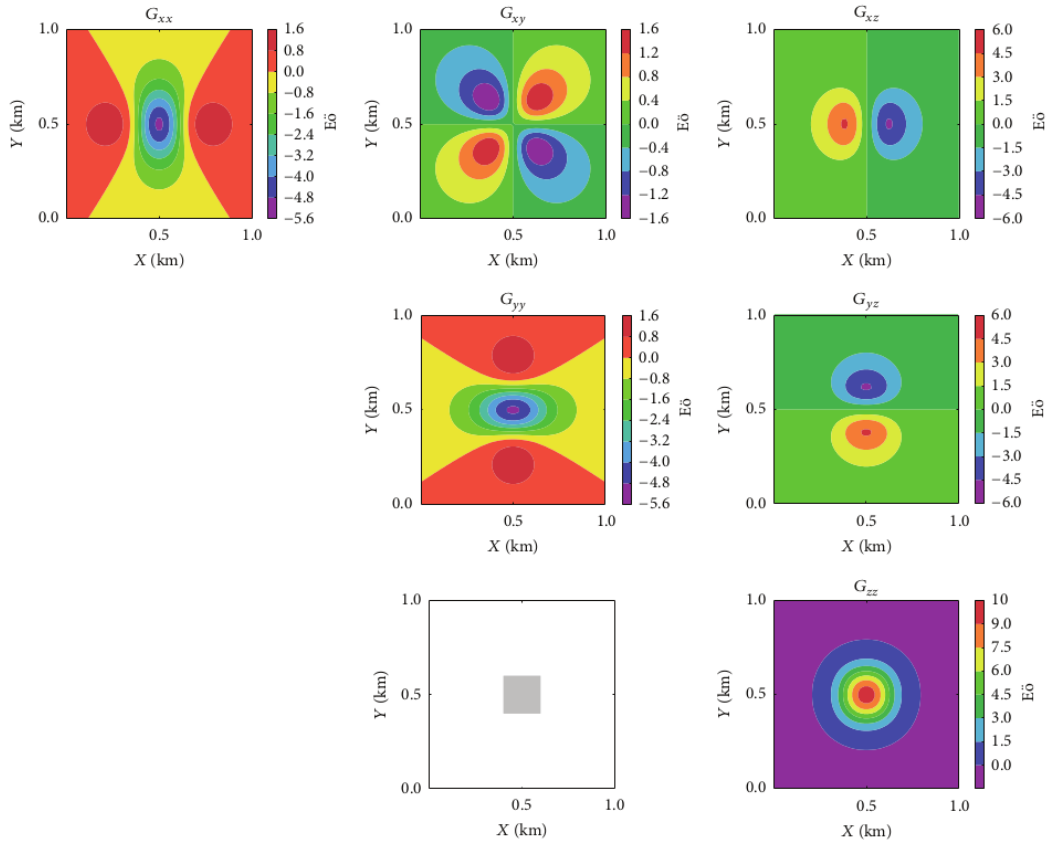
- the horizontal gradient magnitude (*HGM*) of  $g_z$  can be computed from the horizontal derivative components of  $g_z$  and can be used as edge detector or to map the body outline as it verifies the prism boundaries

$$HGM = \sqrt{T_{zx}^2 + T_{zy}^2} = \sqrt{\left(-\frac{\partial g_z}{\partial x}\right)^2 + \left(-\frac{\partial g_z}{\partial y}\right)^2}$$

- the total gradient magnitude (*TGM*) is computed from the three derivatives of vertical component of gravity:

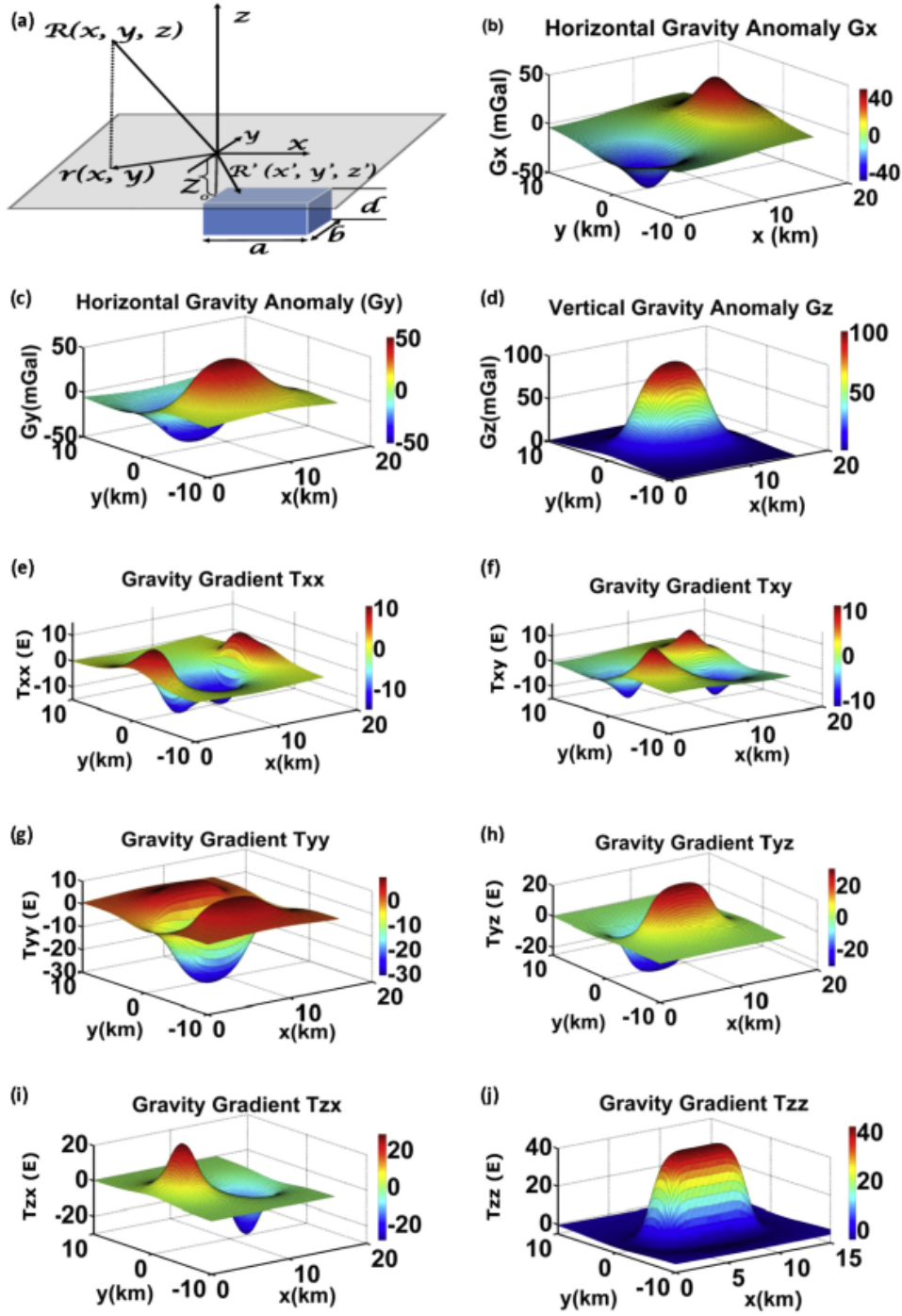
$$TGM = \sqrt{T_{zx}^2 + T_{zy}^2 + T_{zz}^2} = \sqrt{\left(-\frac{\partial g_z}{\partial x}\right)^2 + \left(-\frac{\partial g_z}{\partial y}\right)^2 + \left(-\frac{\partial g_z}{\partial z}\right)^2}$$

**Example 1** The result of calculating the components of a prism measuring  $200\text{m}^3$  at a height of  $0.01\text{ km}$ , with an observation mesh of  $1\text{km}\times 1\text{km}$ , and discretized every  $20\text{m}$  is shown hereunder:

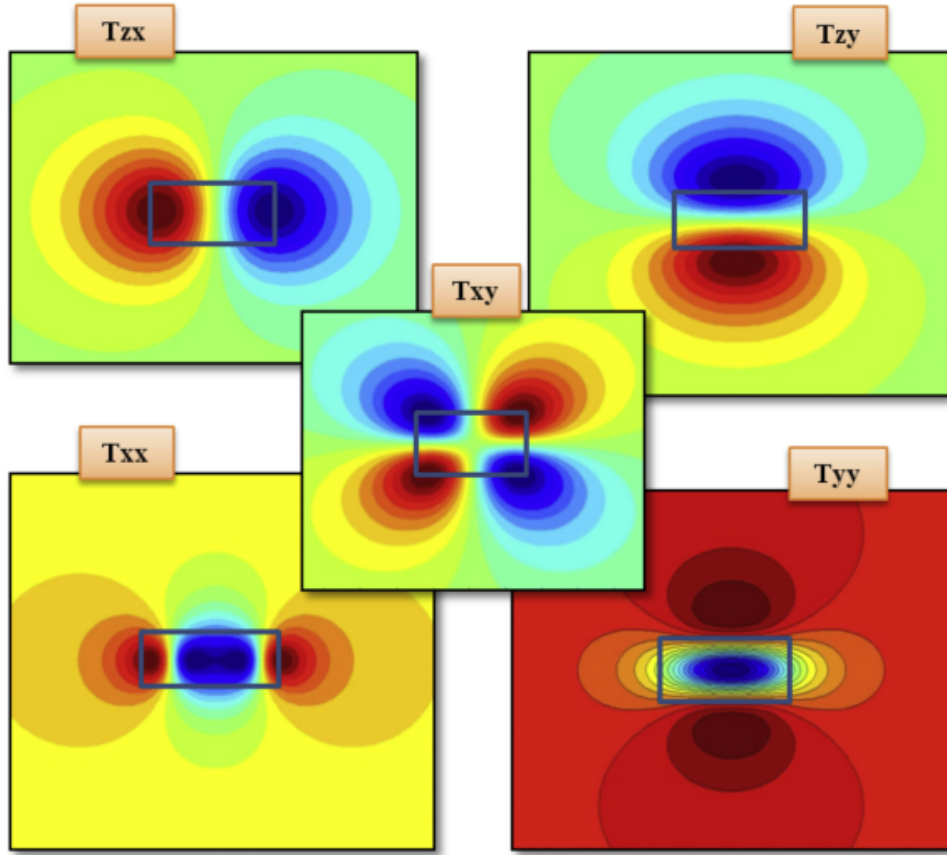


Taken from Arroyo *et al.* (2015) [30]. Gravity gradient response for a prism buried a depth of  $100\text{m}$ ,  
Each side having a length of  $200\text{ m}$  and constant density contrast of  $100\text{ kg/m}^3$ .

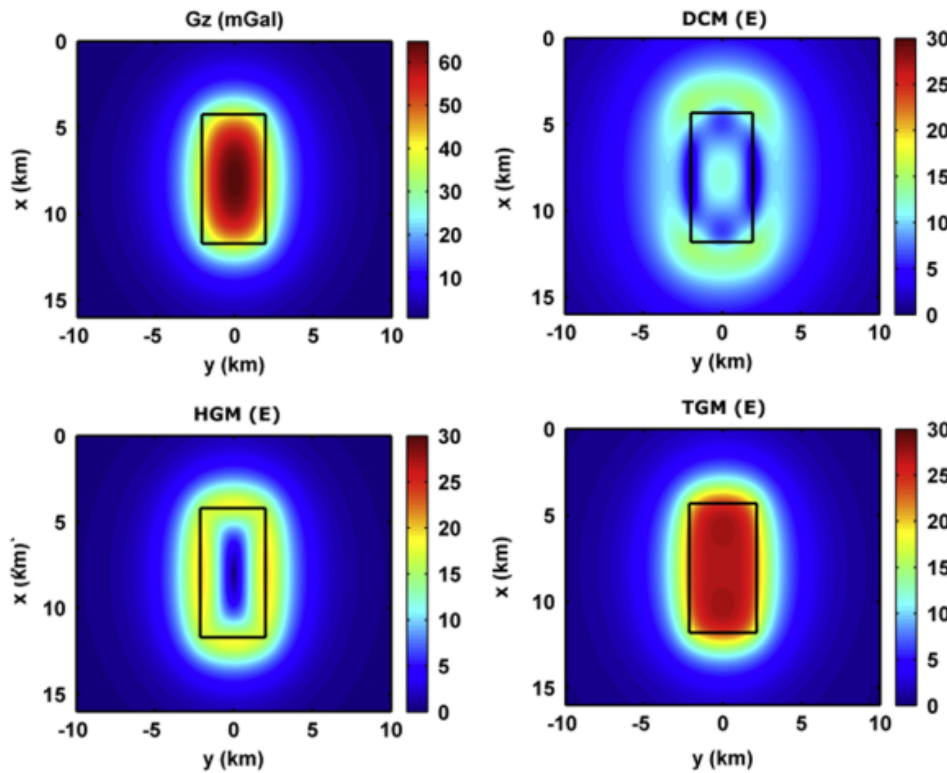
**Example 2** . Buried prisms of size 8x4x1 km along x, y and z directions respectively. Density is 2700. Mapped gravity field and its gradient on a plane of constant 1000m height.



(a) A model containing a prism and (b-d): corresponding gravity vector components and (e-j) GGT components with sampling interval of 0.2 km in x and y directions. Taken from [348]



A map view of complex behavior of gravity gradients for prism model. Taken from [348]



Computed vertical gravity component  $G_z$  and three invariants map of HGM, DCM and TGM for given prism model. HGM = Horizontal Gradient Magnitude, DCM = Differential Curvature Magnitude and TGM = Total Gradient Magnitude. Taken from [348]

#### Relevant Literature:

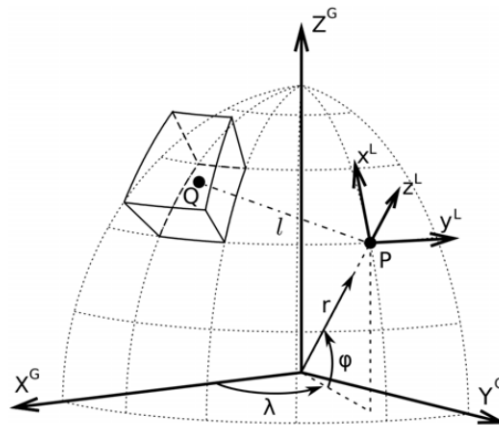
- Analytic Expressions for the Gravity Gradient Tensor of 3D Prisms with Depth-Dependent Density [646]



- New computationally efficient quadrature formulas for triangular prism elements [734]
- 3D Gravity Modeling of Complex Salt Features in the Southern Gulf of Mexico [931].
- Spherical prism gravity effects by Gauss-Legendre quadrature integration [32]
- Perturbing effects of sub-lithospheric mass anomalies in GOCE gravity gradient and other gravity data modelling: Application to the Atlantic-Mediterranean transition zone [424]

## Tesseroids

A tesseroid, or spherical prism, is segment of a sphere.



Taken from [1298]

Tesseroids is a collection of command-line programs for modeling the gravitational potential, acceleration, and gradient tensor. Tesseroids supports models and computation grids in Cartesian and spherical coordinates. <https://tesseroids.readthedocs.io/en/stable/>

### Relevant Literature:

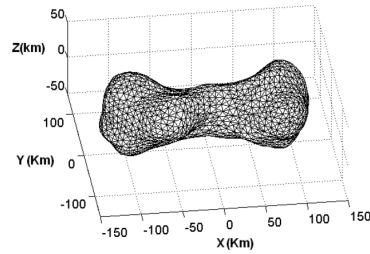
- Forward modeling and inversion of gravitational fields in spherical coordinates, L. Uieda, Phd Thesis [1297]
- Tesseroids: Forward-modeling gravitational fields in spherical coordinates [1296]
- Optimal forward calculation method of the Marussi tensor due to a geologic structure at GOCE height [1298]
- Optimized formulas for the gravitational field of a tesseroid [498]

## Tetrahedra

Werner and Scheeres [1350] (1997) derived analytical expressions for the gravity potential, field and tensor generated by any polyhedron. These can be applied to tetrahedra.

Metherell and Quinn [866] (1986) derived analytical expressions for the gravity field generated by so-called 111-tetrahedra. Note the corrections in Carré, Metherell, and Quinn [209] (1986).

Chanut, Aljbaae, and Carruba [218] (2015) used a mascon approach (point mass approach?) on tetrahedra making up an asteroid:



Taken from Chanut, Aljbaae, and Carruba [218]. Polyhedron model 3D of asteroid (216) Kleopatra. The shape was built with 4092 faces.

This is implemented in [STONE](#) 113.

## Other shapes

### Relevant Literature

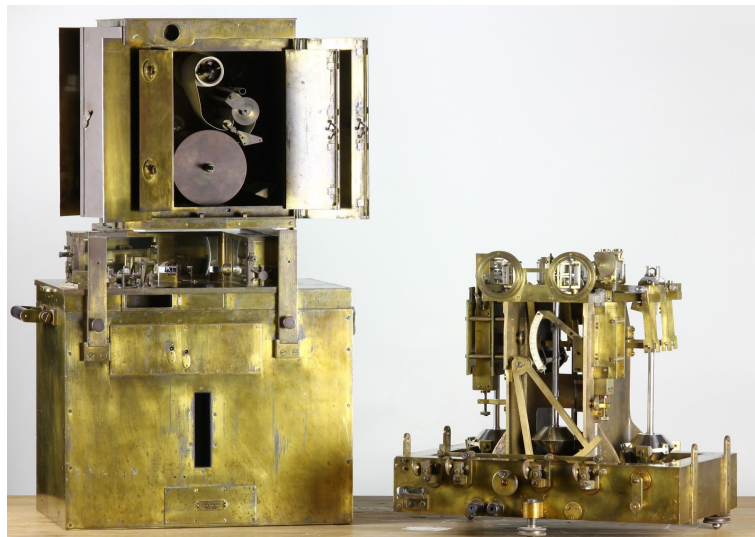
Rapid Gravity Computations for Two-Dimensional Bodies [1233]

Rapid Computation of Gravitational Attraction of Three-Dimensional Bodies of Arbitrary Shape [1232]

## 10.13 Instruments to measure gravity

### Gravity meters

include here pic of Vening Meinesz



Taken from website<sup>20</sup>. The pendulum apparatus of Vening Meinesz, also known as "Het Gouden Kalf" (the Golden Calf). Positioned on the left side is the protective casing with the recording instrument on top. On the right side is the pendulum apparatus with the three pendulums at the back.

See video by Bart Root: <https://youtu.be/SVTJA3KAnck?si=-OZ11PnHQwHy0kE1>

**Absolute gravity measurements** After a time  $t$  an object has fallen by a distance  $x$  in a gravity field  $g$  with  $x = gt^2/2$  so that  $g = 2x/t^2$ .

<sup>20</sup><http://deepearthscience.blogspot.com/2016/06/the-gravimeter-of-professor-vening.html>



by Micro-g LaCoste. The FG5<sup>21</sup> operates by using a free-fall method. An object is dropped inside a vacuum chamber and its position is monitored very accurately using a laser interferometer. Dropping chamber of 33cm. Accuracy of approx.  $2\mu\text{Gal}$ .

## Planes

HALO – the German High Altitude and Long Range Research Aircraft (HALO).



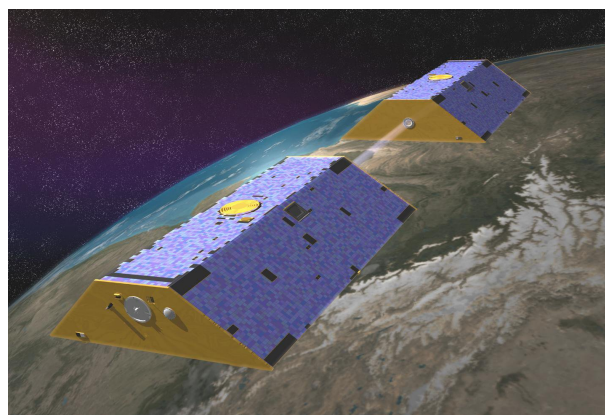
**GEOHALO** –  
Multidisciplinary geoscientific airborne  
mission on the HALO aircraft over Italy  
June 2012



GFZ's airborne gravimeter on HALO:  
CHEKAN - AM

## Satellites

**GRACE** Note that GRACE consists of two satellites which are in a low orbit and the distance between them is accurately measured. Changes in this separation are caused by increases and decreases in gravity.



[https://upload.wikimedia.org/wikipedia/commons/e/e6/GRACE\\_artist\\_concept.jpg](https://upload.wikimedia.org/wikipedia/commons/e/e6/GRACE_artist_concept.jpg)

Examples of applications using GRACE data:

<sup>21</sup><http://microglacoste.com/product/fg5-x-absolute-gravimeter/>

- Inference of mantle viscosity from GRACE and relative sea level data [983]
- Exploring the uncertainty in GRACE estimates of the mass redistributions at the Earth surface: implications for the global water and sea level budgets [96]

**GOCE** Gravity Field and Steady-State Ocean Circulation Explorer (GOCE) was the first of ESA's Living Planet Programme satellites intended to map in unprecedented detail the Earth's gravity field with a spatial resolution up to 80 km. The spacecraft's primary instrumentation was a highly sensitive gravity gradiometer consisting of three pairs of accelerometers which measured gravitational gradients along three orthogonal axes.



Examples of applications using GOCE data:

- GOCE gravitational gradients along the orbit [121]
- Moho Estimation Using GOCE Data: A Numerical Simulation [1057]
- Global Moho from the combination of the CRUST2.0 model and GOCE data [1058]
- Advancements in satellite gravity gradient data for crustal studies [359]
- Sensitivity of GOCE Gravity Gradients to Crustal Thickness and Density Variations: Case Study for the Northeast Atlantic Region [358]
- Mapping the mass distribution of Earth's mantle using satellite-derived gravity gradients [972]
- GOCE gravity gradient data for lithospheric modeling [122]
- Exploration of tectonic structures with GOCE in Africa and across-continentals [131]
- GEMMA: An Earth crustal model based on GOCE satellite data [1056]
- GOCE data, models, and applications: A review [861]
- Geological units and Moho depth determination in the Western Balkans exploiting GOCE data [1101]
- The combined inversion of seismological and GOCE gravity data: New insights into the current state of the Pacific lithosphere and upper mantle [1272]

## 10.14 Gravity anomalies

### Caves and cavities

Underground man-made structures (bunkers)

Mineral deposits

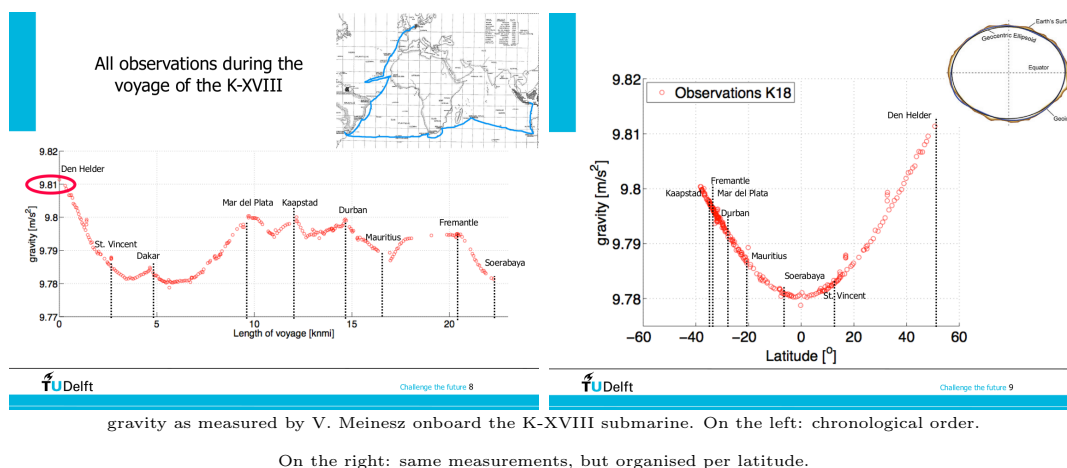
Impact craters

Salt layers

## 10.15 Gravity reductions

In gravity work, more than in any other branch of geophysics, large and (in principle) calculable effects are produced by sources which are not of direct geological interest. These effects are removed by reductions involving sequential calculation of a number of recognized quantities.

**Latitude correction** Because of the shape of the Earth, latitude has a large effect on gravity measurements, as visible in the following figure:



It is obvious that we are not interested in this long wavelength pattern, but rather in the deviations from it.

The formula is as follows:

$$g_n = 978031.85(1.0 + 0.005278895 \sin^2(lat) + 0.000023462 \sin^4(lat))(\text{mgal})$$

where  $lat$  is the latitude.

**Free-air correction** After subtracting the above signal, the observed gravity will be due in part to the height of the gravity station above the sea-level reference surface. An increase in height implies an increase in distance from the Earth's centre of mass and the effect is negative for stations above sea level ( $g \propto r^{-2}$ ). The free-air correction is thus positive and the quantity obtained after applying both the latitude and free-air corrections is termed the free-air anomaly or free-air gravity.

**Bouguer correction**

**Terrain correction** material10.pdf



## 10.16 How not to think about gravity (or Earth Sciences)



Sai

So gravity is really a lie ?

on Sat Like Reply More



Tyler

Not a lie but disorder.... disfunction.... the atoms are shooting all over the place and poles point in any direction.

on Sat Like Reply More



Thomas

Gravity exist but not as we know it. Gravity is an electrical phenomenon called incoherent dielectric acceleration in the form of universal compression radiating from the black sun 😊

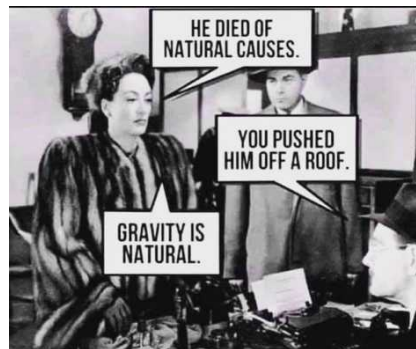
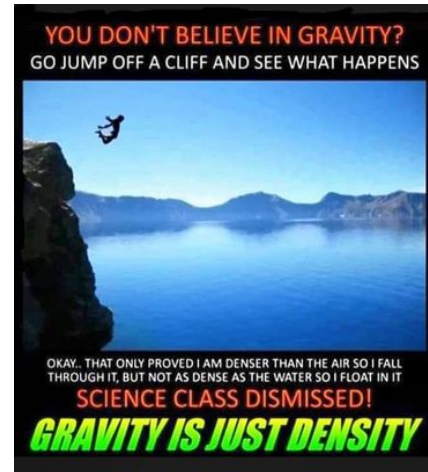
EARTH is a realm, it is not a planet. It is not an object, therefore, it has no edge. Earth would be more easily defined as a system environment. Earth is also a machine, it is a Tesla coil. The sun and moon are powered wirelessly with the electromagnetic field (the Aether). This field also suspends the celestial spheres with electro-magnetic levitation. Electromagnetic levitation disproves gravity because the only force you need to counter is the electromagnetic force, not gravity. The stars are attached to the FIRMAMENT.



jeff

1 year ago

As for gravity one simple question to debunk gravity because the example you gave is density. The heavier something is the harder it falls that's density not gravity. So now to debunk gravity dumbass. How can gravity hold 350 trillion gallons of water to a spinning ball yet gravity can let a fly roam free or when my dog shakes its wet fur why does the water fly everywhere and not stick to the ground ?????? That's common sense not your dumbass shit.



Clarence

Plate tectonics is false. Like with the magnetic pole flip, it lacks a plausible mechanism to produce the energy needed to move kajillions of cubic miles of crust a bazillion miles, and that at a speed high enough to thrust up the Rockies and the Himalayas. I don't buy it.

I don't believe in magic.

9 hrs Like Reply More



Kenny

7 hrs · Facebook for Android · 🌐

I always wondered if aliens are living on the other side of the disc. Maybe they don't know they are upside down. 🤔🤔🤔

**F** Gravity can be debunked with one statement. No matter in this universe can create uniform force because matter is always irregular (even in earth core if it exists it would be fluid form because of the pressure applied) . When feather and steel ball fall at same acceleration in vaccum shows the gravity effect on the object is uniform which is impossible with any matter. My only conclusion is all matter inthis universe is falling at same speed if earth is flat. All my senses tell me earth is infinite horizon. Even if it is a potato then gravity can only be buoyancy effect of matter occupy the space and all the objects try to get to the center of that space. That means gravity varies bases on size of the planet not on the mass of the planet.

Free Speech • 2 days ago

**P** Gravity is 'their' top trump. I personally believe there is a little more to it than density alone, something to do with the speed of light, a reciprocal perhaps. I also think we are upside down, in the anti matter cycle if you prefer. The wrong side of eternity.

Paul • 3 days ago (edited)

**Jacqui** 20 hrs

Has there been any experiment that proves that water can stick to a sphere? If so, I'd like to see it. Also, quick question: if "gravity" can cause the oceans to adhere to the globe, wouldn't the same laws apply to a small ball in front of me? I've tried sticking water to one but it just ain't happening? Do the laws of gravity only apply to large balls? And if so, at what point are balls considered large? (No vulgar replies please - haha)

Like Comment

10

**Cindy** To my knowledge there is no repeatable experiment to prove water sticks to a spinning ball of any size.  
Like Reply • 2 - 19 hrs

**Jerry** Maybe in a vacuum chamber? I imagine, in a vacuum, the water would pool up around the ball and the ball would just spin freely underneath it while the water stays in the same position.  
Like Reply • 3 - 19 hrs - Edited

**Jacqui** Is there a video of this experiment?  
Like Reply • 1 - 19 hrs

**Jerry** I don't know, I just made that hypothesis. I've thought about it before, but I don't think there's been an experiment performed.  
Like Reply • 19 hrs - Edited

**Jacqui** Jerry - do one and video it!  
Like Reply • 1 - 19 hrs

**Jerry** It's a good idea, but I think it's beyond my capability of time and resources. I have a son coming to this earth next month, so he's going to be my priority for a while. Maybe one day.  
Like Reply • 3 - 19 hrs

**Cindy** I would say let's demand an experiment from the scientific community that swears by gravity.  
Like Reply • 4 - 19 hrs

**Cindy** They won't because they can't.  
Like Reply • 3 - 18 hrs

**John** It's all down to relative mass isn't it? and space is supposedly a vacuum, if it exists.  
Like Reply • 3 - 14 hrs

**Tippr** Gravity doesn't exist, water forms to it's container, unless ball is frozen. lol  
Like Reply • 1 - 14 hrs

**We all agree Gravity is nonsense but... Here's my question: is Density and Buoyancy enough? Or is magnetism at play as well?**

**Reece**  
Density and buoyancy are the effect, electromagnetism is the cause. All in line with the Torus field.  
2 hrs Like Reply More

**Tanakah**  
You said correctly... Magnetic fields all around us.. Gravity is a scam  
3 hrs Like Reply More

**Law**  
I don't think buoyancy and density is enough since it doesn't explain the reason things "fall" "down". Incoherent magnetism, dielectric acceleration and electrostatics do explain that in my opinion.  
3 hrs Like Reply More

**Mark**  
EMF plays a vital role on our Plane Flat Earth's Realm, movement of luminaries, light (visible and non visible), energy, eclipses, tides, spectrums, aether, the light described in the Bible prior on creation of luminary Sun, plane of inertia, line and point of inertia of this Realm, torus field and a lot more. All are true science that God created.  
3 hrs Like Reply More

**Mark**  
I agree with your thoughts on this. I feel that density and buoyancy are enough, however very possible that magnetism comes into play somehow. But gravity is definitely a bunch of b.s., they need gravity to explain their claims of objects orbiting other objects. It seems to be at once both a force of adhesion and a force of suspension and rotation. It's the Swiss knife go to reason used for any claims that they cannot explain.  
3 hrs Like Reply More

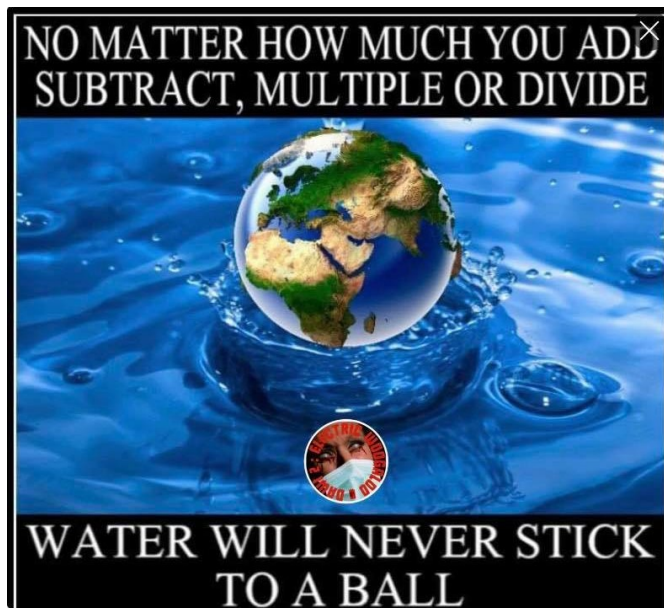
**Jessie**  
Aether...  
3 hrs Like Reply More

**Eddie**  
Even the smallest of creatures, like the bumblebee, use magnetism to hover/fly. Resonance of frequency within the aether, achieves this.  
3 hrs Like Reply More

**Bert**  
Yes magnetism steers the tide  
3 hrs Like Reply More

**Ray**  
Water is not magnetic  
45 mins Like Reply More

**Mike**  
Gravity is a word to misguide us so the true power of our world is never discovered. Would it eventually destroy us once known on a massive scale. Possibly, as those who occupied our flat earth before us are gone with no history that we know of or are told.  
2 hrs Like Reply More





## Americans when Newton discovered gravity :



Sam

Gravity is a cover-up concept for buoyancy and density.



last Sun Like More



Keanu Reeves and Bernie Sanders are sitting 1.7 m apart.  
Keanu has a mass of 76 kg. Bernie has a mass of 64 kg.

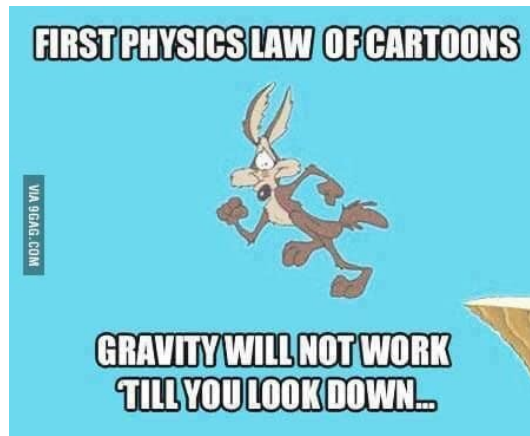


(a) What is the gravitational force that Keanu exerts on Bernie?

N

(b) What is the gravitational force that Bernie exerts on Keanu?

N



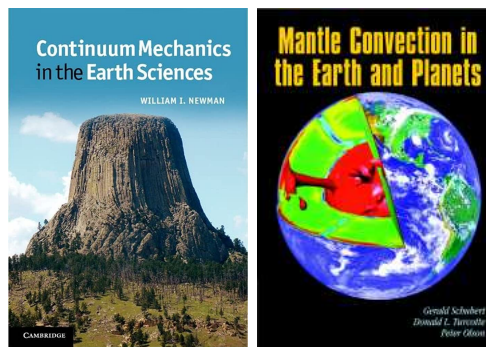


# Chapter 11

## Mantle Dynamics GEO4-1416 syllabus (Utrecht University)

**Remark.** *This handout was written by W. Spakman and is for a large part based on a syllabus by Dr. A.P. van den Berg and Prof. N.J. Vlaar and on material from the book “Mantle convection in the Earth and Planets” by Schubert, Turcotte, and Olson, Cambridge University Press, 2002.*

Additional resources:



Left: William I Newman. *Continuum Mechanics in the Earth Sciences*. Cambridge University Press, 2012. ISBN: 978-0-521-56289-8; Right: G. Schubert, D.L. Turcotte, and P. Olson. *Mantle Convection in the Earth and Planets*. Cambridge University Press, 2001. ISBN: 0-521-70000-0. DOI: 10.1017/CB09780511612879.

In what follows (and in the entire fieldstone document) vectors are denoted by an arrow, e.g.  $\vec{v}$ , while tensors are denoted by a bold font, e.g.  $\boldsymbol{\sigma}$ . It is also my experience that each book/syllabus uses its own set of notations. In particular, and rather to my own confusion, some books call the full stress tensor  $\boldsymbol{\tau}$  and not  $\boldsymbol{\sigma}$ , and call the deviatoric stress  $\boldsymbol{\sigma}$ . Some sources make no difference between the viscous stress tensor and the deviatoric stress tensor. Also it happens (especially in the mathematical literature) that the strain rate tensor is called  $\boldsymbol{D}$  instead of  $\boldsymbol{\epsilon}$  or  $\dot{\boldsymbol{\epsilon}}$ , and that the stress tensor is called  $\boldsymbol{T}$ . All in all be careful when reading additional sources.



|                              |   |                                       |
|------------------------------|---|---------------------------------------|
| Continuum mechanics          | ↔ | continuummechanica                    |
| Coordinate system            | ↔ | coordinatenstelsel                    |
| Cartesian coordinates system | ↔ | cartesische coördinatenstelsel        |
| Scalar                       | ↔ | scalair, scalaire grootheid of scalar |
| Vector                       | ↔ | vector                                |
| Tensor                       | ↔ | tensor                                |
| Axis                         | ↔ | axis                                  |
| Unit vector                  | ↔ | eenheidsvector                        |
| Cross product                | ↔ | vectorprodukt of uitproduct           |
| volume forces                | ↔ | volumekrachten                        |
| Surface forces               | ↔ | oppervlaktekrachten                   |
| Stress tensor                | ↔ | spanningstensor                       |

## 11.0.1 The continuity equation

## 11.1 Review of some essentials of continuum mechanics

Newtonian mechanics deal with particles and rigid (undeformable) bodies on which forces are acting. The application of Newtonian mechanics to realistic media (gases, fluids, solids) is undoable simply because of the many particles (atoms, molecules) involved. Continuum mechanics tackles this problem by assuming that physical fields (e.g. density, temperature, velocity) can be viewed as (piece-wise) continuous functions defined on the time and space coordinates involved in the description of macroscopic matter. The idea is essentially that a tiny cube with sides of, say,  $10^{-8}$  m already contains a sufficient number of atoms (millions) which allows for establishing physically meaningful descriptions of quantities as temperature and density of the cube. In continuum mechanics we are mostly interested in material behavior on much larger scales than  $10^{-8}$  m for which is assumed that physical quantities are smooth functions of time and spatial coordinates.

The forces involved in the deformation of a continuum are postulated to be **body forces**  $\vec{b}$  [ $\text{N m}^{-3}$ ], such that  $\vec{b}dV$  is the force acting on the infinitesimal volume  $dV$ , and surface **tractions**  $\vec{t}^{\vec{n}}$  [ $\text{N m}^{-2}$ ] (e.g. internal friction, applied surface tractions), such that  $\vec{t}^{\vec{n}}dS$  is the force acting on the infinitesimal surface  $dS$ . It is usual to write  $\vec{b} = \rho\vec{g}$ , with  $\rho$  [ $\text{kg m}^{-3}$ ] the mass density and with  $\vec{g}$  [ $\text{ms}^{-2}$ ] the acceleration due to the body force, which in mantle dynamics is gravity.

$\vec{b}$  and  $\vec{t}^{\vec{n}}$  are force densities which after integration over a volume or a surface, respectively, lead to net forces acting on the volume or surface. The traction (or stress vector) is defined as

$$\vec{t}^{\vec{n}} = \lim_{\Delta S \rightarrow 0} \frac{\sum_i f_i^{\vec{n}}}{\Delta S},$$

which expresses the force per unit area working on a tiny surface  $\Delta S$  with unit normal  $\vec{n}$  (defining the orientation of the surface). The forces  $f_i^{\vec{n}}$  can be viewed as the atomic forces [ $\text{N}$ ] that are applied at the  $\vec{n}$ -side of  $\Delta S$  to atoms at the other side of the surface. To maintain equilibrium, by the third law of Newton, the traction applied to the  $-\vec{n}$ -side of the surface is  $\vec{t}^{-\vec{n}} = -\vec{t}^{\vec{n}}$ . Tractions depend on the orientation of the surface. In principle, one can draw an infinite number of oriented surfaces through one point, each associated with a different traction.

Tractions are usually separated into the thermodynamic pressure force  $p\vec{n}$  and the traction  $\vec{\tau}$  related to mechanical deformation:  $\vec{t}^{\vec{n}} = -p\vec{n} + \vec{\tau}$  ( $p > 0$ ). The thermodynamic pressure (a traction always acting perpendicular to any surface) is obtained from the equation of state  $f(\rho, p, T)$  relating thermodynamic quantities density, pressure, and temperature of a continuum. The sign convention in continuum mechanics is that **compression is negative and tension is positive** (in geology this is usually the other way around).

### 11.1.1 Stress

Stress is a second order tensor quantity, which is defined by the following steps:

1. Assume a Cartesian coordinate frame in a point of interest for which the coordinate axis are spanned by three unit orthogonal vectors  $\vec{e}_i$  ( $i = 1, 2, 3$ ),
2. Imagine a tiny cube centered about the origin and with its faces parallel to the coordinate planes,
3. Consider the 3 tractions  $\vec{t}^{\vec{e}_i}$  that are acting on the three positive faces of the cube (i.e. the faces which have the normal vectors  $\vec{e}_i$ ),
4. Lastly, define the components  $\sigma_{ij}$  of the stress tensor  $\boldsymbol{\sigma}$  as

$$\sigma_{ij} = \vec{t}^{\vec{e}_i} \cdot \vec{e}_j \quad (11.1)$$

When the stress tensor is visualized as a matrix, the three rows are the tractions on the positive faces of the cube. The diagonal elements of the stress tensor  $\boldsymbol{\sigma}$  are called normal stresses and the off-diagonal elements are the shear stresses.

Stress is a physical quantity, independent of coordinate frame, but the actual values of components  $\sigma_{ij}$  of the stress tensor can only be computed in a coordinate system. These numbers are dependent on the frame adopted like the components of a flow vector (a first order tensor) are frame dependent. Second order tensors follow (by definition) the coordinate transformation rules of  $3 \times 3$  matrices. From an analysis of force and force-moment balance it can be demonstrated that the stress tensor is symmetric:  $\boldsymbol{\sigma} = \boldsymbol{\sigma}^T$  or  $\sigma_{ij} = \sigma_{ji}$ .

An eigenvalue-eigenvector analysis leads to the three principal stresses  $\sigma_k$  (eigenvalues) and the corresponding three corresponding principal directions  $\vec{q}_k$  (eigenvectors of unit length). The latter span three mutually orthogonal (Cartesian) axes. In the principal-axes frame the stress tensor is diagonal with the three principal stresses as diagonal elements. In the principal-axes frame the tensor components are the maximal normal stresses (tractions perpendicular to the faces of a tiny cube oriented along the principal coordinate planes) compared to the normal stresses in any other coordinate system.

An important relation (the Cauchy relation) exists between the local state of stress (i.e. the stress tensor  $\boldsymbol{\sigma}$ ) and the traction  $\vec{t}^{\vec{n}}$  acting on an (arbitrarily) oriented tiny surface  $\Delta S$  with unit normal  $\vec{n}$ :

$$\boldsymbol{\sigma} \cdot \vec{n} = \vec{t}^{\vec{n}} \quad (11.2)$$

or in components  $\sigma_{ij}n_j = t_i^{\vec{n}}$  (summation convention implied). This relation states how traction can be computed from the local stress and conversely, that from known tractions on independently oriented surfaces the stress tensor can be constructed by solving Eq. (11.2). Note that it is required that  $\vec{n}$  is a unit normal, i.e.  $\vec{n} \cdot \vec{n} = n_j n_j = 1$ .

**Exercise: 1.** Determine the tractions acting on the negative faces of a tiny cube (of which the faces are aligned with the local coordinate axes) when the stress is given.

### 11.1.2 Force balance equation of a continuum at rest

Assume a continuum (gas, liquid, solid) at rest. In this situation the net force acting on the entire continuum is  $\vec{0}$  (Newton). The sum of body forces and applied surface tractions cancel in some way. This also holds for any sub-volume  $V$ . An internal stress field may still exist as a result of the applied forces and surface tractions. The relation between the body force, tractions, and the internal stress

field is derived as follows: Consider an arbitrary sub-volume  $V$  with boundary  $S$ . Internal tractions act on the boundary  $S$  (e.g. to be determined with equation (11.2) from the internal stress field at  $S$ ). The following equation postulates that the total sum of body forces

$$\int_V \rho \vec{g} dV + \int_S \vec{t}^{\vec{n}} dS = \vec{0} \quad (11.3)$$

acting on  $V$  and of tractions on  $S$  leads to a zero net force acting on  $V$ .

Substituting (11.2) in the surface integral and next applying the Divergence theorem<sup>1</sup> one arrives at

$$\int_V \rho \vec{g} dV + \int_V \vec{\nabla} \cdot \boldsymbol{\sigma} dV = \vec{0}$$

Because  $V$  is an arbitrary volume the integrant must equal 0, which leads to the **equilibrium equation** (or momentum conservation equation):

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} = \vec{0} \quad \text{or,} \quad \frac{\partial \sigma_{ij}}{\partial x_j} + \rho g_i = 0 \quad (11.4)$$

This equation holds for any point in the interior of the continuum. Any traction applied at the boundary of the continuum relates to the (local) stress through equation (11.2). Equation (11.4) states that body forces are in equilibrium with the divergence of the stress tensor. Note: Gravity implies spatial variation in stress.

A similar analysis for the equilibrium of torques leads to the symmetry of the stress tensor. In this case the equilibrium equation is

$$\int_V \rho \vec{r} \times \vec{g} dV + \int_S \vec{r} \times \vec{t}^{\vec{n}} dS = \vec{0},$$

where  $\vec{r}$  is the position vector  $\vec{r} = (x_1, x_2, x_3)^T$ . The cross products can be written in index notation using the permutation symbol  $\epsilon_{ijk}$  which equals zero if at least two indices have the same value (e.g.  $\epsilon_{121} = 0$ ), equals +1 if  $ijk$  is an even permutation of 123, and equals -1 if  $ijk$  is an odd permutation of 123. This leads to the following notation of the cross product of two vectors  $\vec{a} \times \vec{b} = \epsilon_{ijk} \vec{e}_i a_j b_k$  and per component  $(\vec{a} \times \vec{b}) = \epsilon_{ijk} a_j b_k$  leading to the torque balance for component  $i$ :

$$\int_V \rho \epsilon_{ijk} x_j g_k dV + \int_S \epsilon_{ijk} x_j t_k^{\vec{n}} dS = 0.$$

**Exercise: 2.** a) Derive equation (11.4)

b) Using a similar approach, prove the symmetry of the stress tensor from the balance of torques (assuming that no internally applied tractions exist)

c) Derive (11.4) by considering the force balance of a tiny cube

### 11.1.3 The material derivative

For a mathematical intro see Appendix W. Consider a continuum with a three-dimensional flow field  $\vec{v}(x_1, x_2, x_3, t)$  dependent on the 3 spatial coordinates  $x_j$  and time  $t$ . For any differentiable scalar function  $T$  defined on these 4 parameters we can write the total differential<sup>2</sup>

$$dT = \frac{\partial T}{\partial t} dt + \sum_j \frac{\partial T}{\partial x_j} dx_j \quad (11.5)$$

<sup>1</sup>[https://en.wikipedia.org/wiki/Divergence\\_theorem](https://en.wikipedia.org/wiki/Divergence_theorem)

<sup>2</sup>See any basic textbook on Calculus.

This equation can be interpreted as follows: Consider a certain point  $(\vec{r}, t)$  in which the scalar function  $T(\vec{r}, t)$  has continuous partial derivatives. The infinitesimal change  $dT$ , which results from going from  $(\vec{r}, t)$  to  $(\vec{r} + d\vec{r}, t + dt)$  in the domain of  $T$  is given by Eq. (11.5). Importantly,  $d\vec{r}$  and  $dt$  can be arbitrarily chosen (including 0 values). The total differential is at the basis of the definition of the so-called directional derivative. Differentiable functions of more than 1 variable can be differentiated in arbitrary directions (in their domain) to find their rate of change in this direction with respect to a specified parameter. Particularly, we can consider the rate of change of  $T$  with the time parameter  $t$  and (spatially) in the direction of the velocity field  $\vec{v}(\vec{r}, t)$ . This time- derivative is easily obtained from (11.5) by coupling  $dt$  and the spatial increment  $d\vec{r}$  such that  $d\vec{r} = \vec{v}dt$ . Substitution in (11.5) leads to the time derivative of  $T$  in the direction of the velocity vector:

$$\frac{DT}{Dt} = \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \frac{\partial T}{\partial t} + \sum_j v_j \frac{\partial T}{\partial x_j} \quad \text{with} \quad \vec{v} = \frac{d\vec{r}}{dt} \quad (11.6)$$

which is called the material derivative of  $T$  giving the rate of change of  $T$  with time in the direction of the flow  $\vec{v}(\vec{r}, t)$  at a certain point  $(\vec{r}, t)$ . The lhs of (11.6) treats  $T$  as a function of  $t$  only in the point  $(\vec{r}(t), t)$  whereas the rhs of (11.6) shows how this can be computed from the partial derivatives of  $T(\vec{r}, t)$  and the local velocity  $\vec{v}(\vec{r}, t)$ . The partial derivative  $\partial T / \partial t$  gives the temporal rate of change at fixed position  $\vec{r}$ , while the second term gives the spatial contribution at fixed time, i.e.  $\vec{v} \cdot \vec{\nabla} T$  expresses the advective contribution (carried with the flow) to  $DT/Dt$ .

Without reference to a particular scalar function the material derivative (operator) is:

$$\frac{D}{Dt} = \frac{\partial}{\partial t} + \vec{v} \cdot \vec{\nabla} \quad (11.7)$$

The material derivative holds for any scalar function, particularly, for the components  $v_i$  of the velocity field leading to the particle acceleration:

$$\frac{D\vec{v}}{Dt} = \frac{\partial \vec{v}}{\partial t} + \vec{v} \cdot \vec{\nabla} \vec{v} \quad (11.8)$$

Note that if the velocity field is time-stationary, i.e.  $\partial_t = 0$ , there is still acceleration. In this case the velocity vector field does not change with time. But, there can still be a spatial variation which gives rise to a stationary acceleration field and material velocity still changes in space, although the velocity is a constant vector in each point.

**Exercise: 3.** Assume 2-D space. Let the temperature field  $T$  be given by

$$T(\vec{r}, t) = T_0 \frac{1}{r} \exp(-t) \quad t > 0, r > 0$$

The temperature field belongs to a flow field given by  $\vec{v}(\vec{r}, t) = \frac{1}{r^2} \exp(-t) \vec{r}$ .

a) Determine the divergence of the velocity field.

b) Compute the particle acceleration.

c) Compute the material derivative of  $T$  at any position and time.

### 11.1.4 The material derivative of a material volume integral

Let  $V$  be a material volume, i.e. a volume that encompasses for all  $t$  the **same flow particles**. This volume is following the flow, possibly being deformed, while there is no material exchange with the region outside  $V$ . Let  $T$  be again a scalar function of the space and time coordinates. The material derivative of the (material) volume integral of  $T$  is (See Appendix W.0.4):

$$\frac{D}{Dt} \left[ \int_{V(t)} T(\vec{r}, t) dV \right] = \int_{V(t)} \frac{DT}{Dt} + T \vec{\nabla} \cdot \vec{v} dV = \int_{V(t)} \frac{DT}{Dt} + T \frac{\partial v_k}{\partial x_k} dV \quad (11.9)$$


**Exercise: 4.** Derive from (11.9) the alternative formula

$$\frac{D}{Dt} \left[ \int_{V(t)} T(\vec{r}, t) dV \right] = \int_{V(t)} \frac{\partial T}{\partial t} + \frac{\partial T v_k}{\partial x_k} dV$$

### 11.1.5 Diffusion processes

Diffusion processes are in many cases described by (empirical) laws of the form  $\vec{a} = -\mathbf{D} \cdot \vec{\nabla} H$  where  $\vec{a}$  is a vector and  $\mathbf{D}$  is the (anisotropic) diffusion (coefficient) tensor, and  $H$  a scalar field. Examples are Fourier's (isotropic) heat flow vector  $\vec{q} = -k \vec{\nabla} T$  where  $k$  is thermal conductivity and  $T$  temperature, or the isotropic diffusion of matter (atoms) given by the mass density flow vector  $\vec{J} = -D \vec{\nabla} c$  where  $D$  is the diffusion coefficient and  $c$  the concentration of the substance.

## 11.2 The basic equations of continuum mechanics

 De algemene bewegingsvergelijking en de continuïteitsvergelijking vormen de basisvergelijkingen van de continuummechanica. Elke deformeerbare stof voldoet eraan.

|                                |   |                                   |
|--------------------------------|---|-----------------------------------|
| The fluid                      | ↔ | de vloeistof                      |
| The continuity equation        | ↔ | de continuïteitsvergelijking      |
| The general equation of motion | ↔ | de algemene bewegingsvergelijking |
| An incompressible material     | ↔ | en onsamendrukbare stof           |
| Strain rate                    | ↔ | de deformatiesnelheid             |
| stress(tensor)                 | ↔ | de spanning(tensor)               |
| velocity gradient              | ↔ |                                   |
| rotation rate                  | ↔ |                                   |
| pressure                       | ↔ | de druk                           |
| The constitutive equation      | ↔ | de constitutievergelijking        |

### 11.2.1 The continuity equation

The continuity equation is also called mass conservation equation or Reynold's transport theorem. Consider an arbitrary material volume  $V$  within a continuum. By definition of a material volume, the mass it contains is conserved,  $M = \text{constant}$ , or  $DM/Dt = 0$ . The mass is given by

$$M = \int_{V(t)} \rho(\vec{r}, t) dV$$

Applying (11.9) to  $DM/Dt = 0$  we find (for arbitrary  $V$ ) the continuity equation as a local expression of mass conservation:

$$\frac{D\rho}{Dt} + \rho \frac{\partial v_k}{\partial x_k} = 0 \quad (11.10)$$

or, when substituting the material derivative:

$$\frac{\partial \rho}{\partial t} + \frac{\partial(\rho v_k)}{\partial x_k} = 0 \quad \left| \quad \frac{\partial \rho}{\partial t} + \vec{\nabla} \cdot (\rho \vec{v}) = 0 \quad (11.11)$$

In incompressible fluids the density cannot change:  $D\rho/Dt = 0$ . Consequently, mass conservation requires that the divergence of the velocity is 0, i.e.

$$\frac{\partial v_k}{\partial x_k} = 0 \quad \left| \quad \vec{\nabla} \cdot \vec{v} = 0 \quad (11.12)$$

Note that an equation like (11.10) can also be derived for any other quantity that is conserved by a material volume.

**Exercise: 5.** Prove that in a flowing medium with density  $\rho$  the following relation holds

$$\frac{D}{Dt} \left[ \int_{V(t)} \rho T dV \right] = \int_{V(t)} \rho \frac{DT}{Dt} dV \quad (11.13)$$

for any differentiable scalar function  $T$  and material volume  $V$ . This formula will be frequently used.

**Exercise: 6.** a) Prove that in an incompressible fluid the cubic-meter content of a material volume does not change (hence, only the shape of boundary of the material volume is allowed to change).

b) Prove that in an incompressible fluid the boundary  $S$  of a material volume obeys the following integral  $\int_S \mathbf{v}_k n_k dS = 0$ . (Interpret this integral)

**Exercise: 7.** Let  $V$  be an imaginary volume fixed in space. Derive the alternative mass conservation law:

$$\int_V \frac{\partial \rho}{\partial t} dV = - \int_S \vec{J} \cdot \vec{n} dS$$

where  $\vec{J} = \rho \vec{v}$  is the mass-density flow. (Interpret this equation).

**Exercise: 8.** We wish to describe the transport of a polluting substance  $X$  carried by a fluid flow. We assume the substance is chemically non-reactive (passive) and dissolved in the fluid. The spatial distribution of  $X$  is given by the concentration function  $c(\vec{r}, t)$  [ $\text{kg m}^{-3}$ ]. Pollutant is also being produced/destroyed according to the function  $H(\vec{r}, t)$  [ $\text{kg s}^{-1} \text{m}^{-3}$ ].

a) Assume for the moment that mass diffusion of  $X$  can be ignored. Derive a conservation law in the form of a differential equation for the concentration function  $c(\vec{r}, t)$ . [Hint: start with computing the total mass of  $X$  (integral form) contained in a material volume  $V$ . Next consider the (material) time derivative of this integral. Is equal to what?]

b) Now assume that mass diffusion of the pollutant is important. This implies material diffusion (not controlled by the fluid flow) across the boundary  $S$  of the control volume  $V$ . Assume that the mass flow density vector  $\vec{J}$  [ $\text{kg s}^{-1} \text{m}^{-2}$ ] is given by  $\vec{J} = -D \vec{\nabla} c$  where  $D$  is the diffusion coefficient. Extend the answer obtained at a) for this situation.

## 11.2.2 The general equation of motion (momentum equation)

The second law of Newton postulates that the sum of applied forces equals the rate of change of the linear momentum of a particle with mass  $m$ :

$$\sum_i \vec{F}_i = \frac{d\vec{p}}{dt}, \quad \vec{p} = m\vec{v}.$$

To arrive at a similar postulate for continuum mechanics, the total linear momentum of an arbitrary material volume is defined as:  $\int_{V(t)} \rho \vec{v} dV$ .

Newton's second law leads to the following postulate of continuum mechanics:

$$\frac{D}{Dt} \int_{V(t)} \rho \vec{v} dV = \int_V \rho \vec{g} dV + \int_S \vec{t} \cdot \vec{n} dS$$



Using (11.13) to evaluate the left side of this equation for each velocity component and applying the derivation following equation (11.3) to the right side we find (as  $V$  is arbitrarily chosen):

$$\rho \frac{D\mathbf{v}_i}{Dt} = \frac{\partial \sigma_{ij}}{\partial j} + \rho g_i \quad \left| \quad \rho \frac{D\vec{\mathbf{v}}}{Dt} = \vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} \right. \quad (11.14)$$

which is the **general equation of motion**. Recall that

$$\frac{D\vec{\mathbf{v}}}{Dt} = \frac{\partial \vec{\mathbf{v}}}{\partial t} + \vec{\mathbf{v}} \cdot \vec{\nabla} \vec{\mathbf{v}}$$

is the material derivative of velocity which renders (11.14) to be a non-linear equation in the unknown velocity field.

### 11.2.3 Velocity gradient, strain rate, and rotation rate

The velocity gradient tensor is  $\vec{\nabla} \vec{\mathbf{v}} = \partial \mathbf{v}_j / \partial x_i$  and can be separated in a symmetric part, the strain rate tensor

$$\dot{\epsilon}_{ij} = \frac{1}{2} \left( \frac{\partial \mathbf{v}_j}{\partial x_i} + \frac{\partial \mathbf{v}_i}{\partial x_j} \right) \quad \left| \quad \dot{\epsilon}(\vec{\mathbf{v}}) = \frac{1}{2} \left( \vec{\nabla} \vec{\mathbf{v}} + (\vec{\nabla} \vec{\mathbf{v}})^T \right)$$

and in an anti-symmetric part called the rotation rate tensor or spin-rate tensor

$$\dot{\omega}(\vec{\mathbf{v}}) = \frac{1}{2} \left( \vec{\nabla} \vec{\mathbf{v}} - (\vec{\nabla} \vec{\mathbf{v}})^T \right)$$

Note that  $\dot{\omega}_{11} = \dot{\omega}_{22} = \dot{\omega}_{33} = 0$ . The strain rate tensor is associated with the rate of deformation (rates of relative length and volume changes and shear) while the spin rate tensor describes an increment of uniform rotation in a continuum (i.e. without internal deformation). Note that  $\vec{\nabla} \cdot \vec{\mathbf{v}} = \partial \mathbf{v}_k / \partial x_k = \dot{\epsilon}_{kk}$  gives the rate of relative volume change during deformation.

**Exercise: 9.** Make a detailed derivation of equation (11.14).

**Exercise: 10.** Assume a velocity field  $\vec{\mathbf{v}} = \vec{\omega} \times \vec{\mathbf{r}}$  with  $\vec{\omega} = [0, 0, \Omega f(x_1, x_2)]^T$  and  $\vec{\mathbf{r}} = (x_1, x_2, 0)^T$  where  $f$  is an unknown function and  $\Omega$  is a constant angular speed (radians/sec). Compute the velocity gradient field, the strain rate field and the rotation rate field. Determine a function  $f$  that leads to incompressible flow and a function  $f$  that leads to flow without shear strain rate.

**Exercise: 11.** Demonstrate that

$$\dot{\omega} = \begin{pmatrix} 0 & -\frac{1}{2}\Omega_3 & \frac{1}{2}\Omega_2 \\ \frac{1}{2}\Omega_3 & 0 & -\frac{1}{2}\Omega_1 \\ -\frac{1}{2}\Omega_2 & \frac{1}{2}\Omega_1 & 0 \end{pmatrix}$$

where

$$\vec{\Omega} = \vec{\nabla} \times \vec{\mathbf{v}}$$

is the so-called vorticity.



### 11.2.4 Pressure and stress

Similar to traction, the total stress  $\boldsymbol{\sigma}(\vec{r}, t)$  is usually separated in the **thermodynamic pressure**  $p$  and the rheological (mechanical) stress  $\boldsymbol{\pi}(\vec{r}, t)$  (also commonly called viscous stress tensor):

$$\sigma_{ij} = -p\delta_{ij} + \pi_{ij} \quad \left| \quad \boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\pi} \right. \quad (11.15)$$

Note:  $\boldsymbol{\pi}$  is *not* the deviatoric stress. Keep reading.

In absence of deforming stresses  $\boldsymbol{\sigma} = -p\mathbf{1}$  and in the equilibrium state (0 inertial force), Eq. (11.14) reduces to the hydrostatic equation

$$0 = -\frac{\partial p}{\partial x_i} + \rho g_i \quad \left| \quad \vec{\nabla} p = \rho \vec{g} \right. \quad (11.16)$$

relating the pressure to the gravitational acceleration.

**Exercise: 12.** Let  $C$  be a line contour in a continuum, which starts at point  $A$  and ends at point  $B$ . The continuum is in hydrostatic equilibrium.

a) Show that the pressure difference between  $B$  and  $A$  equals  $p(B) - p(A) = \int_C \rho \vec{g} \cdot d\vec{r}$  where  $d\vec{r}$  is a line element of  $C$ .

b) Assume that the continuum is a spherically symmetric body: Show that the pressure difference between  $B$  and  $A$  is:  $p(r_B) - p(r_A) = \int_{r_A}^{r_B} \rho g(r) dr$  where quantities only depend on the radius.

c) Show that the density field  $\rho$  should satisfy  $\vec{\nabla} \rho / \vec{\nabla} \Phi$  where  $\Phi$  is the gravitational energy potential implicitly defined by  $\vec{g} = -\vec{\nabla} \Phi$ , i.e. the gradient of  $\rho$  and the gradient of  $\Phi$  are parallel. [Hint: take the curl of the hydrostatic equation]

d) Show that surfaces of constant pressure, density and gravitational potential coincide.

The average pressure is defined as  $\bar{p} = -\frac{1}{3}\sigma_{kk}$ . Using (11.15) we find

$$p - \bar{p} = \frac{1}{3}\pi_{kk} \quad (11.17)$$

which demonstrates that local thermodynamic pressure  $p$  can be perturbed by the isotropic part of mechanical stress. In fluids this situation can occur in locations of local convergence or divergence of flow (see below).

The deviatoric stress  $\boldsymbol{\tau}$  is defined as

$$\tau_{ij} = \sigma'_{ij} = \sigma_{ij} - \frac{1}{3}\sigma_{kk}\delta_{ij} \quad \left| \quad \boldsymbol{\tau} = \boldsymbol{\sigma}' = \boldsymbol{\sigma} - \frac{1}{3}Tr[\boldsymbol{\sigma}]\mathbf{1} \right. \quad (11.18)$$

The deviatoric stress describes the state of stress relative to the ambient (average) pressure. Substituting (11.15) we have

$$\sigma'_{ij} = -(p - \bar{p})\delta_{ij} + \pi_{ij} \quad (11.19)$$

Notice that  $\boldsymbol{\tau} = \boldsymbol{\sigma}' = \boldsymbol{\pi}$  when  $p = \bar{p}$ , i.e. when the mechanical stress does not change the pressure.

**Exercise: 13.** Consider a two-dimensional state of stress.

a) Prove in 2 different ways that the principal values of the deviatoric stress are equal in magnitude but opposite in sign.

b) Demonstrate that only shear traction exists on planes bisecting the principal axis of deviatoric stress

**Exercise: 14.** *Demonstrate that the principal deviatoric stresses  $\sigma'_i$  relate to the principal stresses  $\sigma_i$  as:  $\sigma'_i = \sigma_i + \bar{p}$*

## 11.3 Constitutive equations

The mass conservation equation (11.10) and the momentum equation (11.14) constitute 4 equations in 10 unknowns  $\rho, \vec{v}, \boldsymbol{\sigma}$  (remember: the stress tensor is symmetric so there are only 6 independent terms out of the 9 it contains). Later we will add the energy equation but this also adds the temperature  $T$  as additional unknown. We require knowledge of at least 6 additional independent equations. These are provided for a particular material by specifying the relation between internal kinematics and stress and are called constitutive equations. For real fluids (i.e. fluids that cannot maintain shear stresses) the constitutive relation involves the viscosity as a material parameter. For solids that can deform brittle, elastic, or exhibit fluid behavior, the constitutive equation(s) will in general involve several material parameters. When stress in a solid material exceeds the elastic strength (a stress limit) the material can either break (deform brittle) or enter a regime of so-called ductile behavior where atoms leave their lattice position and occupy new positions elsewhere. Ductile behavior is accommodated by atomic diffusion processes and by dislocation processes (dislocations are geometric disturbances in a crystalline lattice at which fewer atomic bonds exist and which are thus mechanical weakness zones where applied stress will do its work first). Grain boundary processes are also important agents of deformation but are basically determined by diffusion and dislocation processes. Finding relations between stress and strain rate as a function of material properties, temperature, and pressure is the subject of **rheology**.

The thermodynamic pressure  $p$  is taken to be independent of mechanical deformation. The thermodynamic pressure gives the state of stress in a static medium which may have a uniform velocity (either linear, angular or both), i.e.  $\dot{\boldsymbol{\epsilon}} = \mathbf{0}$ . Therefore, constitutive equations take the general form

$$\boldsymbol{\sigma}(\dot{\boldsymbol{\epsilon}}) = -p\mathbf{1} + \boldsymbol{\pi}(\dot{\boldsymbol{\epsilon}}) \quad (11.20)$$

with  $\boldsymbol{\pi}(\mathbf{0}) = \mathbf{0}$ , showing that stress and strain rate can be interdependent.

### 11.3.1 Linear rheology

The most general form of linear rheology leading to a linear viscous fluid (also called a Newtonian fluid) is  $\pi_{ij} = C_{ijkl}\dot{\epsilon}_{kl}$ , where  $C_{ijkl}$  is the anisotropic viscosity tensor. This relation breaks down to the following law for a purely isotropic linearly viscous fluid:

$$\pi_{ij} = \lambda \dot{\epsilon}_{kk} \delta_{ij} + 2\eta \dot{\epsilon}_{ij} \quad \left| \quad \boldsymbol{\pi} = \lambda Tr[\dot{\boldsymbol{\epsilon}}]\mathbf{1} + 2\eta \dot{\boldsymbol{\epsilon}} \quad (11.21)$$

where  $\eta$  is the dynamic viscosity [Pa s] and  $\lambda$  is a viscosity parameter without a specific name<sup>3</sup>.

The mechanical, or rheological, stress  $\pi_{ij}$  quantifies the internal friction of the flow. By computing the trace of the rheological stress we find, using (11.17),

$$p - \bar{p} = \left( \lambda + \frac{2}{3}\eta \right) \dot{\epsilon}_{kk} = \xi \vec{\nabla} \cdot \vec{v} \quad (11.22)$$

where

$$\xi = \lambda + \frac{2}{3}\eta$$

---

<sup>3</sup>It is often called the 'second viscosity coefficient'.

is the **bulk-viscosity**. Buresti [176] (2015) states: “We may then interpret  $\xi \vec{\nabla} \cdot \vec{\nu}$  as the difference between the thermodynamic pressure and the opposite of the average of the normal stresses acting on any three orthogonal planes passing through a point in the fluid, which is usually referred to as the mechanical pressure. This difference is generally considered to be due to the time lag with which the thermodynamic equilibrium condition is reached in a motion that implies an isotropic dilatation of a fluid element.”

Relation (11.22) demonstrates clearly that for a Newtonian fluid the difference between thermodynamic pressure and average pressure is flow induced, i.e. by either local convergence or divergence of flow (i.e.  $\vec{\nabla} \cdot \vec{\nu} \neq 0$ ).

**Remark.** For incompressible fluids  $\vec{\nabla} \cdot \vec{\nu} = 0$ . Then, according to (11.22):  $p = \bar{p}$ . Furthermore, (11.19) leads to  $\sigma'_{ij} = \tau_{ij} = \pi_{ij}$ .

**Remark.** In the case that the bulk viscosity  $\xi$  is assumed 0 (**Stokes condition** - or *Stokes hypothesis*), we have  $p = \bar{p}$  independent of (in)compressibility. Evidently:  $p - \bar{p} \leftrightarrow \xi = 0$  or  $\vec{\nabla} \cdot \vec{\nu} = 0$ . In geodynamics the Stokes hypothesis is assumed so that in practice we never distinguish  $\bar{p}$  from  $p$  and in this case  $\pi = \tau$  so that the equations are formulated as a function of  $p$  and  $\tau$  (also in the compressible case!). Buresti [176] (2015) states: “It implies that the thermodynamic pressure coincides with the mechanical pressure and characterizes the isotropic part of the complete stress tensor; furthermore, the viscous stress tensor becomes a purely deviatoric tensor and corresponds to the deviatoric part of  $\sigma$ . In other words, assuming the validity of this hypothesis is equivalent to state that isotropic dilatations of an elementary volume of fluid do not produce viscous stresses.”

Note that in some books/publications one can find a different formulation of the Stokes hypothesis. For example in Carey and Oden [208] (1986) we find that “Stokes’ hypothesis is that  $\pi$  is a function only of the deformation rate tensor  $\dot{\epsilon}$ ”.

In terms of bulk viscosity and dynamic viscosity Eq. (11.21) reads

$$\pi_{ij} = \left( \xi - \frac{2}{3}\eta \right) \dot{\epsilon}_{kk} \delta_{ij} + 2\eta \dot{\epsilon}_{ij} \quad \left| \quad \pi = \left( \xi - \frac{2}{3}\eta \right) \dot{\epsilon}_{kk} \mathbf{1} + 2\eta \dot{\epsilon}. \quad (11.23)$$

By using the equation of deviatoric strain rate,

$$\dot{\epsilon}'_{ij} = \dot{\epsilon}_{ij} - \frac{1}{3} \dot{\epsilon}_{kk} \delta_{ij} \quad \left| \quad \dot{\epsilon}' = \dot{\epsilon} - \frac{1}{3} \dot{\epsilon}_{kk} \mathbf{1} \quad (11.24)$$

we obtain as alternative expression for Eq. (11.21)

$$\pi_{ij} = \xi \dot{\epsilon}_{kk} \delta_{ij} + 2\eta \dot{\epsilon}'_{ij} \quad \left| \quad \pi = \xi \dot{\epsilon}_{kk} \mathbf{1} + 2\eta \dot{\epsilon}' \quad (11.25)$$

which explicitly shows the role of bulk viscosity and dynamic viscosity. In particular, for an incompressible fluid we have

$$\tau_{ij} = \pi_{ij} = 2\eta \dot{\epsilon}'_{ij} = 2\eta \dot{\epsilon}_{ij} \quad \left| \quad \tau = \pi = 2\eta \dot{\epsilon}' = 2\eta \dot{\epsilon} \quad (11.26)$$

From Eq. (11.25) we can write

$$\sigma = (-p + \xi \vec{\nabla} \cdot \vec{\nu}) \mathbf{1} + 2\eta \dot{\epsilon}'$$

**Exercise: 15.** Prove that for a Newtonian fluid  $\sigma'_{ij} = 2\eta \dot{\epsilon}'_{ij}$  where the prime denotes the deviators of stress and strain rate. (Interpret this general equation)

### 11.3.2 Non-linear rheology

Microphysical processes in solids like atomic diffusion and dislocation motion lead to permanent deformation (macroscopic flow). Pure diffusion, either along grain boundaries or through the bulk of a grain, is a prime example of a Newtonian deformation mechanism. Dislocation glide (low-temperature creep or exponential creep) and diffusion assisted dislocation climb (power law creep) are examples of deformation mechanisms with a nonlinear relation between stress and strain rate. Microphysical considerations (theory) combined with lab experiments (practice) lead to the following (simplified) constitutive equation relating strain rate to rheological stress:

$$\dot{\epsilon}_{ij} = A^{-1} \tau^{n-1} \tau_{ij} \quad (11.27)$$

where  $n \geq 1$  is the stress exponent,  $A$  is a material parameter and  $\tau = (\frac{1}{2} \tau_{ij} \tau_{ij})^{1/2}$  is the effective stress, i.e. the second invariant of rheological stress (i.e. a scalar stress quantity which value is independent of coordinate frame). In analogy with linear viscosity in an incompressible fluid we define the viscosity function

$$\eta = \frac{\tau_{ij}}{2\dot{\epsilon}_{ij}}. \quad (11.28)$$

Combined with (11.27), the viscosity function of power law rheology reads:

$$\eta(\tau) = \frac{1}{2} A \tau^{1-n}. \quad (11.29)$$

Notice that the viscosity decreases as stress increases. The stress (internal friction) is due to body forces and/or applied tractions (cf. Eq. (11.14)). In regions of high internal stress, the viscosity decreases in a power law rheology which causes increasing strain rates (according to  $\tau_{ij} = 2\eta\dot{\epsilon}_{ij}$ ) leading to a localization of deformation.

**Exercise: 16.** *Derive the following relations:*

- a)  $\tau = (A\dot{\epsilon})^{1/n}$
  - b)  $\tau_{ij} = A^{1/n} \dot{\epsilon}^{-1+1/n} \dot{\epsilon}_{ij}$
  - c)  $\eta(\dot{\epsilon}) = \frac{1}{2} A^{1/n} \dot{\epsilon}^{-1+1/n}$
- where  $\dot{\epsilon}$  is the effective strain rate.

In the upper mantle  $n \simeq 3 - 5$  while in the lower mantle  $n \simeq 1 - 3$  (but in both cases not known for sure!). In theory, if  $\dot{\epsilon} \rightarrow 0$  then  $\eta \rightarrow \infty$  which leads to stagnating flow. In practice, however, other (competing) deformation mechanism will provide higher strain rates keeping the viscosity finite. The total strain rate is the sum of strain rate contributions from the different mechanisms active under the same rheological stress:

$$\dot{\epsilon}_{ij} = \sum_k \dot{\epsilon}_{ij}^{(k)} = \frac{1}{2} \sum_k (\eta^{(k)})^{-1} \tau_{ij} = \frac{1}{2} \tau_{ij} \sum_k (\eta^{(k)})^{-1} = \frac{\tau_{ij}}{2\eta_{eff}}$$

where the effective viscosity is

$$\eta_{eff} = \sum_k \frac{1}{\eta^{(k)}} \quad (11.30)$$

A constitutive equation, more detailed than (11.27), encompassing both power law creep and pure diffusion creep is (see for example Karato and Wu [673] (1993)):

$$\dot{\epsilon}_{ij} = A \left( \frac{b}{d} \right)^m \exp \left( -\frac{Q + pV}{RT} \right) \left( \frac{\tau}{\mu} \right)^{n-1} \tau_{ij} \quad (11.31)$$

with

- $b$ : length of Burgers (dislocation) vector ( $\sim 5 \cdot 10^{-10}\text{m}$ )
- $d$ : grain size (0.001-0.1 m)
- $Q$ : activation energy (related to atomic bonding)
- $V$ : activation volume (associated with atomic diffusion)
- $R$ : Universal gas constant 8.31444 J/(mol K)
- $p, T$ : pressure and temperature
- $\mu$ : elastic shear modulus (only used to scale stress) (80 GPa)

The combination of  $n = 1$  and  $m = 2$  or  $3$ , gives a Newtonian fluid resulting from pure atomic diffusion.  $n > 1$  and  $m = 0$  relates to various dislocation (power law) creep mechanisms.

Furthermore, feedback relations can exist between grain growth (due to dynamic crystallization) and stress, e.g.

$$d = Kb \left( \frac{\tau}{\mu} \right)^{-q} \quad K \sim 19 \quad (11.32)$$

The above account of deformation laws (constitutive equations) is by no means exhaustive and only presented to give examples of nonlinear relations between stress and strain rate involving a viscosity function and, in practice, leading to the notion of effective viscosity derived from a superposition of competing deformation mechanisms.

**Exercise: 17.** Discuss the effects of temperature and pressure (as a function of depth) for the constitutive equation (11.31)

## 11.4 The Navier-Stokes equation

Recalling the general equation of motion (11.14)

$$\rho \frac{D\mathbf{v}_i}{Dt} = \rho g_i + \frac{\partial \sigma_{ij}}{\partial x_j} \quad \left| \quad \rho \frac{D\vec{\mathbf{v}}}{Dt} = \rho \vec{g} + \vec{\nabla} \cdot \boldsymbol{\sigma}$$

and the separation of the total stress in the thermodynamic pressure and mechanical stress (11.15)

$$\sigma_{ij} = -p\delta_{ij} + \pi_{ij}, \quad \left| \quad \boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\pi} \quad (11.33)$$

we find after substitution:

$$\rho \frac{D\mathbf{v}_i}{Dt} = \rho g_i - \frac{\partial p}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} \quad \left| \quad \rho \frac{D\vec{\mathbf{v}}}{Dt} = \rho \vec{g} - \vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\pi} \quad (11.34)$$

Adopting the constitutive equation (11.23) for linear rheology

$$\pi_{ij} = \left( \xi - \frac{2}{3}\eta \right) \dot{\epsilon}_{kk} \delta_{ij} + 2\eta \dot{\epsilon}_{ij} \quad \left| \quad \boldsymbol{\pi} = \left( \xi - \frac{2}{3}\eta \right) \dot{\epsilon}_{kk} \mathbf{1} + 2\eta \dot{\boldsymbol{\epsilon}} \quad (11.35)$$

and assuming that viscosities do not depend on the spatial coordinates we find that

$$\rho \frac{D\mathbf{v}_i}{Dt} = \rho g_i - \frac{\partial p}{\partial x_i} + \eta \vec{\nabla}^2 \mathbf{v}_i + \left( \xi + \frac{1}{3}\eta \right) \frac{\partial}{\partial x_i} \vec{\nabla} \cdot \vec{\mathbf{v}} \quad \left| \quad \rho \frac{D\vec{\mathbf{v}}}{Dt} = \rho \vec{g} - \vec{\nabla} p + \eta \vec{\nabla}^2 \vec{\mathbf{v}} + \left( \xi + \frac{1}{3}\eta \right) \vec{\nabla} (\vec{\nabla} \cdot \vec{\mathbf{v}}) \quad (11.36)$$

This is the **Navier-Stokes equation** for a fluid with constant viscosities.

check this equation!

**Exercise: 18.** Derive (11.36)

**Exercise: 19.** Consider a flat layered lithosphere-asthenosphere system with thickness  $L$  and  $h$  respectively. Assume the lithosphere is rigid and moving with a horizontal velocity  $\mathbf{v}_0$ . Further, assume a constant viscosity of the asthenosphere and an incompressible fluid. Take the  $z$ -coordinate positive down with  $z = -L$  corresponding to the top of the lithosphere,  $z = 0$  with the base of the lithosphere and  $z = h$  with the base of the asthenosphere.

- Derive the velocity profile with depth  $z$  by solving the Navier-Stokes equation using a zero material derivative of velocity (an assumption which is valid for mantle convection). The velocity at the base of the asthenosphere should be taken 0 (as a boundary condition).
- The horizontal pressure gradient is still an unconstrained parameter in the solution of a). Assume that the net amount of mass going through a vertical cross section is zero (a mass balance constraint). Use this constraint to determine the horizontal pressure gradient and determine the velocity profile.
- Determine for the model under b) the shear stress at the base of the lithosphere. Determine its value using  $4 \cdot 10^{19}$  Pas for the viscosity, 100 km thickness for the lithosphere and for the asthenosphere and a lithosphere velocity of  $5 \text{ cm yr}^{-1}$ .

## 11.5 Density perturbations as a driving force for mantle convection

The general equation of motion obtained previously

$$\rho \frac{D\mathbf{v}_i}{Dt} = -\frac{\partial p}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} + \rho g_i \quad \left| \quad \rho \frac{D\vec{v}}{Dt} = -\vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\pi} + \rho \vec{g} \right. \quad (11.37)$$

can be rewritten in terms of an equation relative to the hydrostatic reference state of the Earth's mantle. We define the hydrostatic reference state as the undeformed state in which no density perturbations exist. The hydrostatic, or static, pressure is  $p_0(r)$  and the density field in the hydrostatic state is  $\rho_0(r)$ . For the acceleration of gravity in the hydrostatic state we take  $g_i^0(r)$ . The hydrostatic quantities obey the equilibrium equation:

$$0 = -\frac{\partial p_0}{\partial x_i} + \rho_0 g_i^0. \quad \left| \quad \vec{0} = -\vec{\nabla} p_0 + \rho_0 \vec{g}^0 \right. \quad (11.38)$$

We separate quantities in the dynamic state in a hydrostatic contribution and a perturbation:

$$\begin{aligned} p(\vec{r}) &= p_0(\vec{r}) + \tilde{p}(\vec{r}) \\ \rho(\vec{r}) &= \rho_0(\vec{r}) + \tilde{\rho}(\vec{r}) \\ \vec{g}(\vec{r}) &= \vec{g}_0(\vec{r}) + \vec{\tilde{g}}(\vec{r}) \end{aligned}$$

where  $\tilde{p}(\vec{r})$  is the so-called dynamic pressure. Note that the mechanical stress  $\boldsymbol{\pi}$  is by itself a perturbation with respect to the hydrostatic state. After substitution in the equation of motion the hydrostatic terms cancel and we find:

$$(\rho_0 + \tilde{\rho}) \frac{D\mathbf{v}_i}{Dt} = -\frac{\partial \tilde{p}}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} + \tilde{\rho} g_i^0 + \tilde{\rho} \tilde{g}_i + \rho_0 \tilde{g}_i \quad (11.39)$$

The term  $\tilde{\rho} D\mathbf{v}_i/Dt$  and the the last two terms can often be neglected as a small second order perturbation leading to the perturbation equation

$$\rho_0 \frac{D\mathbf{v}_i}{Dt} = -\frac{\partial \tilde{p}}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} + \tilde{\rho} g_i^0 \quad \left| \quad \rho_0 \frac{D\vec{\mathbf{v}}}{Dt} = -\vec{\nabla} \tilde{p} + \vec{\nabla} \cdot \boldsymbol{\pi} + \tilde{\rho} \vec{g} \right. \quad (11.40)$$

This equation shows explicitly that density perturbations and not the total density field are the driving forces for mantle convection.

Gravity inside the Earth relates to the density field according to the differential (Poisson) equation:

$$\vec{\nabla} \cdot \vec{g} = -\vec{\nabla}^2 U = -4\pi \mathcal{G} \rho \quad (11.41)$$

where  $\mathcal{G}$  is the universal constant of gravity. With  $\vec{g} = -\vec{\nabla} U$ , we obtain the **Poisson equation**:

$$\vec{\nabla}^2 U = 4\pi \mathcal{G} \rho \quad (11.42)$$

or, for the dynamic quantities deviating from the hydrostatic state

$$\vec{\nabla}^2 \tilde{U} = 4\pi \mathcal{G} \tilde{\rho} \quad (11.43)$$

**Exercise: 20.** Assume a spherically symmetric non-rotating Earth in hydrostatic equilibrium. In spherical coordinates the divergence of a vector field  $\vec{a}(r)$ , which only depends on the radius  $r$  is

$$\vec{\nabla} \cdot \vec{a} = \frac{1}{r^2} \frac{d}{dr} (r^2 a_r)$$

- a) Prove that the acceleration of gravity at radius  $r$  only depends on the mass contained in the sphere of radius  $R$ .
- b) Assume that the mass of the Earth's core is  $M_c$ . Assume a linear density profile for the crust and mantle and determine the acceleration of gravity as a function of the radius in the mantle.

## 11.6 Two-dimensional formulation for incompressible fluids: the stream function approach

Although 2-D formulations of flow problems may seem restrictive at first glance, it is useful for many applications in which variation in the omitted dimension are small. For instance, slab subduction of a laterally long subduction zone can be modeled in 2-D perpendicular to the strike of a subduction zone. Apart from applications, a full 3-D approach is not always necessary in order to get insight into fundamental problems like the onset of convection, or studies of the influence on the flow pattern of certain model parameters (e.g. viscosity function, internal heat production), or of specific model attributes (e.g. phase transitions, boundary layers). Of course, before the age of high-speed computers, flow calculations were done mostly analytically which is easier in 2-D than in 3-D.

In 2-D we have the 4 unknowns:  $\mathbf{v}_x, \mathbf{v}_y, \rho$  and  $p$ <sup>4</sup>. By assuming incompressibility we can satisfy this equation trivially by defining the stream function  $\psi(x, y)$  as follows:

$$\begin{aligned} \mathbf{v}_x &= \frac{\partial \psi}{\partial y} \\ \mathbf{v}_y &= -\frac{\partial \psi}{\partial x} \end{aligned} \quad (11.44)$$

---

<sup>4</sup>Here too we assume Stokes hypothesis, so that  $\bar{p} = p$  and  $\boldsymbol{\pi} = \boldsymbol{\tau}$

Please check that we have indeed  $\vec{\nabla} \cdot \vec{v} = 0$ . This definition does not impose any restriction on the velocity field or the stream function (apart from being a differentiable function).

Furthermore, assume (for now) that the viscosity is constant and that  $\rho D\vec{v}/Dt \sim 0$  for mantle convection (this will be demonstrated later). Then, the Navier-Stokes equation (11.36) reduces to  $-\vec{\nabla}p + \eta \vec{\nabla}^2 \vec{v} + \rho \vec{g} = \vec{0}$  and then :

$$\begin{aligned} 0 &= -\frac{\partial p}{\partial x} + \eta \left( \frac{\partial^3 \psi}{\partial y^3} + \frac{\partial^3 \psi}{\partial x^2 \partial y} \right) + \rho g_x \\ 0 &= -\frac{\partial p}{\partial y} - \eta \left( \frac{\partial^3 \psi}{\partial x^3} + \frac{\partial^3 \psi}{\partial y^2 \partial x} \right) + \rho g_y \end{aligned} \quad (11.45)$$

The pressure terms in (11.45) can be removed by first differentiating the first line w.r.t.  $y$  and the second line w.r.t.  $x$ ,

$$\begin{aligned} 0 &= -\frac{\partial^2 p}{\partial x \partial y} + \eta \left( \frac{\partial^4 \psi}{\partial y^4} + \frac{\partial^4 \psi}{\partial x^2 \partial y^2} \right) + \frac{\partial(\rho g_x)}{\partial y} \\ 0 &= -\frac{\partial^2 p}{\partial y \partial x} - \eta \left( \frac{\partial^4 \psi}{\partial x^4} + \frac{\partial^4 \psi}{\partial y^2 \partial x^2} \right) + \frac{\partial(\rho g_y)}{\partial x} \end{aligned} \quad (11.46)$$

and next by subtracting the resulting equations, leading to:

$$0 = \eta \left( \frac{\partial^4 \psi}{\partial x^4} + 2 \frac{\partial^4 \psi}{\partial x^2 \partial y^2} + \frac{\partial^4 \psi}{\partial y^4} \right) + \frac{\partial(\rho g_x)}{\partial y} - \frac{\partial(\rho g_y)}{\partial x} \quad (11.47)$$

or

$$\vec{\nabla}^4 \psi = \frac{1}{\eta} \left( -\frac{\partial(\rho g_x)}{\partial y} + \frac{\partial(\rho g_y)}{\partial x} \right) \quad (11.48)$$

where

$$\vec{\nabla}^4 = \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right)$$

Equation (11.48) is the inhomogeneous bi-harmonic equation.

**Exercise: 21.** Show that (11.48) reduces to the homogeneous bi-harmonic equation  $\vec{\nabla}^4 \psi = 0$  if density is independent of the spatial coordinates. NB: we assumed  $\vec{\nabla} \cdot \vec{v} = 0$  which requires, because of the continuity equation, that density is constant in time.

**Exercise: 22.** Consider incompressible flow in three dimensions and (implicitly) define the vector potential  $\vec{\psi}$  as  $\vec{v} = -\vec{\nabla} \times \vec{\psi}$ .

a) Demonstrate that  $\vec{\nabla} \cdot \vec{v} = 0$ .

b) Assume a Newtonian fluid with constant viscosity and assume that  $\rho D\vec{v}_i/Dt \sim 0$ .

Demonstrate that

$$\vec{\nabla} \times \vec{\nabla} \times \vec{\nabla} \times \vec{\nabla} \times \vec{\psi} = -\frac{1}{\eta} \vec{\nabla} \times (\rho \vec{g})$$

(Use the identity  $\vec{\nabla}^2 \vec{v} = \vec{\nabla}(\vec{\nabla} \cdot \vec{v}) - \vec{\nabla} \times \vec{\nabla} \times \vec{v}$  and take the curl of the N-S equation)

## 11.6.1 Application of the stream function approach: Post-glacial rebound

We consider the restoration of the Earth's surface to equilibrium shape in the aftermath of global de-glaciation. Relatively sudden removal of the huge ice-caps, covering parts of the northern and



southern hemisphere, left a depression in the Earth surface which, as a result of horizontal pressure gradients, is being restored to isostatic equilibrium. The uplift resulting from the last glaciation period started some 8000 years ago and is as yet not complete. To model the isostatic rebound we adopt an approximate analytical approach originally due to Haskell [551] (1935) - see also Mitrovica [882] (1996). This analysis will lead to a first estimate of the viscosity of the mantle.

Assume an infinite 2-D half-space filled with an incompressible Newtonian fluid of constant viscosity and with density independent of spatial coordinates. The  $y$ -axis is taken positive downward. The equations to solve are (11.48) for the stream function and next (11.45) for the pressure. As density is constant the right-hand-side of (11.48) is zero. Furthermore, we assume that the horizontal component of gravity is 0 which allows for a straightforward solution of (11.45) when the stream function is known. Because (11.48) is a linear differential equation we assume a spectral approach in which surface deformation is prescribed as a harmonic function with a specific wavelength  $\lambda$  :

$$w_k(x, t) = w_0(t) \cos(kx) \quad (11.49)$$

with the wave number  $k = 2\pi/\lambda$ .

In this model the surface deflection results from sudden ice-unloading at  $t = 0$  and  $w_0(t) \rightarrow 0$  if  $t \rightarrow \infty$ . Because of the linearity of (11.48) the 'loading' function  $w_k(x, t)$  will lead to a stream function of the form  $\psi(x, y) = A(y) \cos kx + B(y) \sin kx$  (i.e. harmonic with separation of variables). A more detailed analysis than is given below will show that  $A(y) = 0$ . For simplicity, we take

$$\psi(x, y) = Y(y) \sin kx \quad (11.50)$$

Substitution in (11.48) leads to

$$\frac{d^4 Y}{dy^4} - 2k^2 \frac{d^2 Y}{dy^2} + k^4 Y = 0 \quad (11.51)$$

This is an ordinary differential equation (ODE) with constant coefficients. Substituting  $Y(y) = Y_0 \exp(my)$  gives

$$m^4 - 2k^2 m^2 + k^4 = (m^2 - k^2)^2 = 0$$

hence  $m = \pm k$  with multiplicity 2. The multiplicity requires two additional solutions  $y \exp(\pm ky)$ . The general solution is then

$$Y(y) = A \exp(-ky) + By \exp(-ky) + C \exp(ky) + Dy \exp(ky) \quad (11.52)$$

To arrive at a specified solution for the problem we have to consider boundary conditions. The first condition is that when  $y \rightarrow \infty$  then  $\psi \rightarrow 0$ . This guarantees a finite solution at infinity and satisfies the idea that at large depth the velocity field should approach zero. Substitution in (11.52) leads to  $C = D = 0$  and to the stream function solution

$$\psi(x, y) = (A + By) \exp(-ky) \sin kx \quad (11.53)$$

The second boundary condition concerns the motion of the surface. We expect that the deflection  $w_k(x, t)$  primarily leads to a vertical motion for small wave number  $k$  (large wavelength  $\lambda$ ) and assume that horizontal motion is negligibly small<sup>5</sup>. The condition  $\mathbf{v}_x(x, y) = 0$  for  $y = w_k(x, t)$  is only a crude approximation. Still, for Haskell's problem it helps finding an analytical solution. Using (11.44) we find for the velocity components:

$$\mathbf{v}_x = (B - k(A + By)) \exp(-ky) \sin kx \quad (11.54)$$

$$\mathbf{v}_y = -k(A + By) \exp(-ky) \cos kx \quad (11.55)$$

---

<sup>5</sup>From modern GPS research on surface motions in Scandinavia we know that horizontal surface motions are between 10-20% of the vertical motion

**Exercise: 23.** Derive equations (11.54) and (11.55).

The condition  $\mathbf{v}_x(x, w_k(x, t)) = 0$  leads to  $B \simeq kA$  where we used that  $kw_k = 2\pi w_k/\lambda \ll 1$  because  $w_k$  is on the order of 1 km while  $\lambda$  is on the order of 3000 km. This leads to the following results:

$$\psi(x, y) = A(1 + ky) \exp(-ky) \sin kx \quad (11.56)$$

$$\mathbf{v}_x = -Ak^2 y \exp(-ky) \sin kx \quad (11.57)$$

$$\mathbf{v}_y = -Ak(1 + ky) \exp(-ky) \cos kx \quad (11.58)$$

The last boundary condition concerns the zero pressure at the deformed surface  $y = w_k(x, t)$ . To solve for the pressure the solution (11.56) is first substituted in (11.45) which gives

$$0 = -\frac{\partial p}{\partial x} + 2\eta Ak^3 \exp(-ky) \sin kx + \rho g_x \quad (11.59)$$

$$0 = -\frac{\partial p}{\partial y} + 2\eta Ak^3 \exp(-ky) \cos kx + \rho g_y \quad (11.60)$$

**Exercise: 24.** Derive equations (11.59) and (11.60).

With  $g_x \simeq 0$ , integration of (11.60) leads to

$$p(x, y) = -2\eta Ak^2 \exp(-ky) \cos kx + \rho g_y y + p_0 \quad (11.61)$$

**Exercise: 25.** Derive Eq. (11.61) and determine that  $p_0 = 0$ .

Using that  $kw_k \ll 1$ , the condition  $p(x, w_k(x, t)) = 0$  gives for the surface deflection

$$w_k(x, t) = \frac{2A\eta k^2}{\rho g_y} \cos kx \quad (11.62)$$

This solution explicitly relates the amplitude of the surface deflection to the wave length  $\lambda = 2\pi/k$  and the principal model parameters density and viscosity. However, there is still an undetermined coefficient  $A$  which must describe the time behavior of the surface deformation (recall that  $w_k(x, t) \rightarrow 0$  if  $t \rightarrow \infty$ ). Time has played no role in the derivation because we have basically solved a time-stationary process. Time stationary problems are characterized by a static velocity field (time-constant velocity vectors at each position). But, the velocity field still describes material flow and time is a parameter if we wish to follow a particle in the flow. Particularly, the surface particles move with the vertical velocity  $\mathbf{v}_y(x, w_k(x, t))$  and thus

$$\frac{\partial w_k}{\partial t} = \mathbf{v}_y(x, w_k(x, t)) = -Ak(1 + kw_k) \exp(-kw_k) \cos kx \simeq -Ak \cos kx \quad (11.63)$$

Using (11.62) we can eliminate  $A$  and arrive at the differential equation

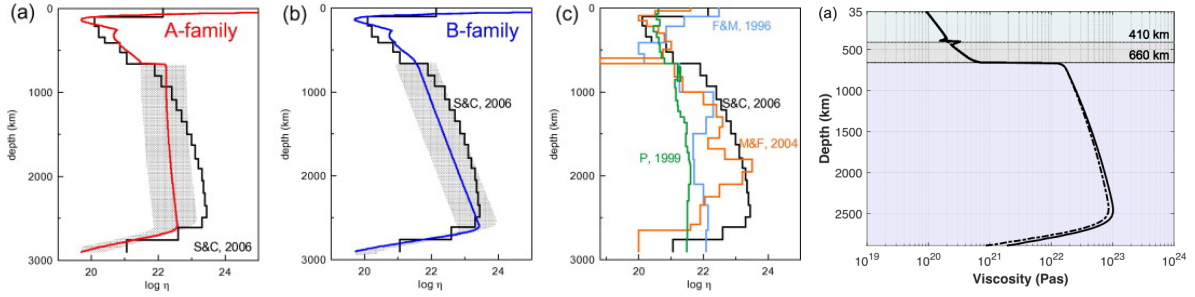
$$\frac{\partial w_k}{\partial t} = -\frac{\rho g_y}{2\eta k} w_k \quad (11.64)$$

which has the solution

$$w_k(x, t) = w_k(x, 0) \exp(-t/t_r) \quad (11.65)$$

where the relaxation time is  $t_r = 4\pi\eta/\rho g_y \lambda$ .

From geological observations of (e.g. beach uplift, river incisions), data curves of  $w(x, t)$  have been obtained which leads to relaxation times of about 4400 yr for wave lengths of about 3000 km. Using values of  $3300 \text{ kg m}^{-3}$  for density and  $10 \text{ m s}^{-2}$  for the acceleration of gravity, we find for the viscosity  $\eta \simeq 10^{21} \text{ Pa s}$ , a number that still stands today as an average value for the 600 to 1000 km of the mantle.



Left: Taken from Ciskova *et al.* [259] (2012); Right: Taken from Neuharth and Mittelstaedt [934] (2023).

Modern modeling concentrates on the complex global problem of Global Isostatic Adjustment (GIA) which also involves the global variation of sea-level which is coupled to ice-sheet formation and melting. Both are functions of time and surface coordinates (topography!). Positive feedback exist between water extraction from the oceans (unloading; attraction of mantle flow) and ice-sheet creation elsewhere (loading; pushing the mantle away). Inverse processes occur during and after ice cap melting. Much data is available on relative sea-level changes, but much less about actual ice sheet (un-)loading histories. Results obtained from different GIA modeling strategies still show disagreement on the detail of viscosity change with depth which also results from the relatively poor sensitivity of surface motions for the detail of viscosity layering. The interested student can find its way in the literature through papers of J. Mitrovica, D. Peltier, and K. Lambeck.

## 11.7 The energy equation

Mantle convection is driven by density perturbations relative to a hydrostatic state. In part, the density perturbations  $\Delta\rho$  are due to thermal perturbations  $\Delta T$  where the connection is given by an equation of state involving the thermal expansion coefficient  $\alpha$ . The full description of mantle convection requires also an equation describing the temperature field of flow. This equation is called the energy equation and involves contributions from adiabatic (de)compression, heat dissipation due to friction in the flow (viscous dissipation), heat conduction and advection, and heat production (including phase changes).

The derivation of the energy equation starts with the first law of thermodynamics which equates the change  $\Delta E$  in total energy of a system to the work  $\Delta W$  done by thermo-mechanical processes and the heat input  $\Delta Q$  to the system resulting from heat flow and heat production (we disregard here the energy contribution from chemical and electro- magnetic processes):

$$\Delta E = \Delta W + \Delta Q \quad (11.66)$$

For processes developing continuously in time we can write (11.66) in terms of the power

$$\frac{DE}{Dt} = \frac{DW}{Dt} + \frac{DQ}{Dt} \quad (11.67)$$

To apply the first law of thermodynamics to the mechanical deformation of a continuum, consider a material volume  $V$  with boundary  $S$ . The power<sup>6</sup> developed by mechanical forces  $\rho\vec{g}$  and the boundary tractions  $\vec{t}^n$  is

$$\frac{DW}{Dt} = \int_V \rho g_i v_i dV + \int_S v_i t_i^n dS \quad \left| \quad \frac{DW}{Dt} = \int_V \rho \vec{g} \cdot \vec{v} dV + \int_S \vec{v} \cdot \vec{t}^n dS \quad (11.68)$$

By substituting the Cauchy relation (11.2) (i.e.  $t_i^n = \sigma_{ij} n_j$  or  $\vec{t}^n = \boldsymbol{\sigma} \cdot \vec{n}$ ), with  $\vec{n}$  the outward pointing normal on  $S$ , and next applying the divergence theorem, the surface integral is transformed to a volume integral:

$$\begin{aligned} \frac{DW}{Dt} &= \int_V \rho \vec{v} \cdot \vec{g} dV + \int_S \vec{v} \cdot (\boldsymbol{\sigma} \cdot \vec{n}) dS \\ &= \int_V \rho \vec{v} \cdot \vec{g} dV + \int_S (\vec{v} \cdot \boldsymbol{\sigma}) \cdot \vec{n} dS \\ &= \int_V \rho \vec{v} \cdot \vec{g} dV + \int_V \vec{\nabla} \cdot (\vec{v} \cdot \boldsymbol{\sigma}) dV \\ &= \int_V \vec{v} \cdot (\rho \vec{g}) dV + \int_V \vec{v} \cdot \vec{\nabla} \cdot \boldsymbol{\sigma} dV + \int_V \vec{\nabla} \vec{v} : \boldsymbol{\sigma} dV \end{aligned}$$

After rearranging terms we get

$$\frac{DW}{Dt} = \int_V \left[ v_i \left( \rho g_i + \frac{\partial \sigma_{ij}}{\partial x_j} \right) + \sigma_{ij} \frac{\partial v_i}{\partial x_j} \right] dV \quad \left| \quad \frac{DW}{Dt} = \int_V \left[ \vec{v} \cdot (\rho \vec{g} + \vec{\nabla} \cdot \boldsymbol{\sigma}) + \boldsymbol{\sigma} : \vec{\nabla} \vec{v} \right] dV \quad (11.69)$$

The general equation of motion (11.14) is used in (11.69) to arrive at the second equality. Next, we

---

<sup>6</sup>[https://en.wikipedia.org/wiki/Power\\_\(physics\)](https://en.wikipedia.org/wiki/Power_(physics))

can apply (11.13) to the first term in the last equality which leads to

$$\begin{aligned}
\frac{DW}{Dt} &= \int_V \left[ \vec{v} \cdot (\rho \vec{g} + \vec{\nabla} \cdot \boldsymbol{\sigma}) + \boldsymbol{\sigma} : \vec{\nabla} \vec{v} \right] dV \\
&= \int_V \left[ \vec{v} \cdot \left( \rho \frac{D\vec{v}}{Dt} \right) + \boldsymbol{\sigma} : \vec{\nabla} \vec{v} \right] dV \\
&= \int_V \left[ \rho \frac{D(\frac{1}{2} \vec{v} \cdot \vec{v})}{Dt} + \boldsymbol{\sigma} : \vec{\nabla} \vec{v} \right] dV \\
&= \frac{D}{Dt} \int_V \rho \frac{1}{2} \vec{v} \cdot \vec{v} dV + \int_V \boldsymbol{\sigma} : \vec{\nabla} \vec{v} dV
\end{aligned} \tag{11.70}$$

or,

$$\begin{aligned}
\frac{DW}{Dt} &= \int_V \left[ v_i \left( \rho g_i + \frac{\partial \sigma_{ij}}{\partial x_j} \right) + \sigma_{ij} \frac{\partial v_i}{\partial x_j} \right] dV \\
&= \int_V \left( \rho v_i \frac{Dv_i}{Dt} + \sigma_{ij} \frac{\partial v_i}{\partial x_j} \right) dV \\
&= \int_V \left( \frac{1}{2} \rho \frac{D(v_i v_i)}{Dt} + \sigma_{ij} \frac{\partial v_i}{\partial x_j} \right) dV \\
&= \frac{D}{Dt} \int_V \frac{1}{2} \rho v_i v_i dV + \int_V \sigma_{ij} \frac{\partial v_i}{\partial x_j} dV
\end{aligned} \tag{11.71}$$

The first term in (11.71) describes the power input from kinetic energy of the flow and the second, the power input resulting from viscous dissipation.

The term  $DQ/Dt$  in (11.67) concerns the heat flow  $\vec{q}(\vec{r}, t)$  across the boundary  $S$  and the heat production  $H(\vec{r}, t)$  within the volume  $V$ . The heat flow  $\vec{q}$  has dimensions  $\text{J s}^{-1} \text{m}^{-2}$  and caused by thermal gradients in the medium as described by the Fourier law  $\vec{q} = -k \vec{\nabla} T$  where  $k$  is the thermal conductivity. The heat production  $H$  (or consumption as in an endothermic phase change) has units  $\text{J s}^{-1} \text{kg}^{-1}$ . The total rate of heat input is obtained by integration of the heat production over the volume and by integration of the normal component of heat flow over the surface  $S$

$$\frac{DQ}{Dt} = \int_V \rho H dV - \int_S \vec{q} \cdot \vec{n} dS \tag{11.72}$$

$$= \int_V \rho H dV - \int_V \vec{\nabla} \cdot \vec{q} dV \tag{11.73}$$

$$= \int_V \left( \rho H + \vec{\nabla} \cdot (k \vec{\nabla} T) \right) dV \tag{11.74}$$

Outward directed heat flow causes a negative contribution to  $DQ/Dt$  as heat is flowing out of the volume. Because the normal  $\vec{n}$  is defined as outward pointing on  $S$  a minus sign is therefore required in front of the heat flow integral in the first equality. In the second equality, the Fourier law for heat conduction is substituted. Application of the divergence theorem leads to the last equality.

The total energy  $E$  in (11.67) is now separated into two contributions: the kinetic energy and the internal energy<sup>7</sup>  $e$  ( $\text{J kg}^{-1}$ )<sup>8</sup>:

$$E = \int_V \left( \frac{1}{2} \rho v_i v_i + \rho e \right) dV \quad \Bigg| \quad E = \int_V \left( \frac{1}{2} \rho \vec{v} \cdot \vec{v} + \rho e \right) dV \tag{11.75}$$

<sup>7</sup>[https://en.wikipedia.org/wiki/Internal\\_energy](https://en.wikipedia.org/wiki/Internal_energy)

<sup>8</sup> $[e] = L^2 T^{-2}$

Taking the time derivative leads to

$$\begin{aligned}\frac{DE}{Dt} &= \frac{D}{Dt} \left( \int_V \frac{1}{2} \rho \vec{v} \cdot \vec{v} dV \right) + \frac{D}{Dt} \left( \int_V \rho e dV \right) \\ &= \frac{D}{Dt} \left( \int_V \frac{1}{2} \rho \vec{v} \cdot \vec{v} dV \right) + \left( \int_V \rho \frac{De}{Dt} dV \right)\end{aligned}\quad (11.76)$$

where we applied equation (11.13) to rewrite the last integral.

We are now ready to create the energy balance of Eq. (11.67): Combining equations (11.76), (11.71), and (11.74) results in an integral equation over the material volume  $V$  in which the kinetic energy terms cancel. As  $V$  is chosen arbitrary we find the **energy (or heat) equation**

$$\rho \frac{De}{Dt} = \sigma_{ij} \frac{\partial v_i}{\partial x_j} + \frac{\partial}{\partial x_j} \left( k \frac{\partial T}{\partial x_j} \right) + \rho H \quad \left| \quad \rho \frac{De}{Dt} = \boldsymbol{\sigma} : \vec{\nabla} \vec{v} + \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H \right. \quad (11.77)$$

which states that the rate of change of internal energy equals the sum of viscous heat dissipation, thermal conduction, and heat production.

**Remark.** *Looking at units:* We have the internal energy  $e$  given in  $J/kg$ , i.e.  $[e] = ML^2T^{-2}/M = L^2T^{-2}$ . Then  $[\rho De/Dt] = ML^{-3}(L^2T^{-2})T^{-1} = ML^{-1}T^{-3}$ ,  $[\boldsymbol{\sigma} : \vec{\nabla} \vec{v}] = ML^{-1}T^{-2}L^{-1}LT^{-1} = ML^{-1}T^{-3}$ ,  $[\vec{\nabla} \cdot (k \vec{\nabla} T)] = L^{-1}LMT^{-3}\theta^{-1}L^{-1}\theta = ML^{-1}T^{-3}$  and  $[\rho H] = ML^{-1}T^{-3}$  so that  $[H] = L^2T^{-3}$ .

Taking the usual separation of the stress tensor in the thermodynamic stress and the mechanical stress, i.e  $\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\pi}$ , we get

$$\rho \frac{De}{Dt} + p \frac{\partial v_j}{\partial x_j} = \pi_{ij} \frac{\partial v_i}{\partial x_j} + \frac{\partial}{\partial x_j} \left( k \frac{\partial T}{\partial x_j} \right) + \rho H \quad (11.78)$$

or

$$\rho \frac{De}{Dt} + p \vec{\nabla} \cdot \vec{v} = \boldsymbol{\pi} : \vec{\nabla} \vec{v} + \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H$$

Note that under the Stokes hypothesis (see before), we have

$$\rho \frac{De}{Dt} + p \vec{\nabla} \cdot \vec{v} = \boldsymbol{\tau} : \vec{\nabla} \vec{v} + \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H$$

From thermodynamic considerations<sup>9</sup> (not treated here) we have the relation  $de = TdS - pdv$  between internal energy  $e$ , the entropy  $S$ , and the specific volume  $v = 1/\rho$ . This relation is used to transform Eq. (11.78) into the entropy form of the energy equation:

$$\frac{De}{Dt} = T \frac{DS}{Dt} - p \frac{D(1/\rho)}{Dt} = T \frac{DS}{Dt} + \frac{1}{\rho^2} p \frac{D\rho}{Dt} = T \frac{DS}{Dt} + \frac{1}{\rho^2} p (-\rho \vec{\nabla} \cdot \vec{v}) \quad (11.79)$$

where we have used the continuity equation. We then obtain

$$\rho \frac{De}{Dt} = \rho T \frac{DS}{Dt} - p \vec{\nabla} \cdot \vec{v}$$

so that in the end

$$\rho T \frac{DS}{Dt} = \boldsymbol{\pi} : \vec{\nabla} \vec{v} + \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H \quad (11.80)$$

An **adiabatic** state is a state of reversible processes with no heat exchange with surroundings and is defined by  $S=\text{constant}$  in which case the lhs of (11.80) is 0. This state is not compatible with

<sup>9</sup>See for instance [https://en.wikipedia.org/wiki/Maxwell\\_relations](https://en.wikipedia.org/wiki/Maxwell_relations)

viscous heat dissipation, heat conduction, or internal heat production, as each of these processes would lead to a non-zero contribution at the rhs of (11.80). A fluid particle flowing in an adiabatic mantle assumes at any time the temperature and pressure of the ambient mantle. Therefore there is no conductive heat exchange in the adiabatic state.

Another useful thermodynamic relation is<sup>10</sup>

$$\frac{DS}{Dt} = \left( \frac{\partial S}{\partial T} \right)_p \frac{DT}{Dt} + \left( \frac{\partial S}{\partial p} \right)_T \frac{Dp}{Dt}$$

For the first term we use the definition of the heat capacity ( $\text{J K}^{-1}$ ) at constant pressure

$$C_p = T \left( \frac{\partial S}{\partial T} \right)_p$$

while for the second one we use

$$\left( \frac{\partial S}{\partial p} \right)_T = - \left( \frac{\partial(1/\rho)}{\partial T} \right)_p = \frac{1}{\rho^2} \left( \frac{\partial \rho}{\partial T} \right)_p$$

and the definition of the thermal expansion coefficient ( $\text{K}^{-1}$ )

$$\alpha = -\frac{1}{\rho} \left( \frac{\partial \rho}{\partial T} \right)_p$$

so that in the end we have the thermodynamical relationship

$$\begin{aligned} \frac{DS}{Dt} &= \frac{C_p}{T} \frac{DT}{Dt} - \left( \frac{\partial(1/\rho)}{\partial T} \right)_p \frac{Dp}{Dt} \\ &= \frac{C_p}{T} \frac{DT}{Dt} + \frac{1}{\rho^2} \left( \frac{\partial \rho}{\partial T} \right)_p \frac{Dp}{Dt} \\ &= \frac{C_p}{T} \frac{DT}{Dt} - \frac{\alpha}{\rho} \frac{Dp}{Dt} \end{aligned}$$

**Remark.** Let us look again at the dimensions of these quantities. The entropy  $S$  is in  $\text{J K}^{-1}$ , or  $[S] = ML^2T^{-2}\theta^{-1}$ , and we know that  $C_p$  has the same unit. We have  $[\alpha] = \theta^{-1}$ ,  $[\rho] = ML^{-3}$  and  $[dp] = ML^{-1}T^{-2}$  so we indeed recover  $[\alpha dp/\rho] = \theta^{-1}ML^{-1}T^{-2}M^{-1}L^3 =$ . All in well.

We have

$$\frac{DS}{Dt} = \frac{C_p}{T} \frac{DT}{Dt} - \frac{\alpha}{\rho} \frac{Dp}{Dt}$$

Applying this relation to (11.80) leads to the temperature form of the heat equation

$$\rho C_p \frac{DT}{Dt} - \alpha T \frac{Dp}{Dt} = \boldsymbol{\tau} : \vec{\nabla} \vec{v} + \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H \quad (11.81)$$

The lhs of this equation is zero for an adiabatic state. In this case we find that

$$\frac{dT_a}{dp} = \frac{\alpha T_a}{\rho C_p} \quad (11.82)$$

which gives how temperature changes due to pure adiabatic (de)compression.

In the end we solve the energy equation in this form:

---

<sup>10</sup>[https://en.wikipedia.org/wiki/Relations\\_between\\_heat\\_capacities](https://en.wikipedia.org/wiki/Relations_between_heat_capacities)

$$\rho C_p \frac{DT}{Dt} - \vec{\nabla} \cdot k \vec{\nabla} T = \alpha T \frac{Dp}{Dt} + \Phi + \rho H \quad (11.83)$$

where  $\Phi$  is the shear heating. Note that it couples temperature, velocity (and its derivative the strain rate) and pressure.

In many publications the shear hearing is denoted by  $\Phi$  and is given by  $\Phi = \tau_{ij} \partial_j \mathbf{v}_i = \boldsymbol{\tau} : \vec{\nabla} \vec{\mathbf{v}}$  where  $\boldsymbol{\tau}$  is the deviatoric stress tensor. In what follows I use the index notation as it makes for easier derivations:

$$\begin{aligned} \Phi &= \tau_{ij} \partial_j \mathbf{v}_i \\ &= 2\eta \dot{\epsilon}_{ij}^d \partial_j \mathbf{v}_i \\ &= 2\eta \frac{1}{2} (\dot{\epsilon}_{ij}^d \partial_j \mathbf{v}_i + \dot{\epsilon}_{ji}^d \partial_i \mathbf{v}_j) \\ &= 2\eta \frac{1}{2} (\dot{\epsilon}_{ij}^d \partial_j \mathbf{v}_i + \dot{\epsilon}_{ij}^d \partial_i \mathbf{v}_j) \\ &= 2\eta \dot{\epsilon}_{ij}^d \frac{1}{2} (\partial_j \mathbf{v}_i + \partial_i \mathbf{v}_j) \\ &= 2\eta \dot{\epsilon}_{ij}^d \dot{\epsilon}_{ij} \\ &= 2\eta \dot{\boldsymbol{\epsilon}}^d : \dot{\boldsymbol{\epsilon}} \\ &= 2\eta \dot{\boldsymbol{\epsilon}}^d : \left( \dot{\boldsymbol{\epsilon}}^d + \frac{1}{3} (\vec{\nabla} \cdot \vec{\mathbf{v}}) \mathbf{1} \right) \\ &= 2\eta \dot{\boldsymbol{\epsilon}}^d : \dot{\boldsymbol{\epsilon}}^d + 2\eta \dot{\boldsymbol{\epsilon}}^d : \mathbf{1} (\vec{\nabla} \cdot \vec{\mathbf{v}}) \\ &= 2\eta \dot{\boldsymbol{\epsilon}}^d : \dot{\boldsymbol{\epsilon}}^d \end{aligned} \quad (11.84)$$

Finally (in Cartesian coordinates)

$$\Phi = \boldsymbol{\tau} : \vec{\nabla} \vec{\mathbf{v}} = 2\eta \dot{\boldsymbol{\epsilon}}^d : \dot{\boldsymbol{\epsilon}}^d = 2\eta ((\dot{\epsilon}_{xx}^d)^2 + (\dot{\epsilon}_{yy}^d)^2 + 2(\dot{\epsilon}_{xy}^d)^2) \quad (11.85)$$

See Schubert & Yuen (1978) [1142] for an analysis of shear heating instability in the upper mantle.

Let us quickly look at the  $Dp/Dt = \partial_t p + \vec{\mathbf{v}} \cdot \vec{\nabla} p$  term. Often the term  $\partial p / \partial t$  is neglected and the pressure is assumed to be mostly hydrostatic in this term so that  $\vec{\nabla} p = -\rho \vec{g}$  which yields a much simpler formulation.

See discussions about shear heating (“viscous dissipation”) – meaning and application– in Froidevaux [420] (1973), Stein [1198] (1978), Bird and Yuen [90] (1979), Sleep, Stein, Geller, and Gordon [1173] (1979), Winter [1365] (1987), Masek and Duncan [838] (1998).

A mantle in the adiabatic state satisfies the equilibrium equation  $\vec{\nabla} p = \rho \vec{g}$ . Assuming spherical symmetry we have  $dp/dr = -\rho g$ . Substitution in (11.82) gives

$$\frac{dT_a}{dr} = -\frac{\alpha T_a g}{C_p}, \quad (11.86)$$

i.e. the **adiabatic temperature gradient**.

Along similar lines, for a mantle in motion we can approximate the second term on the lhs of (11.81) (the adiabatic (de)compression term) as

$$\alpha T \frac{dp}{dt} \simeq -\alpha T \frac{\rho \vec{g} \cdot d\vec{r}}{dt} = -\alpha \rho T \vec{g} \cdot \vec{\mathbf{v}} \simeq -\alpha \rho T g_r \mathbf{v}_r \quad (11.87)$$

where  $g_r$  and  $\mathbf{v}_r$  are the radial components of gravity and velocity, respectively.



**Exercise: 26.** Determine the temperature as a function of depth for a particle that is being transported through the center of a vertical mantle upwelling. Assume adiabatic conditions and a constant vertical flow velocity.

**Exercise: 27.** Show that for a Newtonian fluid the dissipation term  $\Phi = \pi_{ij} \partial v_i / \partial x_j$  can be written as  $\Phi = \pi_{ij} \dot{\epsilon}_{ij} = 2\eta \dot{\epsilon}'_{ij} \dot{\epsilon}'_{ij} + \xi (\dot{\epsilon}_{kk})^2$ . What can one conclude for the two viscosities?

**Exercise: 28.** Demonstrate that for 2D laminar flow (as in exercise 19) the dissipation function can be written as  $\Phi = \eta (\partial v_x / \partial z)^2$ . Use this to evaluate the dissipative heat production in the flow of exercise 19a assuming zero horizontal pressure gradient (Couette flow). Calculate a numerical value of this heat production using values of  $h = 200$  km,  $\eta = 10^{21}$  Pa s,  $v_0 = 1$  cm yr<sup>-1</sup>. Compare this to estimates of radiogenic heat production in the upper mantle of  $8.4 \cdot 10^{-9}$  J s<sup>-1</sup> m<sup>-3</sup>.

**Exercise: 29.** The Navier-Stokes equation for perturbations relative to a hydrostatic reference state is in Cartesian coordinates

$$\rho_0 \frac{Dv_i}{Dt} = -\frac{\partial \tilde{p}}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} + \tilde{\rho} g_i^0$$

(see equation 11.40). The density perturbation is driving the flow of a medium contained in the material volume  $V$ . Assume that gravity is only working vertical  $\vec{g} = (0, 0, g)^T$  and that the velocity field is  $\vec{v} = (u, v, w)^T$ . Assume incompressible flow and that the inertial term can be neglected.

a) Derive the following energy conservation law relating the dissipation of gravitational energy into the energy released due to frictional flow:

$$\int_V \tilde{\rho} g w \, dV = \int_V \pi_{ij} \frac{\partial v_i}{\partial x_j} \, dV$$

Hint: Take the inner product of the equation of motion with the velocity field and integrate the result over  $V$ . Assume an impermeable boundary  $S$  of  $V$  (i.e.  $\vec{v} \cdot \vec{n} = 0$ ) and that the boundary is shear stress free (free slip) or has no-slip ( $\vec{v} = \vec{0}$ ).

b) Show by substitution of the Newtonian rheology that the rhs of this equation is positive and show that the lhs of the equation equals the created power due to the change in gravitational potential energy as a result of  $\tilde{\rho}$  w.r.t.  $\rho_0$ .

**Exercise: 30.** a) Derive the heat equation for a static medium ( $\vec{v} = \vec{0}$ ) from the conservation law of thermal energy. (Hint: Create the heat balance for a fixed control volume based on internal energy, heat production and heat flow through the boundary of the volume).

b) Next, assume there is a flow field  $\vec{v}$  and extend the result under a) with a term involving the flow density of thermal energy  $\vec{J} = \rho C_p T \vec{v}$  (assume  $\rho C_p$  to be constant and  $\vec{\nabla} \cdot \vec{v} = 0$ ).

The result is the heat equation for an incompressible medium in which adiabatic compression and viscous dissipation are neglected.

## 11.8 The equation of state

An equation of state relates basic thermodynamic parameters such as density, temperature, pressure, or entropy. For an isochemical fluid, we can choose two independent thermodynamic parameters on which all other depend. Here we choose temperature and pressure as independent parameters, e.g.  $\rho(T, p)$ . Mantle convection is usually studied relative to some motionless hydrostatic reference state either with a conductive or with an adiabatic geotherm. The perturbations with respect to this state are related to the convective state of the mantle. Therefore, we write  $\rho(T, p) = \rho_0(T_0, p_0) + \tilde{\rho}(\tilde{T}, \tilde{p})$  with  $T = T_0 + \tilde{T}$  and  $p = p_0 + \tilde{p}$  where the 0-subscript denotes reference state quantities and the 'tilde' variables the perturbations with respect to the reference state. The reference state quantities do not depend on time but can depend on the spatial coordinates. The usual assumption is that  $\tilde{\rho}(\tilde{T}, \tilde{p})$  depends linearly on the perturbations in temperature and pressure. In this case one can use the thermodynamic relation

$$d\rho = \left( \frac{\partial \rho}{\partial T} \right)_P dT + \left( \frac{\partial \rho}{\partial p} \right)_T dp = -\alpha \rho dT + K_T^{-1} \rho dp$$

where  $\alpha$  is the thermal expansion coefficient and  $K_T$  the isothermal incompressibility, or isothermal bulk modulus<sup>11</sup> (units: Pa). Assuming linearity, this relation is up-scaled to the macroscopic mantle to obtain the equation of state

$$\rho(\vec{r}, t) = \rho_0(\vec{r})(1 - \alpha \tilde{T} + K_T^{-1} \tilde{p}) \quad (11.88)$$

hence

$$\tilde{\rho} = -\alpha \rho_0 \tilde{T} + K_T^{-1} \rho_0 \tilde{p}$$

A simple extension to a 2-phase medium consisting of materials with reference densities  $\rho_0$  and  $\rho_1$  is obtained by assuming that the dependence of density of  $T$  and  $p$  is the same for both phases. Then, only an addition factor is needed to account for the density of a mixed composition

$$\rho(\vec{r}, t) = \rho_0(\vec{r}) \left( 1 + \Gamma \left( \frac{\rho_1 - \rho_0}{\rho_0} \right) \right) (1 - \alpha \tilde{T} + K_T^{-1} \tilde{p}) \quad (11.89)$$

Taking  $\rho_1 > \rho_0$ , the phase distribution  $\Gamma$  can assume values between 0 and 1.

**Exercise: 31.** Assume an isothermal (i.e.  $\tilde{T} = 0$ ) and incompressible fluid ( $K_T^{-1} = 0$ ) and derive from the conservation law of mass, the conservation law  $D\Gamma/Dt = 0$  for the phase distribution function.

The reference quantities satisfy equations of a motionless reference state. The definition of the reference state depends on the problem studied. Usually the reference temperature is chosen constant or to follow a mantle adiabat or a conductive geotherm. For a simple reference state with constant thermodynamic parameters the following equations suffice:

$$\vec{\nabla} p_0 = -\rho_0 \vec{\nabla} U_0 \quad (11.90)$$

$$\vec{\nabla}^2 U_0 = 4\pi \mathcal{G} \rho_0 \quad (11.91)$$

$$\vec{\nabla} \cdot (k \vec{\nabla} T_0) = 0 \quad \frac{DT_0}{Dt} = 0 \quad (11.92)$$

More complex reference states involving depth variable thermodynamic parameters require internally consistent relations between thermodynamic parameters such as  $\alpha$ ,  $K_T$ ,  $C_p$ ,  $k$  and temperature, density, and pressure.

<sup>11</sup>The bulk modulus of a substance is a measure of the resistance of a substance to bulk compression. It is defined as the ratio of the infinitesimal pressure increase to the resulting relative decrease of the volume

**Exercise: 32.** Show that an adiabatic geotherm approximately satisfies the heat equation for the reference state.

### 11.8.1 The complete set of perturbation equations

The following equations (11.93), (11.94), (11.95), (11.96), (11.97) give the relevant equations in terms of the perturbations w.r.t. the adiabatic reference state as defined with equations (11.90), (11.91) and (11.92). Some terms then cancel in the general equations; no approximations have been made.

The equation of state is

$$\rho(\vec{r}, t) = \rho_0(\vec{r}) \left( 1 - \alpha \tilde{T} + K_T^{-1} \tilde{p} \right) \quad (11.93)$$

The heat equation in terms of  $\tilde{T}$

$$\rho C_p \frac{D\tilde{T}}{Dt} = \frac{\partial}{\partial x_j} \left( k \frac{\partial \tilde{T}}{\partial x_j} \right) + \rho H + \pi_{ij} \frac{\partial \mathbf{v}_i}{\partial x_j} + \alpha T \frac{Dp}{Dt} \quad (11.94)$$

The equation of motion is in terms of  $\tilde{p}$ ,  $\tilde{\rho}$ ,  $\tilde{T}$ :

$$\rho \frac{D\mathbf{v}_i}{Dt} = - \frac{\partial \tilde{p}}{\partial x_i} + \frac{\partial \pi_{ij}}{\partial x_j} - \rho_0 (\alpha \tilde{T} + K_T^{-1} \tilde{p}) g_i \quad (11.95)$$

The Poisson equation in terms of  $\tilde{U}$  and  $\tilde{\rho}$ :

$$\vec{\nabla}^2 \tilde{U} = 4\pi \mathcal{G} \tilde{\rho} \quad (11.96)$$

and, finally, the continuity equation in terms of  $\tilde{\rho}$  reads

$$\frac{\partial \tilde{\rho}}{\partial t} + \frac{\partial (\rho_0 + \tilde{\rho}) \mathbf{v}_j}{\partial x_j} = 0 \quad (11.97)$$

## 11.9 Scaling of equations

An important subject is the scaling of equations by means of scaling the relevant parameters with appropriate estimates (expected values) for the convection problem at hand. Here, we discuss scaling of parameters for the study of whole mantle convection, however, other studies such as boundary layer modeling or subduction modeling may require different scaling parameters. The idea behind scaling of equations is to find out the relative importance (contribution) of separate terms. Usual scaling parameters for mantle convection are:

$$x_i = x'_i h \quad t = t' \frac{h^2}{\kappa} \quad \mathbf{v}_i = \mathbf{v}'_i \frac{\kappa}{h} \quad \tilde{T} = \tilde{T}' \Delta T_m \quad \tilde{p} = \tilde{p}' \frac{\eta \kappa}{h^2} \quad \rho = \rho' \rho_m$$

where  $h$  is the mantle thickness,  $\kappa = k/(\rho C_p)$  is the thermal diffusivity,  $\Delta T_m$  is the temperature difference between the Core-Mantle boundary and the surface,  $\eta$  is the dynamic viscosity<sup>12</sup> and  $\rho_m$  is a scaling density.

Note that on Earth, we have  $\alpha \sim 3 \cdot 10^{-5} \text{ K}^{-1}$ ,  $K_T \sim 3 \cdot 10^{12} \text{ Pa}$ ,  $\Delta T_m \sim 10^3 \text{ K}$ ,  $C_p \sim 10^3 \text{ J kg}^{-1} \text{ K}^{-1}$ ,  $k \sim 4 \text{ W m}^{-1} \text{ K}^{-1}$ ,  $\kappa \sim 1.2 \times 10^{-6} \text{ m}^2 \text{ s}^{-1}$ .

The primed quantities are **dimensionless** and are of the order 1 if the scaling is correct for the problem at hand<sup>13</sup>. Application of scaling to the equation of state (11.93) leads to

$$\rho' = \rho'_0 \left( 1 - \alpha \Delta T_m \tilde{T}' + \frac{K_T^{-1} \eta \kappa}{h^2} \tilde{p}' \right) \quad (11.98)$$

<sup>12</sup>Note that the viscosity varies by many orders of magnitude, i.e. between  $\sim 10^{17}$  and  $\sim 10^{26}$  in the mantle and lithosphere, so finding the appropriate scaling value is necessary.

<sup>13</sup>Except for the viscosity, obviously

The coefficients of the primed-quantities determine the relative importance of the terms (with respect to 1). We have  $\alpha\Delta T_m \sim 3 \cdot 10^{-2}$  and  $K_T^{-1}\eta\kappa/h^2 \sim 2.5 \cdot 10^{-10}$ , which leads us to conclude that thermal perturbations of density are much more important than pressure perturbations (assuming the equation of state and scaling are correct for the convection problem). We note that the scaling factor of pressure is about 100 Pa, i.e. perturbations of hydrostatic pressure are expected to be small compared to the rheological stress (1-100 MPa).

The coefficient  $K_T^{-1}\eta\kappa/h^2$  can be written as  $\eta C_p/k \cdot k^2 K_T^{-1}/C_p^2/\rho_m h^2 = \text{Pr} \cdot \text{M}^2$  where  $\text{Pr}$  is the Prandtl number<sup>14</sup> and  $\text{M}$  the Mach number<sup>15</sup>.  $\text{M}$  gives the ratio between flow speed and sound speed which is about  $10^{-16}$  for mantle convection. The Prandtl number is about  $10^{24}$  and will appear in the coefficient of the inertial term of the dimensionless equation of motion.

Let us now turn to the scaling of the continuity equation (11.97) which leads to some basic insight into the problem of (in-)compressibility. The dimensionless version is

$$\frac{\partial \tilde{\rho}'}{\partial t'} + \frac{\partial(\rho'_0 + \tilde{\rho})\mathbf{v}'_j}{\partial x'_j} = 0 \quad ; \quad \rho'_0 = \frac{\rho_0(\vec{r})}{\rho_m} \quad (11.99)$$

**Exercise: 33.** Derive Eq. (11.99).

Substitution of the equation of state (11.93) and taking the limit  $Pr \cdot M^2 \rightarrow \infty$  leads to

$$-\rho'_0\alpha\Delta T_m \frac{\partial \tilde{T}'}{\partial t'} + \frac{\partial}{\partial x'_j} \left[ (\rho'_0 - \rho'_0\alpha\Delta T_m \tilde{T}')\mathbf{v}'_j \right] = 0 \quad (11.100)$$

which gives the conservation law based on thermal perturbations of density only. In the limit that the perturbation terms are very small we arrive at

$$\frac{\partial}{\partial x'_j} (\rho'_0\mathbf{v}'_j) = 0 \quad (11.101)$$

which is the equation of anelastic conservation of mass. Compared to (11.97) we effectively replaced the density by the reference density. The time derivative of density has disappeared. This derivative is primarily related to the propagation of seismic waves. This process occurs at totally different time scales than mantle convection. The anelastic conservation of mass (11.101) is being used in convection modeling of compressible fluids, see for instance the equations of the ASPECT code<sup>16</sup>:

Specifically, we consider the following set of equations for velocity  $\mathbf{u}$ , pressure  $p$  and temperature  $T$ , as well as a set of advected quantities  $c_i$  that we call *compositional fields*:

$$-\nabla \cdot \left[ 2\eta \left( \varepsilon(\mathbf{u}) - \frac{1}{3}(\nabla \cdot \mathbf{u})\mathbf{1} \right) \right] + \nabla p = \rho \mathbf{g} \text{ in } \Omega, \quad (1)$$

$$\nabla \cdot (\rho \mathbf{u}) = 0 \text{ in } \Omega, \quad (2)$$

$$\begin{aligned} \rho C_p \left( \frac{\partial T}{\partial t} + \mathbf{u} \cdot \nabla T \right) - \nabla \cdot k \nabla T &= \rho H \\ &+ 2\eta \left( \varepsilon(\mathbf{u}) - \frac{1}{3}(\nabla \cdot \mathbf{u})\mathbf{1} \right) : \left( \varepsilon(\mathbf{u}) - \frac{1}{3}(\nabla \cdot \mathbf{u})\mathbf{1} \right) \\ &+ \alpha T (\mathbf{u} \cdot \nabla p) \\ &+ \rho T \Delta S \left( \frac{\partial X}{\partial t} + \mathbf{u} \cdot \nabla X \right) \text{ in } \Omega, \end{aligned} \quad (3)$$

$$\frac{\partial c_i}{\partial t} + \mathbf{u} \cdot \nabla c_i = q_i \text{ in } \Omega, i = 1 \dots C \quad (4)$$

where  $\varepsilon(\mathbf{u}) = \frac{1}{2}(\nabla \mathbf{u} + \nabla \mathbf{u}^T)$  is the symmetric gradient of the velocity (often called the *strain rate*)<sup>[1]</sup>.

In this set of equations, (1) and (2) represent the compressible Stokes equations in which  $\mathbf{u} = \mathbf{u}(\mathbf{x}, t)$  is the velocity field and  $p = p(\mathbf{x}, t)$  the pressure field. Both fields depend on space  $\mathbf{x}$  and time  $t$ . Fluid flow is driven by the gravity force that acts on the fluid and that is proportional to both the density of the fluid and the strength of the gravitational pull.

Taken from the ASPECT website. Note that the equations are not dimensionless.

<sup>14</sup>[https://en.wikipedia.org/wiki/Prandtl\\_number](https://en.wikipedia.org/wiki/Prandtl_number)

<sup>15</sup>[https://en.wikipedia.org/wiki/Mach\\_number](https://en.wikipedia.org/wiki/Mach_number)

<sup>16</sup><https://aspect-documentation.readthedocs.io/en/latest/user/methods/basic-equations/index.html>

In Solheim and Peltier [1177] we read:

“The anelastic-liquid approximation [e.g., Jarvis and McKenzie [637] (1980)] is applied to the system [mass, momentum, energy conservation equations]. This involves setting  $\partial\rho/\partial t = 0$  in the mass conservation equation, assuming  $g$ ,  $C_p$ ,  $\alpha$ ,  $k$ ,  $K_T$  and  $\kappa$  to be known functions of radius and replacing  $\rho$  by  $\rho_r$  everywhere in [the equations] except in the body force term of the momentum equation. Because the mantle has essentially infinite Prandtl number, the inertial force term in the momentum conservation equation may be neglected [986]. Furthermore, owing to the extremely small velocities associated with the mantle convection process, the pressure distribution is very nearly that of a fluid in hydrostatic equilibrium and we may then safely assume

$$\alpha T \frac{Dp}{Dt} \simeq -\alpha T \rho g u_r$$

where  $u_r$  is the radial component of the velocity. These approximations have been discussed in greater detail by Solheim and Peltier [1178] (1990). ”

## 11.10 The Boussinesq approximation

The Boussinesq<sup>17</sup> approximation is ubiquitous in computational geodynamics [1186, 1387]. This approximation leads to a set of simplified equations that are easier to solve for some analytical problems and in numerical modeling of convective flow. The following approximations are being made:

- Neglect the effect of density variations with respect to a reference state except in terms related to the driving force of convection
- assume that the divergence of the velocity field is 0 (in this case  $\boldsymbol{\pi} = \boldsymbol{\tau}$ ).
- Only concern temperature changes resulting from diffusion and advection (i.e. neglect terms related to adiabatic compression, and heat dissipation)

This simplifies Eqs. (11.93),(11.94),(11.95),(11.96),(11.97) to

$$\rho(\vec{r}, t) = \rho_0(1 - \alpha \tilde{T}) \quad (11.102)$$

$$\rho_0 C_p \frac{D\tilde{T}}{Dt} = \frac{\partial}{\partial x_j} \left( k \frac{\partial \tilde{T}}{\partial x_j} \right) + \rho_0 H \quad (11.103)$$

$$\rho_0 \frac{D\mathbf{v}_i}{Dt} = -\frac{\partial \tilde{p}}{\partial x_i} + \frac{\partial \tau_{ij}}{\partial x_j} - \rho_0 \alpha \tilde{T} g_i \quad (11.104)$$

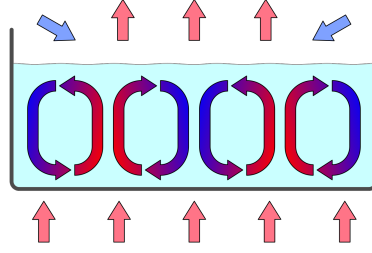
$$\vec{\nabla}^2 \tilde{U} = -4\pi \mathcal{G} \rho_0 \alpha \tilde{T} \quad (11.105)$$

$$\frac{\partial \mathbf{v}_j}{\partial x_j} = 0 \quad (11.106)$$

### 11.10.1 The Rayleigh-Bénard convection

This is the classical example of a laterally unlimited fluid layer of thickness  $h$  in a gravity field. We derive here the pertinent equations in the Boussinesq approximation. In the next section these equations are being used to study the problem of onset of convection.

<sup>17</sup>[https://en.wikipedia.org/wiki/Joseph\\_Valentin\\_Boussinesq](https://en.wikipedia.org/wiki/Joseph_Valentin_Boussinesq)



Convection cells in a gravity field. Taken from [https://en.wikipedia.org/wiki/Rayleigh-Benard\\_convection](https://en.wikipedia.org/wiki/Rayleigh-Benard_convection)

Gravity is directed along the positive  $z$ -axis. The top and bottom of the layer are kept at a constant temperature:

$$T(x, z = 0, t) = T_0 \quad \text{and} \quad T(x, z = h, t) = T_0 + \Delta T \quad (\Delta T > 0) \quad (11.107)$$

The mechanical boundary conditions on the top and bottom are impermeability and shear stress free (also called free slip):

$$\mathbf{v}_z(x, z = 0, t) = 0 \quad (11.108)$$

$$\mathbf{v}_z(x, z = h, t) = 0 \quad (11.109)$$

$$\tau_{xz}(x, z = 0, t) = 0 \quad (11.110)$$

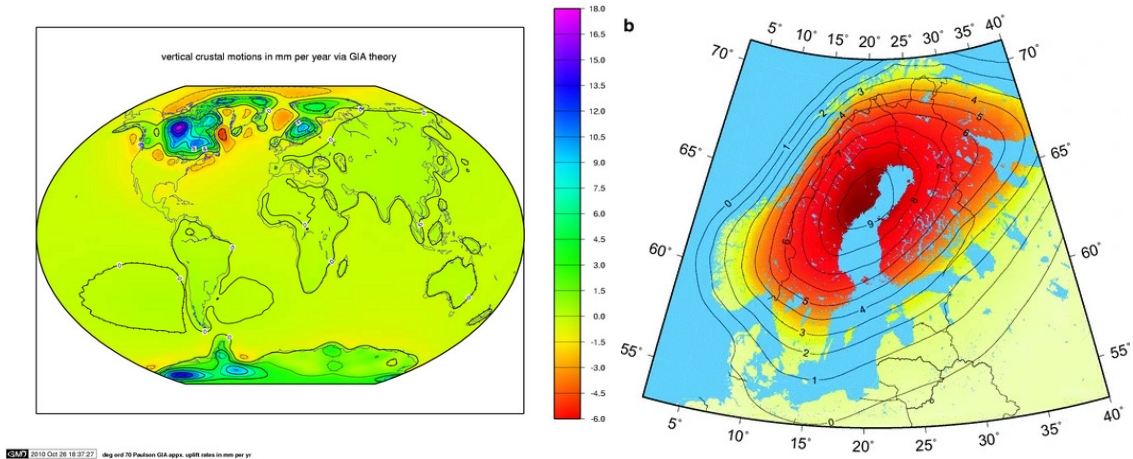
$$\tau_{xz}(x, z = h, t) = 0 \quad (11.111)$$

$$\tau_{yz}(x, z = 0, t) = 0 \quad (11.112)$$

$$\tau_{yz}(x, z = h, t) = 0 \quad (11.113)$$

**Exercise: 34.** Show for an incompressible Newtonian fluid that at  $z = 0$  the traction  $\vec{t}^n$  is given by  $t_x^n = t_y^n = 0$  and  $t_z^n = -p + 2\eta\partial\mathbf{v}_z/\partial z$ . What can you tell about the horizontal components of velocity?

The condition of zero vertical velocity is only approximately valid at the Earth's surface. Vertical motions can now be measured by geodetic techniques and are generally one order of magnitude less than horizontal motions. Also from the geological past we know that horizontal motions have had much larger amplitude than vertical motions. If the Earth would have a pure fluid behavior (has it?) then mantle flow can induce **dynamic surface topography**. The vertical velocity  $\mathbf{v}_z$  equals the local time derivative  $\partial d/\partial t$  of the vertical surface deflection  $d(t)$ . In numerical solutions (with impermeable boundaries along which pressure variation can accumulate) the deflection is often computed a posteriori by equating  $\rho g d \sim -p + 2\eta\partial\mathbf{v}_z/\partial z$ .



Left: A model of present-day mass change due to post-glacial rebound and the reloading of the ocean basins with seawater. Blue and purple areas indicate rising due to the removal of the ice sheets. Yellow and red areas indicate falling as mantle material moved away from these areas in order to supply the rising areas, and because of the collapse of the forebulges around the ice sheets. Taken from Wikipedia<sup>18</sup>, after Paulson, Zhong, and Wahr [983]; Right: The absolute land uplift in Fennoscandia in mm/year. The Finnish Geodetic Institute, taken from Kakkuri [664].

In the following we assume an isoviscous Newtonian fluid with constant thermal conductivity, constant heat capacity and no heat production and we adopt the Boussinesq approximation. Then we have the following equations for Rayleigh-Benard convection:

$$\tilde{\rho} = -\rho_0 \alpha \tilde{T} \quad (11.114)$$

$$\frac{D\tilde{T}}{Dt} = \kappa \vec{\nabla}^2 \tilde{T} \quad (11.115)$$

$$\rho_0 \frac{D\mathbf{v}_i}{Dt} = -\frac{\partial \tilde{p}}{\partial x_i} + \eta \vec{\nabla}^2 \mathbf{v}_i - \rho_0 \alpha \tilde{T} g_i \quad (11.116)$$

$$\vec{\nabla}^2 \tilde{U} = -4\pi \mathcal{G} \rho_0 \alpha \tilde{T} \quad (11.117)$$

$$\frac{\partial \mathbf{v}_j}{\partial x_j} = 0 \quad (11.118)$$

The reference temperature is assumed to be constant or following an adiabat or is a time stationary conductive geotherm. Reference pressure is computed from the hydrostatic equation  $\vec{\nabla} p = \rho_0 \vec{g}$ .

We assume a 2-D situation in which the velocity vector is denoted by  $\vec{\mathbf{v}} = (u, 0, w)^T$ . The mechanical boundary condition becomes  $\tau_{xz}(x, z = 0) = \tau_{xz}(x, z = h) = 0$  and  $\eta \partial u / \partial z(x, z = 0) = \eta \partial u / \partial z(x, z = h) = 0$ . The scaling of the equation of motion gives

$$\frac{1}{\text{Pr}} \frac{D\mathbf{v}'_i}{Dt'} = -\frac{\partial \tilde{p}'}{\partial x'_i} + (\nabla')^2 \mathbf{v}'_i - \text{Ra} \tilde{T}' \delta_{zi} \quad (11.119)$$

where the Prandtl number is  $\text{Pr} = \eta C_p / k = \eta / \rho_0 \kappa$  and the Rayleigh number is  $\text{Ra} = \rho_0 \alpha g \Delta T h^3 / \eta \kappa$ . For the Earth's mantle  $\text{Pr} \sim 10^{24}$  and  $\text{Ra} \sim 10^7$  which leads to the conclusion that the inertial term can be neglected. Please check Section 2.11.2 for the complete adimensionalisation of the temperature-dependent Navier-Stokes equations. We arrive at the non-dimensional equation (dropping the primes for convenience):

$$0 = -\vec{\nabla} \tilde{p} + \vec{\nabla}^2 \vec{\mathbf{v}} - \text{Ra} \tilde{T} \vec{e}_z \quad (11.120)$$

in which  $\text{Ra}$  is the only free parameter. Recall that  $\tilde{p}$  is the pressure anomaly with regards to  $p_0$ . The Rayleigh number determines the magnitude of the force driving convection (see computer practical).

Because we consider a 2-D situation of an incompressible fluid we can adopt the stream function approach with  $u = \partial \Psi / \partial z$  and  $w = -\partial \Psi / \partial x$  which leads to

$$\vec{\nabla}^4 \Psi = -\text{Ra} \frac{\partial \tilde{T}}{\partial x} \quad (11.121)$$

(to be compared with (11.48)) with boundary conditions  $\Psi(x, z = 0) = \Psi(x, 1) = 0$  and  $\partial^2 \Psi / \partial z^2(x, y = 0) = \partial^2 \Psi / \partial z^2(x, y = 1) = 0$ .

The energy equation becomes after scaling

$$\frac{D\tilde{T}}{Dt} = \vec{\nabla}^2 \tilde{T} \quad (11.122)$$

<sup>18</sup>[https://en.wikipedia.org/wiki/Post-glacial\\_rebound](https://en.wikipedia.org/wiki/Post-glacial_rebound)



Equation (11.121) and (11.122) are coupled non-linear equations which usually need to be solved numerically.

Recall that the temperature  $\tilde{T}$  is the deviation from either a constant temperature or from adiabatic temperature  $\tilde{T} = T - T_a$ . Advective transport of adiabatic temperature does not lead to temperature changes with the surroundings and therefore  $\tilde{T}$  is the temperature associated with convective flow. In convective flow that assumes a (non-adiabatic) conductive geotherm in the background we can separate  $\tilde{T}$  as  $\tilde{T} = \tilde{T}_0 + \tilde{T}_1$  where  $\tilde{T}_0$ , is the geotherm of the conductive reference state and  $\tilde{T}_1$  is the perturbation of the temperature field.

**Exercise: 35.** Derive the stationary conductive reference temperature profile  $\tilde{T}_0'(z') = z'$  using the temperature boundary conditions given earlier (with  $T_{surface} = 0$ ).

### 11.10.2 Linear stability analysis (the onset of convection problem)

We are now ready to solve a linear stability problem which will provide fundamental insight in the role of the Rayleigh number in convection.

The system is a layer of fluid between  $y = 0$  and  $y = h$ , with boundary conditions  $T(x, y = 0) = T_b$  and  $T(x, y = h) = 0$ , characterized by  $\rho_0$ ,  $C_p$ ,  $k$ ,  $\eta_0$  which are assumed to be constant (in space and time).

The Stokes equation is  $\vec{\nabla} \cdot \boldsymbol{\sigma} + \rho \vec{g} = \vec{0}$ . The components of this equation on the  $x$ - and  $y$ -axis are:

$$\begin{aligned} (\vec{\nabla} \cdot \boldsymbol{\sigma})_x &= -\rho \vec{g} \cdot \vec{e}_x = 0 \\ (\vec{\nabla} \cdot \boldsymbol{\sigma})_y &= -\rho \vec{g} \cdot \vec{e}_y = \rho g_0 \end{aligned}$$

since  $\vec{g}$  and  $\vec{e}_y$  are in opposite directions ( $\vec{g} = -g_0 \vec{e}_y$ , with  $g_0 > 0$ ).

Following Eq. (11.48), the stream function formulation of the incompressible isoviscous Stokes equation is

$$\eta_0 \nabla^4 \Psi = \frac{\partial \rho g_y}{\partial x} - \frac{\partial \rho g_x}{\partial y} = \frac{\partial \rho g_y}{\partial x} = -g_0 \frac{\partial \rho}{\partial x}$$

since  $g_x = 0$  and  $g_y = \vec{g} \cdot \vec{e}_y = -g_0$ . Assuming a linearised density field with regards to temperature  $\rho(T) = \rho_0(1 - \alpha T)$  we have

$$\frac{\partial \rho}{\partial x} = -\rho_0 \alpha \frac{\partial T}{\partial x}$$

and then

$$\vec{\nabla}^4 \Psi = \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{\partial T}{\partial x} \quad (11.123)$$

For small perturbations of the conductive state<sup>19</sup>  $T_c(y) = (1 - y/h)T_b$  we define the temperature perturbation  $\tilde{T}(x, y)$  such that

$$T(x, y, t) = T_c(y) + \tilde{T}(x, y, t)$$

Note that the temperature perturbation  $\tilde{T}$  must satisfy the homogeneous boundary conditions  $\tilde{T}(x, y = 0) = 0$  and  $\tilde{T}(x, y = h) = 0$ . We then have<sup>20</sup>:

$$\boxed{\vec{\nabla}^4 \Psi = \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{\partial \tilde{T}}{\partial x}} \quad (11.124)$$

<sup>19</sup>The conductive state temperature is defined as the solution of the steady state diffusion equation  $\Delta T_c = 0$  subjected to the desired boundary conditions at the top and at the bottom.

<sup>20</sup>This is the same equation as in Turcotte & Schubert, eq 6.310.



In the absence of heat production, the temperature equation (in the Boussinesq approx.) is

$$\rho_0 C_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = k \Delta T \quad \Rightarrow \quad \rho_0 C_p \left( \frac{\partial(T_c + \tilde{T})}{\partial t} + \vec{v} \cdot \vec{\nabla}(T_c + \tilde{T}) \right) = k \Delta(T_c + \tilde{T}) \quad (11.125)$$

First we start by acknowledging that  $T_c$  does not depend on time, so that  $\partial T_c / \partial t = 0$ . Then,  $\Delta T_c = 0$  since  $T_c$  is a linear function of  $y$ . Finally, we assume the nonlinear term  $\vec{v} \cdot \vec{\nabla} \tilde{T}$  to be second order (temperature perturbations and coupled velocity changes are assumed to be small). In the end, defining the heat diffusion coefficient  $\kappa$  as  $\kappa = k / \rho_0 C_p$ , the energy equation can be simplified further as follows:

$$\frac{\partial \tilde{T}}{\partial t} + \vec{v} \cdot \vec{\nabla} T_c = \kappa \Delta \tilde{T}$$

Using the relationship between velocity and stream function  $\vec{v} = (u, v) = (\partial_y \Psi, -\partial_x \Psi)$  and since  $\vec{\nabla} T_c = -(T_b/h) \vec{e}_y$  then

$$\vec{v} \cdot \vec{\nabla} T_c = \begin{pmatrix} \partial_y \Psi \\ -\partial_x \Psi \end{pmatrix} \cdot \begin{pmatrix} 0 \\ -T_b/h \end{pmatrix} = \frac{T_b}{h} \frac{\partial \Psi}{\partial x}$$

and finally<sup>21</sup>:

$$\boxed{\frac{\partial \tilde{T}}{\partial t} - \kappa \Delta \tilde{T} = -\frac{T_b}{h} \frac{\partial \Psi}{\partial x}} \quad (11.126)$$

Looking at these equations, we immediately think about a separation of variables approach to solve these equations. Both equations showcase the Laplace operator  $\Delta$ , and the eigenfunctions of the biharmonic operator and the Laplace operator are the same. We then pose that  $\Psi$  and  $\tilde{T}$  can be written<sup>22</sup>:

$$\tilde{T}(x, y, t) = \tilde{T}_0 \exp(pt) [a_k \cos(k_x x) + b_k \sin(k_x x)] [c_k \cos(k_y y) + d_k \sin(k_y y)]$$

$$\Psi(x, y, t) = \Psi_0 \exp(pt) [\alpha_k \cos(k_x x) + \beta_k \sin(k_x x)] [\delta_k \cos(k_y y) + \gamma_k \sin(k_y y)]$$

where  $\tilde{T}_0$  and  $\Psi_0$ ,  $a_k, b_k, c_k, d_k$  and  $\alpha_k, \beta_k, \delta_k, \gamma_k$  are constants. We then of course have

$$\begin{aligned} \vec{\nabla}^2 \tilde{T} &= \frac{\partial^2 \tilde{T}}{\partial x^2} + \frac{\partial^2 \tilde{T}}{\partial y^2} \\ &= \tilde{T}_0 \exp(pt) \{ [-k_x^2 a_k \cos(k_x x) - k_x^2 b_k \sin(k_x x)] [c_k \cos(k_y y) + d_k \sin(k_y y)] \} \\ &\quad + \tilde{T}_0 \exp(pt) \{ [a_k \cos(k_x x) + b_k \sin(k_x x)] [-k_y^2 c_k \cos(k_y y) - k_y^2 d_k \sin(k_y y)] \} \\ &= -(k_x^2 + k_y^2) \tilde{T} \end{aligned} \quad (11.127)$$

and a similar expression for  $\Psi$ .

The boundary conditions on  $\tilde{T}$  are  $\tilde{T}(y=0) = \tilde{T}(y=h) = 0$ . From the first one it follows immediately that  $c_k = 0$ . From the second, we arrive at  $k_y h = n\pi$ , which yields  $\sin(n\pi y/h)$  where  $n$  is an integer. We then arrive at the following expression for the temperature  $\tilde{T}$ :

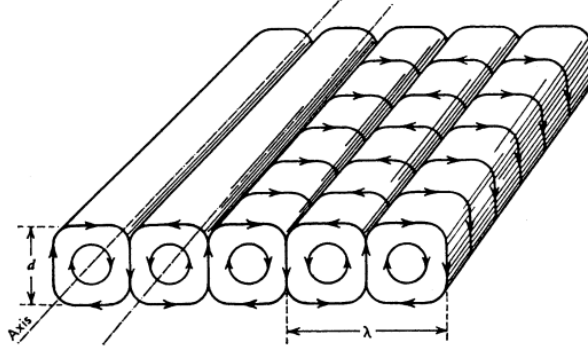
$$\tilde{T}(x, y, t) = \tilde{T}_0 \exp(pt) [a_k \cos(k_x x) + b_k \sin(k_x x)] \sin\left(n\pi \frac{y}{h}\right)$$

where  $n$  is an integer number.

At this stage we make an important assumption: at  $t = 0$  we then only consider a single horizontal periodic perturbation with wavelength  $\lambda$  as depicted below. In such a case we expect that it would in fact lead to the formation of the following convection cells:

<sup>21</sup>This is the same equation as in Turcotte & Schubert, eq 6.309.

<sup>22</sup>T & S actually very much postulate the final form of these quantities without (enough?) justification. I would like to revisit this in the future and better support this assertion.



Note that the number of cells left and right of those shown is infinite. Source unknown. The coordinate  $x = 0$  is set to the left vertical dashed line for convenience.

The boundary conditions are free slip at the top and at the bottom, i.e.  $v(y = 0) = v(y = h) = 0$ . Also, by symmetry of the perturbation we see that  $u = 0$  at each 'side' (i.e. on the vertical dashed lines of the figure above), i.e. for  $x = 0$  and  $x = \lambda$ .

Let us now turn to the vertical  $y$  component of the velocity:

$$\begin{aligned}
 v &= -\frac{\partial \Psi}{\partial x} \\
 &= -\frac{\partial}{\partial x} \{ \Psi_0 \exp(pt) [\alpha_k \cos(k_x x) + \beta_k \sin(k_x x)] [\delta_k \cos(k_y y) + \gamma_k \sin(k_y y)] \} \\
 &= -\Psi_0 \exp(pt) [-\alpha_k k_x \sin(k_x x) + \beta_k k_x \cos(k_x x)] [\delta_k \cos(k_y y) + \gamma_k \sin(k_y y)] \quad (11.128)
 \end{aligned}$$

The boundary condition at the bottom is  $v(y = 0) = 0$ , so that  $\delta_k = 0$  here again. The boundary condition at the top is  $v(y = h) = 0$ , so that  $k_y h = n\pi$  as before. Then

$$\Psi(x, y, t) = \Psi_0 \exp(pt) [\alpha_k \cos(k_x x) + \beta_k \sin(k_x x)] \sin\left(n\pi \frac{y}{h}\right)$$

Turning now to the horizontal component of the velocity:

$$\begin{aligned}
 u &= \frac{\partial \Psi}{\partial y} \\
 &= \frac{\partial}{\partial y} \left\{ \Psi_0 \exp(pt) [\alpha_k \cos(k_x x) + \beta_k \sin(k_x x)] \sin\left(n\pi \frac{y}{h}\right) \right\} \\
 &= \Psi_0 \exp(pt) [\alpha_k \cos(k_x x) + \beta_k \sin(k_x x)] \frac{n\pi}{h} \sin\left(n\pi \frac{y}{h}\right) \quad (11.129)
 \end{aligned}$$

Using now the 'side' boundary conditions:  $u(x = 0) = 0$  yields  $\alpha_k = 0$  and  $u(x = \lambda) = 0$  yields  $k_y \lambda = 2\pi$  so that in the end:

$$\boxed{\Psi(x, y, t) = \Psi_0 \exp(pt) \sin\left(\frac{2\pi}{\lambda} x\right) \sin\left(\frac{n\pi}{h} y\right)} \quad (11.130)$$

Looking at the biharmonic equation (11.126), its rhs is  $\sim \frac{\partial \Psi}{\partial x}$ . Then, the  $x$  dependency of this term will be  $\cos(2\pi x/\lambda)$ . The lhs term of (11.126) is proportional to  $\tilde{T}$  (see Eq. (11.127)), i.e. proportional to  $a_k \cos(k_x x) + b_k \sin(k_x x)$ . For these equations to be compatible, we must set  $b_k = 0$  and we then obtain<sup>23</sup>

<sup>23</sup>Taking  $n = 1$  and remembering that Turcotte & Schubert have the domain between  $y - h/2$  and  $y = h/2$ , these expressions are identical to Eqs. 6.311 and 6.312 of the book. Also one could have assigned  $\partial T/\partial x$  on the sides for symmetry reasons and have obtained the same expression.

$$\boxed{\tilde{T}(x, y, t) = \tilde{T}_0 \exp(pt) \cos\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right)} \quad (11.131)$$

where  $a_k$  has been 'absorbed' in  $\tilde{T}_0$ .

Then the two framed PDEs above, Eq. (11.124) and Eq. (11.126), when coupled with Eq. (11.130) and Eq. (11.131), become:

$$\begin{aligned} \nabla^4 \Psi &= \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{\partial \tilde{T}}{\partial x} \\ \Rightarrow \nabla^2 \left[ \left( -\frac{4\pi^2}{\lambda^2} - \frac{n^2 \pi^2}{h^2} \right) \Psi \right] &= \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{\partial \tilde{T}}{\partial x} \\ \Rightarrow \left( -\frac{4\pi^2}{\lambda^2} - \frac{n^2 \pi^2}{h^2} \right)^2 \Psi &= \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{\partial \tilde{T}}{\partial x} \\ \Rightarrow \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right)^2 \Psi_0 \exp(pt) \sin\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right) &= \frac{\rho_0 g_0 \alpha}{\eta_0} \cdot -\frac{2\pi}{\lambda} \tilde{T}_0 \exp(pt) \sin\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right) \\ \Rightarrow \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right)^2 \Psi_0 &= -\frac{\rho_0 g_0 \alpha}{\eta_0} \frac{2\pi}{\lambda} \tilde{T}_0 \end{aligned} \quad (11.132)$$

$$\begin{aligned} \frac{\partial \tilde{T}}{\partial t} - \kappa \Delta \tilde{T} &= -\frac{T_b}{h} \frac{\partial \Psi}{\partial x} \\ \Rightarrow p \tilde{T} - \kappa \left( -\frac{4\pi^2}{\lambda^2} - \frac{n^2 \pi^2}{h^2} \right) \tilde{T} &= -\frac{T_b}{h} \cdot \frac{2\pi}{\lambda} \Psi_0 \exp(pt) \cos\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right) \\ \Rightarrow \left[ p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right) \right] \tilde{T}_0 \exp(pt) \cos\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right) &= -\frac{T_b}{h} \frac{2\pi}{\lambda} \Psi_0 \exp(pt) \cos\left(\frac{2\pi}{\lambda}x\right) \sin\left(n\pi\frac{y}{h}\right) \\ \Rightarrow \left[ p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right) \right] \tilde{T}_0 &= -\frac{T_b}{h} \frac{2\pi}{\lambda} \Psi_0 \end{aligned} \quad (11.133)$$

We are then left with two equations:

$$\begin{aligned} \left[ p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right) \right] \tilde{T}_0 &= -\frac{T_b}{h} \frac{2\pi}{\lambda} \Psi_0 \\ \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right)^2 \Psi_0 &= -\frac{\rho_0 g_0 \alpha}{\eta_0} \frac{2\pi}{\lambda} \tilde{T}_0 \end{aligned}$$

which we can cast as

$$\begin{pmatrix} p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right) & \frac{T_b}{h} \frac{2\pi}{\lambda} \\ -\frac{\rho_0 g_0 \alpha}{\eta_0} \frac{2\pi}{\lambda} & -\left( \frac{4\pi^2}{\lambda^2} + \frac{n^2 \pi^2}{h^2} \right)^2 \end{pmatrix} \begin{pmatrix} \tilde{T}_0 \\ \Psi_0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

The determinant of the matrix should be zero to have non-trivial solutions<sup>24</sup> for the amplitude factors

---

<sup>24</sup>Let us consider the following matrix

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \cdot \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \Rightarrow \quad \begin{pmatrix} ac & bc \\ 0 & ad - bc \end{pmatrix} \cdot \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

where we have multiplied the first row by  $c$  and the second row by  $a$  and subtract row 1 from row 2. The lower right term  $ad - bc$  is the determinant and we find that it must be equal to zero since  $y$  is not zero.

(i.e.  $\tilde{T}_0 = 0$  and  $\Psi_0 = 0$  which is not helpful). This leads to the condition:

$$\begin{aligned}
Det &= \left[ p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right) \right] \cdot - \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2 + \frac{\rho_0 g_0 \alpha}{\eta_0} \frac{2\pi}{\lambda} \cdot \frac{T_b}{h} \frac{2\pi}{\lambda} \\
&= - \left[ p + \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right) \right] \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2 + \frac{\rho_0 g_0 \alpha T_b}{h \eta_0} \frac{4\pi^2}{\lambda^2} \\
&= -p \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2 - \kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \frac{\rho_0 g_0 \alpha T_b}{h \eta_0} \frac{4\pi^2}{\lambda^2}
\end{aligned}$$

The determinant is zero for

$$\begin{aligned}
p \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2 &= -\kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \frac{\rho_0 g_0 \alpha T_b}{h \eta_0} \frac{4\pi^2}{\lambda^2} \\
p &= \frac{-\kappa \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \frac{\rho_0 g_0 \alpha T_b}{h \eta_0} \frac{4\pi^2}{\lambda^2}}{\left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2} \\
&= \kappa \frac{-\left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \frac{\rho_0 g_0 \alpha T_b}{h \kappa \eta_0} \frac{4\pi^2}{\lambda^2}}{\left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2} \\
&= \frac{\kappa}{h^6} \frac{-h^6 \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \frac{\rho_0 g_0 \alpha T_b h^3}{\kappa \eta_0} \frac{4\pi^2 h^2}{\lambda^2}}{\left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2} \\
&= \frac{\kappa}{h^2} \frac{-h^6 \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^3 + \text{Ra} \frac{4\pi^2 h^2}{\lambda^2}}{h^4 \left( \frac{4\pi^2}{\lambda^2} + \frac{n^2\pi^2}{h^2} \right)^2} \\
&= \frac{\kappa}{h^2} \frac{-\left( \frac{4\pi^2 h^2}{\lambda^2} + n^2 \pi^2 \right)^3 + \text{Ra} \frac{4\pi^2 h^2}{\lambda^2}}{\left( \frac{4\pi^2 h^2}{\lambda^2} + n^2 \pi^2 \right)^2} \tag{11.134}
\end{aligned}$$

where we have used the Rayleigh number of the system defined as

$$\text{Ra} = \frac{\rho_0 g_0 \alpha T_b h^3}{\eta_0 \kappa}$$

The coefficient  $p$  inside  $\exp(pt)$  present in both temperature and stream function expressions determines the stability of the system: if it is negative, the system is stable and both  $\Psi$  and  $\tilde{T}$  will decay to zero (return to conductive state). If  $p = 0$ , then the system is meta-stable, and if  $p > 0$  then the system is unstable and the perturbations will grow.

In case of the stable regime an initial temperature perturbation will die out (e.g. because conduction wins from advection). In the case of the unstable regime convection will occur with an exponential growth factor. The intermediate, marginally stable, regime is the transition between convection and no convection for which our linearization applies.

The threshold is then  $p = 0$  and the corresponding critical Rayleigh number  $\text{Ra}_c$  is<sup>25</sup>:

$$\text{Ra}_c = \frac{\left( \frac{4\pi^2 h^2}{\lambda^2} + n^2 \pi^2 \right)^3}{\frac{4\pi^2 h^2}{\lambda^2}}$$

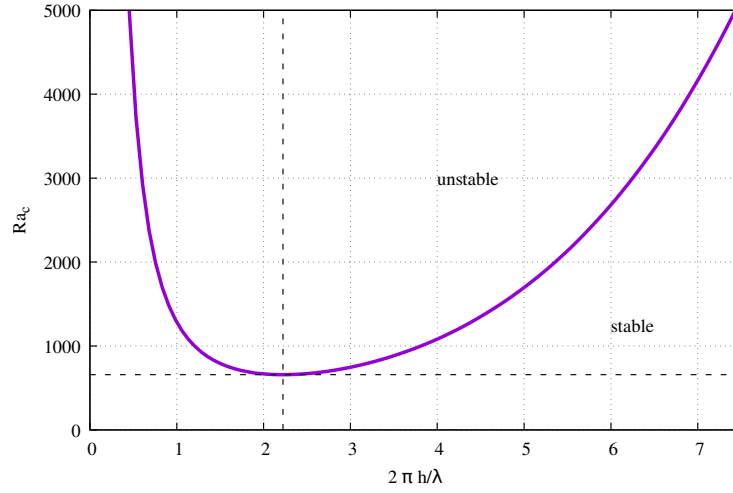
---

<sup>25</sup>this is eq 6.319 of T&S for n=1

Let us denote  $\underline{h} = 2\pi h/\lambda$  the dimensionless thickness of the layer. The critical Rayleigh number is then a function of  $\underline{h}$ :

$$\text{Ra}_c(\underline{h}) = \frac{(\underline{h}^2 + n^2\pi^2)^3}{\underline{h}^2}$$

It is plotted on the following figure for  $n = 1$ :



Critical Rayleigh number  $\text{Ra}_c$  for the onset of convection in a layer heated from below with stress-free boundaries as a function of dimensionless wavenumber  $2\pi h/\lambda$  and for  $n = 1$ . For a system with  $\text{Ra} = 2000$  then convection cannot occur for  $2\pi h/\lambda < 0.8$  and  $2\pi h/\lambda > 5.4$ . The dashed lines indicate the minimum critical Rayleigh number and its corresponding  $\underline{h}$  value. Unstable means that perturbations will grow and yield convection, while stable means that perturbations will diffuse away. Gnuplot script in `images/chapter.md`.

The minimum critical Rayleigh number is given by

$$\left. \frac{\partial \text{Ra}_c}{\partial (2\pi h/\lambda)} \right|_{n=1} = 0$$

We find<sup>26</sup> that the value of the wavelength corresponding to the smallest value of the critical Rayleigh number is  $\lambda = 2\sqrt{2}h$ , or  $\underline{h} = \pi/\sqrt{2} \simeq 2.22$  and substitution of this value for the wavelength gives the critical Rayleigh number

$$\text{Ra}_c = \frac{27}{4}\pi^4 \simeq 657.5$$

This solves the linearised onset of convection problem in the sense that an unstable layering (cold above hot) only starts convecting after a critical Rayleigh number has been overcome, e.g., by an increased  $\Delta T$

These numbers hold for a model with boundaries that are isothermal, impermeable, and free slip. Adopting other boundary conditions leads to different critical numbers. For instance, in the extreme of having fixed (no slip) boundaries one obtains  $\text{Ra} \simeq 1707.8$  and  $\lambda \simeq 2.016h$  demonstrating that it is more difficult to initiate convection compared to free slip boundaries.

When conducting a similar analysis in a spherical shell, minimum critical Rayleigh numbers prove to be much larger. Free slip (rigid) conditions at the surface and bottom 'CMB' boundary lead to  $\text{Ra}_{c,\min} \sim 14,000(35,000)$  with a critical wave length of spherical harmonic degree  $L = 3(4)$ . As the main difference between a flat layer and a spherical shell is the geometry, apparently, convection in a spherical shell experiences strong geometrical constraints (less "space" to flow near the bottom than near the top of the layer and more cooling at the surface compared to less heat input at the bottom).

For realistic values of the physical parameters defining the Rayleigh number and a realistic layer thickness, the only quantity that changes the Rayleigh number is the temperature difference  $\Delta T$

<sup>26</sup>eq 6.320 of T&S

between top and bottom. **The critical minimum Rayleigh number thus determines the critical  $\Delta T$  below which no convection occurs and above which convection is enhanced.** The analysis above is only valid in the linear regime, i.e. near the critical Rayleigh number.

Estimates for the Rayleigh number of the Earth's mantle vary between  $5 \cdot 10^5$  (upper mantle) to  $6 \cdot 10^7$  for the whole mantle which is by many factors larger than the minimum critical Rayleigh numbers that follow from experiments as described above. The mantle is in a state of vigorous convection (on the geological time scale). Thermal expansion and thermal diffusivity are decreasing with depth while viscosity is likely increasing with depth. The net effect may be that the Rayleigh number for the lower mantle is less than  $\sim 10^7$ .

The linear stability analysis for the onset of convection can also be carried out for a fluid layer heated uniformly from within and cooled from above. The lower boundary is assumed to be insulating, i.e. no heat flows across the boundary. In this case the appropriate Rayleigh number for a fluid layer heated from within is

$$\text{Ra}_H = \frac{\alpha \rho_0^2 g H h^5}{k \eta \kappa}$$

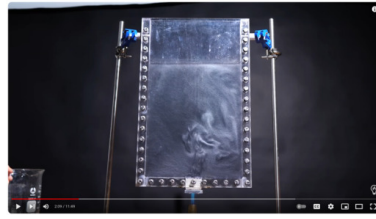
where  $H$  is the rate of internal heat generation per unit mass. For no-slip velocity boundary conditions, the minimum critical Rayleigh number is 2772, and the associated value of  $2\pi h/\lambda$  is 2.63; for free-slip conditions, the minimum  $\text{Ra}_c = 867.8$ , and the associated value of  $2\pi h/\lambda$  is 1.79.

Additional resources:

- D.L. Turcotte and G. Schubert. *Geodynamics, 3rd edition*. Cambridge University Press, 2014. ISBN: 9780521186230, Section 6.19
- D. Bercovici and G. Schubert. *Treatise on geophysics: Mantle dynamics. Vol. 7*. Elsevier, 2007, Section 2.4.4
- G. Schubert, D.L. Turcotte, and P. Olson. *Mantle Convection in the Earth and Planets*. Cambridge University Press, 2001. ISBN: 0-521-70000-0. DOI: 10.1017/CB09780511612879, chapter 7
- Pelletier book, chapter 7.2

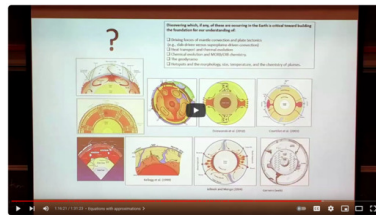
## 11.11 Video resources

- *The bizarre patterns that emerge when you heat ANY fluid* by Steve Mould



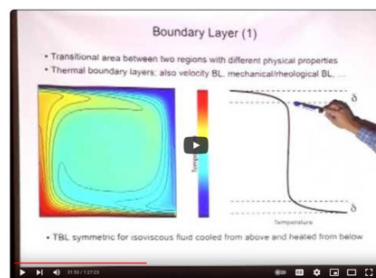
[https://youtu.be/kuLX76g7Fec?si=1DSMkXICnr2\\_yMXv](https://youtu.be/kuLX76g7Fec?si=1DSMkXICnr2_yMXv)

- *Geodynamics 1: Large-Scale Mantle Convection and Numerical Modeling of it* by Allen McNamara at CIDER 2014



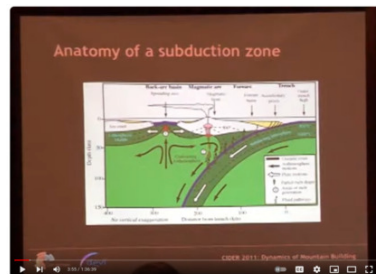
<https://www.youtube.com/watch?v=JyHdFKCIYDE>

- *Geodynamics/Subduction Lab* by Tobias Hoeink



<https://www.youtube.com/watch?v=8cpl5c6lrXA>

- *Subduction* by Dave Stegman



[https://youtu.be/f2sJJ08bqqU?si=qBHUOMdVHo8\\_woNW](https://youtu.be/f2sJJ08bqqU?si=qBHUOMdVHo8_woNW)

## 11.12 Computer practicals

### 11.12.1 Introduction

In parameterized convection models for planetary thermal evolution the heat transport characteristics of the convecting mantle are formulated in a pseudo steady state approximation. This is done by parameterization of the surface heatflux as a function of convective vigor, through the non-dimensional Nusselt number  $\text{Nu}$ . In the convective regime  $\text{Nu}$  is usually expressed in terms of the Rayleigh number  $\text{Ra}$  as  $\text{Nu} \sim C \text{Ra}^\beta$ , where  $C$  is a constant depending on the domain geometry (please read section 1 of Wolstencroft, Davies, and Davies [1368] (2009) for more information, check also Plumley and Julien [1005] (2019) and references therein, also Korenaga [723] (2003)).

In this lab exercise you will investigate the characteristics of steady-state Rayleigh-Bénard convection and determine the relation between the Nusselt and Rayleigh number experimentally, by means of numerical modelling. In particular you will measure the heatflow through the top surface of a 2D model of a convecting layer, as a function of the Rayleigh number, expressed in the temperature contrast across the convecting layer. This is done by a series of modelling experiments where the coupled equations for thermal convection are solved numerically using finite element methods.

The following sections contain descriptions of the numerical model and the experiments to be done.

### 11.12.2 Reminder of the governing model equations

In this computerlab you will perform experiments with numerical solutions of the coupled equations describing thermal convection in an incompressible viscous fluid with infinite Prandtl number<sup>27</sup>.

In what follows, the assumption is made that geological materials can be treated as fluids (with special properties) within the realm of continuum fluid mechanics and under the Stokes hypothesis. A Boussinesq approximation is applied, neglecting density variations in the equations except in the buoyancy term of the momentum conservation equation. We consider two-dimensional problems.

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \rho(T)\vec{g} = \vec{0} \quad (11.135)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (11.136)$$

$$\boldsymbol{\sigma} = -p\mathbf{1} + \boldsymbol{\tau} \quad (11.137)$$

$$\boldsymbol{\tau} = 2\eta\dot{\boldsymbol{\epsilon}}(\vec{v}) \quad (11.138)$$

$$\dot{\boldsymbol{\epsilon}}(\vec{v}) = \frac{1}{2} \left( \vec{\nabla}\vec{v} + (\vec{\nabla}\vec{v})^T \right) \quad (11.139)$$

$$\rho_0 C_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot (k \vec{\nabla} T) \quad (11.140)$$

$$\rho(T) = \rho_0(1 - \alpha(T - T_0)) \quad (11.141)$$

Equation (11.135) is the momentum conservation equation and Eq. (11.136) is the mass conservation equation for incompressible fluids. One can resolve the stress tensor  $\boldsymbol{\sigma}$  into its spherical part  $-p\mathbf{1}$  and its stress deviation  $\boldsymbol{\tau}$  (see Eq. (11.137)), where the deviatoric stress tensor is proportional to the strain rate tensor  $\dot{\boldsymbol{\epsilon}}$  (see Eq.(11.138)) through the dynamic viscosity  $\eta$ . Finally Eq. (11.139) relates the strain rate tensor to the velocity field.

Equations (11.135), (11.136), (11.137), (11.138) and (11.139) all together lead to the following

---

<sup>27</sup>In heat transfer problems, the Prandtl number controls the relative thickness of the momentum and thermal boundary layers. When  $\text{Pr}$  is small, it means that the heat diffuses quickly compared to the velocity (momentum).



form of the Stokes equations:

$$\vec{\nabla} \cdot [\eta(\vec{\nabla}\vec{v} + \vec{\nabla}\vec{v}^T)] - \vec{\nabla}p + \rho\vec{g} = \vec{0} \quad (11.142)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (11.143)$$

Equation (11.142) is an elliptic equation characterized by the fact that changes in buoyancy and constitutive relationships *anywhere* in the domain have an immediate influence on the entire domain.

| symbol                                 | meaning and dimension                                 |
|----------------------------------------|-------------------------------------------------------|
| $\vec{g}$                              | gravity acceleration vector ( $\text{m s}^{-2}$ )     |
| $L_x, L_y$                             | domain size (m)                                       |
| $p$                                    | pressure (Pa)                                         |
| $\boldsymbol{\tau}$                    | deviatoric stress vector (Pa)                         |
| $\vec{v} = (u, v, w)$                  | velocity ( $\text{m s}^{-1}$ )                        |
| $\dot{\boldsymbol{\epsilon}}(\vec{v})$ | strain-rate tensor ( $\text{s}^{-1}$ )                |
| $\eta$                                 | viscosity (Pa s)                                      |
| $\rho, \rho_0$                         | mass density ( $\text{kg m}^{-3}$ )                   |
| $\boldsymbol{\sigma}$                  | stress tensor (Pa)                                    |
| $k$                                    | heat conductivity ( $\text{W m}^{-1} \text{K}^{-1}$ ) |
| $C_p$                                  | heat capacity ( $\text{J K}^{-1}$ )                   |
| $\alpha$                               | thermal expansion ( $\text{K}^{-1}$ )                 |

### 11.12.3 Numerical solution of the equations

Introduced in the late 1950s, the finite element method (FEM) [604, 1430, 1431, 1432] has emerged as one of the most powerful numerical methods so far devised.

A thorough mathematical treatment of the finite element formulation of the equations governing the physics of the system is beyond the scope of this computer practical, and has been exposed in Section 7 and in various textbooks such as [341] or [507].

We wish to study the system at steady state, which means that it no more changes in time. In practice, only the heat transport equation contains a time (derivative) term so we actually wish to solve these three (coupled) equations:

$$\vec{\nabla} \cdot [\eta(\vec{\nabla}\vec{v} + \vec{\nabla}\vec{v}^T)] - \vec{\nabla}p + \rho(T)\vec{g} = \vec{0} \quad (11.144)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (11.145)$$

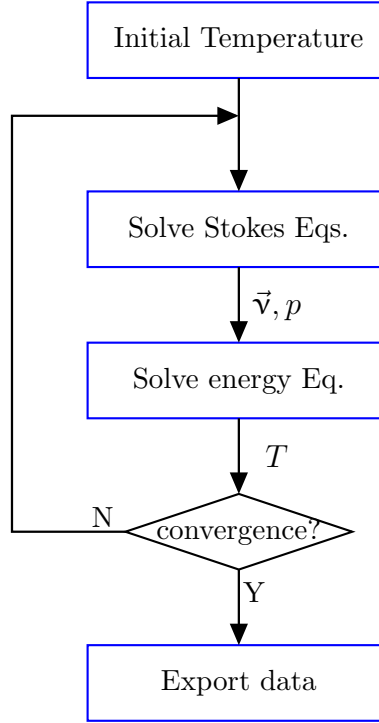
$$\rho_0 C_p \vec{v} \cdot \vec{\nabla} T = k \Delta T \quad (11.146)$$

The main problem is clearly visible in the third one: with respect to the velocity and temperature unknowns this is a nonlinear equation! We therefore have to design a strategy (an algorithm) which will allow us to solve these equations.

One simple approach is as follows: the equations are not solved in a coupled manner. but rather the obtention of a new set of variables ( $\boldsymbol{v}, p, T$ ) is the product of a two-stage process:

1. assume temperature known, solve for velocity (and pressure) field (the first two equations)
2. assume velocity known, solve for temperature (the last equation).

These two steps need to be repeated as long as the system has not converged to a steady solution. The code therefore implements iterations and these iterations stop when either the maximum number of iterations `nstep` is reached or convergence has been reached (i.e. temperature and velocity do not change substantially between two consecutive iterations). The structure of the code is then as follows:



#### 11.12.4 Two-dimensional convection in a unit box

This benchmark deals with the 2-D thermal convection of a fluid of infinite Prandtl number in a rectangular closed cell. In what follows, we will focus on the case 1a, 1b, and 1c experiments as shown in Blankenbach et al. [95] (1989): steady convection with constant viscosity in a square box.

The temperature is fixed to zero on top and to  $\Delta T$  at the bottom, with reflecting symmetry at the sidewalls (i.e.  $\partial_x T = 0$ ) and there are no internal heat sources. Free-slip conditions are implemented on all boundaries.

The Rayleigh number is given by

$$\text{Ra} = \frac{\alpha g_y \rho_0 \Delta T h^3}{\kappa \nu} = \frac{\alpha g_y \Delta T h^3 \rho^2 C_p}{k \eta} \quad (11.147)$$

In what follows, I use the following parameter values:  $L_x = L_y = 1$ ,  $\rho_0 = C_P = k = \eta = 1$ ,  $T_0 = 0$ ,  $\alpha = 10^{-4}$ ,  $g = 10^4 \text{Ra}$ .

The initial temperature field is given by

$$T(x, y) = (1 - y) - 0.01 \cos(\pi x / L_x) \sin(\pi y / L_y) \quad (11.148)$$

The perturbation in the initial temperature fields leads to a perturbation of the density field and sets the fluid in motion. Depending on the initial Rayleigh number, the system ultimately reaches a steady state after some time.

The root mean square of the velocity field in the whole domain is defined as follows:

$$\mathbf{v}_{rms} = \left( \frac{1}{V_\Omega} \int_\Omega |\vec{v}|^2 dV \right)^{1/2} = \left( \frac{1}{L_x L_y} \int_\Omega (u^2 + v^2) dV \right)^{1/2} \quad (11.149)$$

The Nusselt number (i.e. the mean surface temperature gradient over mean bottom temperature) is computed as follows [95]:

$$\text{Nu} = -L_y \frac{\int_0^{L_x} \frac{\partial T}{\partial y}(y = L_y) dx}{\int_0^{L_x} T(y = 0) dx} \quad (11.150)$$

Note that in our case the denominator is equal to  $L_x$  since the temperature at the bottom is prescribed to be 1.

Finally, the steady state root mean square velocity  $\mathbf{v}_{rms}$  and Nusselt number measurements are indicated in the following table alongside those of Blankenbach et al. [95] (1989) and Tackley [1227] (1994). (Note that this benchmark was also carried out and its results published in many other publications [1283, 5, 455, 309, 769] but since they did not provide a complete set of measurement values, they are not included in the table.)

|             |                    | Blankenbach <i>et al.</i> [95] | Tackley [1227] |
|-------------|--------------------|--------------------------------|----------------|
| $Ra = 10^4$ | $\mathbf{v}_{rms}$ | $42.864947 \pm 0.000020$       | 42.775         |
|             | Nu                 | $4.884409 \pm 0.000010$        | 4.878          |
| $Ra = 10^5$ | $\mathbf{v}_{rms}$ | $193.21454 \pm 0.00010$        | 193.11         |
|             | Nu                 | $10.534095 \pm 0.000010$       | 10.531         |
| $Ra = 10^6$ | $\mathbf{v}_{rms}$ | $833.98977 \pm 0.00020$        | 833.55         |
|             | Nu                 | $21.972465 \pm 0.000020$       | 21.998         |

Steady state Nusselt number Nu and  $\mathbf{v}_{rms}$  measurements as reported in the literature.

### 11.12.5 Obtaining the python code

The code is to be downloaded as follows in a terminal:

```
wget https://raw.githubusercontent.com/cedrict/fieldstone/master/python_codes/md/stone_new.py
```

If you do not know what a terminal is or if you are using Windows, simply copy

```
https://raw.githubusercontent.com/cedrict/fieldstone/master/python_codes/md/stone_new.py
```

in the address bar of your web browser, select all, paste it in a file on your computer which you save as `stone.py` in a dedicated folder.

You can run the code in a terminal, in Anaconda, Spyder, etc ...

### 11.12.6 Experiments

To conduct the exercise you can change the following parameters (and run the code until convergence):

- `Lx`: horizontal extent of the domain (do not change `Ly`)
- `Ra_nb`: the Rayleigh number  $Ra$
- `tol_ss`: the steady state detection tolerance
- `nelx,nely`: the number of elements in  $x$  and  $y$  directions
- `nstep`: the maximum number of iterations to reach steady state
- `top_bc_noslip`: flag to switch no slip boundary conditions at top boundary (default is free slip)
- `bot_bc_noslip`: flag to switch no slip boundary conditions at bottom boundary (default is free slip)

Results are written out to `.ascii` files and to `.vtu` files. You can produce plots with python (matplotlib), gnuplot or even excel, as long as these are not pixelated in your report, that they are labeled, captioned, and their axes too. You will find in Appendix [O](#).

Have a thorough look at the code, read **all** instructions, and carry out the following tasks:

1. Determine analytically the expected value of the Nusselt number when there is no convection.
2. Determine the Nusselt number at steady state for a range of Rayleigh numbers, starting from a subcritical value. Produce a plot of  $Nu$  against  $Ra$  using double logarithmic axes. Determine the critical Rayleigh number.
3. The code produces data files containing ‘snapshots’ of the resulting numerical solution of the temperature and velocity fields in a suitable format (`.vtu`) for visualization with graphics program *paraview*. Produce colorplots with *paraview* of the temperature field, for three contrasting Rayleigh number cases, and discuss them.
4. Determine the logarithmic slope or powerlaw index  $\beta$  defined in the introduction.
5. Produce such a  $Nu - Ra$ -plot for various grid resolutions. How can you explain the differences in the results ? Produce a plot of  $Nu$  as a function of the grid spacing.
6. Look at how the  $v_{rms}$  values at steady state depend on  $Ra$ .
7. For three contrasting  $Ra$  values, plot the temperature profiles (data to be found in *T\_profile.ascii*) on a single plot and discuss the obtained figure.
8. Set the number of elements in the horizontal directions to 16. Choose  $Ra = 10^5$  and progressively increase the number of points in the vertical direction. Report on the variation of the  $Nu$  number at steady state as a function of the vertical resolution.
9. Estimate the value of the critical Rayleigh number from your Nusselt number plot and investigate the difference with the value found in Rayleigh’s linear stability analysis for a layer of depth  $h$  and infinite horizontal extent,  $Ra_C = (27/4)\pi^4$  (see Section 11.10.2).
10. Change the initial temperature profile to something more random, repeat some of these experiments. What can you conclude ?
11. Change the top and bottom boundary conditions from free-slip to no-slip. How does this influence  $Ra_c$  ?
12. When convection occurs and a steady state is reached the depth-averaged temperature curves showcases two boundary layers. Measure their thickness as a function of  $Ra$ .
13. Bonus: Explore the effect of the aspect ratio of the domain on  $Ra_c$  and the slope  $\beta$ .
14. Bonus: change the viscosity function so that the viscosity is 1 in the lower half of the domain and  $10^m$  in the upper half with  $m > 1$ . Explore & discuss ...

Relevant sources:

- Chapter 2 and the first part of this chapter present the physical equations in more detail.
- Stone 88 shows examples of mantle-scale convection.
- <https://youtu.be/YIN9Dcq31x0>
- <https://youtu.be/5SPCU1sFGGc>

- <https://youtu.be/ln7QBN0IRTs>
- <https://youtu.be/d4AS1FmdarU>

# Chapter 12

## Manufactured solutions & numerical benchmarks

### 12.1 The method of manufactured solutions

mms.tex

The method of manufactured solutions is a relatively simple way of carrying out code verification. In essence, one postulates a solution for the PDE at hand (as well as the proper boundary conditions), inserts it in the PDE and computes the corresponding source term. The same source term and boundary conditions will then be used in a numerical simulation so that the computed solution can be compared with the (postulated) true analytical solution.

Examples of this approach are to be found in Donea and Huerta [341], Burstedde et al. [189], Bochev, Dohrmann, and Gunzburger [101], Popov, Lobanov, Popov, Popov, and Gerya [1013], Popov, Lobanov, Popov, and Gerya [1012], Lobanov, Popov, Popov, and Gerya [804], Blinova, Makeev, and Popov [97], Thieulot and Bangerth [1260].

#### 12.1.1 The repository

I have created in folder `mms` a template python script for incompressible isoviscous isothermal Stokes flow. All one has to do is to provide the velocity and pressure, the strain rate tensor components and its spatial derivatives.

```
def u_th(x,y):
 return 0

def v_th(x,y):
 return 0

def p_th(x,y):
 return 0

def dpdx_th(x,y):
 return 0

def dpdy_th(x,y):
 return 0

def exx_th(x,y):
 return 0

def exy_th(x,y):
```

```

 return 0

def eyy_th(x,y):
 return 0

def dexxdx(x,y):
 return 0

def dexydx(x,y):
 return 0

def dexydy(x,y):
 return 0

def deyydy(x,y):
 return 0

def bx(x,y):
 return dpdx_th(x,y)-2*dexxdx(x,y)-2*dexydy(x,y)

def by(x,y):
 return dpdy_th(x,y)-2*dexydx(x,y)-2*deyydy(x,y)

def vrms_th():
 return 0

def eta(x,y):
 return 1

```

$$\begin{aligned}
 u(x,y) &= \\
 v(x,y) &= \\
 p(x,y) &=
 \end{aligned}$$

$$\begin{aligned}
 \partial_x u(x,y) &= \\
 \partial_y u(x,y) &= \\
 \partial_x v(x,y) &= \\
 \partial_y v(x,y) &=
 \end{aligned}$$

$$\begin{aligned}
 \dot{\epsilon}_{xx}(x,y) &= \\
 \dot{\epsilon}_{yy}(x,y) &= \\
 \dot{\epsilon}_{xy}(x,y) &=
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial p}{\partial x}(x,y) &= \\
 \frac{\partial p}{\partial y}(x,y) &=
 \end{aligned}$$

$$\begin{aligned}
\partial_x \dot{\epsilon}_{xx}(x, y) &= \\
\partial_x \dot{\epsilon}_{xy}(x, y) &= \\
\partial_y \dot{\epsilon}_{xy}(x, y) &= \\
\partial_y \dot{\epsilon}_{yy}(x, y) &=
\end{aligned}$$

$$v_{rms} = \sqrt{\frac{1}{L_x L_y} \iint (u^2 + v^2) dx dy} =$$

### 12.1.2 Manufactured solution in Donea and Huerta [341] (book)

mms\_dohu03.tex

Taken from [341]. We consider a two-dimensional problem in the square domain  $\Omega = [0, 1] \times [0, 1]$ , which possesses a closed-form analytical solution. The problem consists of determining the incompressible flow velocity field  $\vec{v} = (u, v)$  and the pressure  $p$  such that

$$\vec{\nabla} \cdot (2\eta \dot{\epsilon}(\vec{v})) - \vec{\nabla} p + \vec{b} = \vec{0} \quad \text{in } \Omega \quad (12.1)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad \text{in } \Omega \quad (12.2)$$

$$\vec{v} = \vec{0} \quad \text{on } \Gamma_D \quad (12.3)$$

where the fluid viscosity is taken as  $\eta = 1$ . The components of the body force  $\vec{b}$  are prescribed as

$$\begin{aligned}
b_x &= (12 - 24y)x^4 + (-24 + 48y)x^3 + (-48y + 72y^2 - 48y^3 + 12)x^2 \\
&\quad + (-2 + 24y - 72y^2 + 48y^3)x + 1 - 4y + 12y^2 - 8y^3 \\
b_y &= (8 - 48y + 48y^2)x^3 + (-12 + 72y - 72y^2)x^2 \\
&\quad + (4 - 24y + 48y^2 - 48y^3 + 24y^4)x - 12y^2 + 24y^3 - 12y^4
\end{aligned}$$

With this prescribed body force, the exact solution is

$$\begin{aligned}
u(x, y) &= x^2(1 - x)^2(2y - 6y^2 + 4y^3) \\
&= x^2(1 - x)^2 2y(1 - 3y + 2y^2) \\
&= x^2(1 - x)^2 2y(y - 1)(2y - 1) \\
v(x, y) &= -y^2(1 - y)^2(2x - 6x^2 + 4x^3) \\
&= -y^2(1 - y)^2 2x(1 - 3x + 2x^2) \\
&= -y^2(1 - y)^2 2x(x - 1)(2x - 1) \\
p(x, y) &= x(1 - x) - 1/6
\end{aligned}$$

Note that the pressure obeys  $\int_{\Omega} p \, dV = 0$ . One can turn to the spatial derivatives of the fields:

$$\dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} = (2x - 6x^2 + 4x^3)(2y - 6y^2 + 4y^3) \quad (12.4)$$

$$\dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} = -(2x - 6x^2 + 4x^3)(2y - 6y^2 + 4y^3) \quad (12.5)$$

$$\dot{\epsilon}_{xy} = \frac{1}{2} \left( \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) = \frac{1}{2} (x^2(1 - x)^2(2 - 12y + 12y^2) - y^2(1 - y)^2(2 - 12x + 12x^2)) \quad (12.6)$$

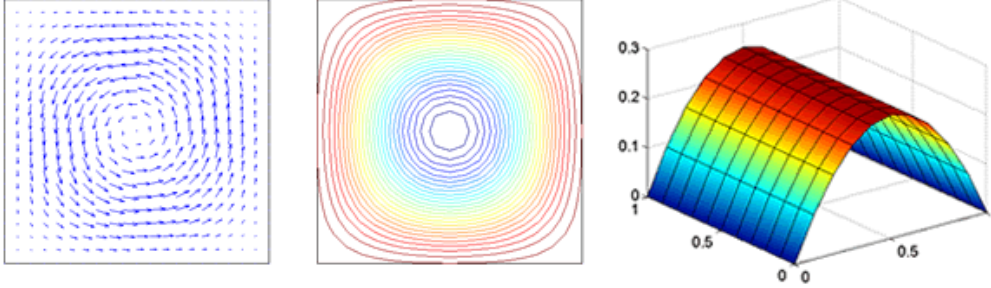


with of course  $\vec{\nabla} \cdot \vec{v} = 0$  and

$$\frac{\partial p}{\partial x} = 1 - 2x \quad (12.7)$$

$$\frac{\partial p}{\partial y} = 0 \quad (12.8)$$

The velocity and pressure fields look like:



[http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible\\_Ex1.htm](http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible_Ex1.htm)

Then the velocity magnitude is given by

$$|\vec{v}|(x, y) = \sqrt{u^2 + v^2} \quad (12.9)$$

$$= \sqrt{[x^2(1-x)^2 2y(1-3y+2y^2)]^2 + [-y^2(1-y)^2 2x(1-3x+2x^2)]^2} \quad (12.10)$$

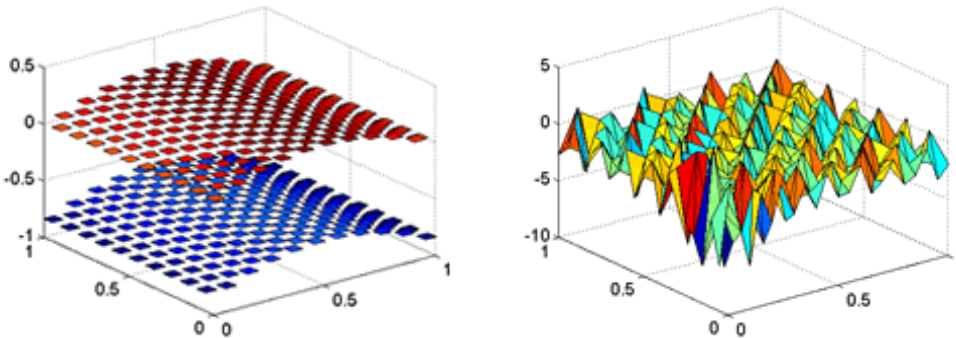
$$= \sqrt{x^4(1-x)^4 4y^2(1-3y+2y^2)^2 + y^4(1-y)^4 4x^2(1-3x+2x^2)^2} \quad (12.11)$$

$$= \sqrt{4x^2y^2 \sqrt{x^2(1-x)^4(1-3y+2y^2)^2 + y^2(1-y)^4(1-3x+2x^2)^2}} \quad (12.12)$$

$$= 2xy \sqrt{x^2(1-x)^4(1-3y+2y^2)^2 + y^2(1-y)^4(1-3x+2x^2)^2} \quad (12.13)$$

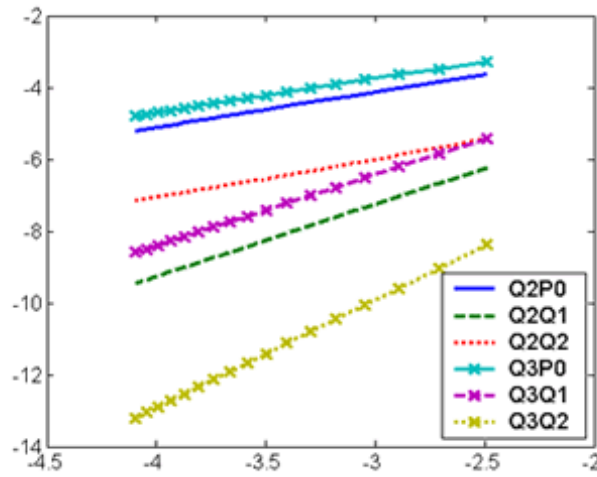
This expression is unfortunately not very useful for later postprocessing...

As shown in [341], If the LBB condition is not satisfied, spurious oscillations spoil the pressure approximation. Figures below show results obtained with a mesh of  $20 \times 20$   $Q_1 \times P_0$  (left) and  $P_1 \times P_1$  (right) elements:



[http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible\\_Ex1.htm](http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible_Ex1.htm) ]]

Taking into account that the proposed problem has got analytical solution, it is easy to analyze convergence of the different pairs of elements:



[http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible\\_Ex1.htm](http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible_Ex1.htm)

One can also compute the stress components:

$$\sigma_{xx} = 2x^2(2x-2)(4y^3-6y^2+2y) + 4x(-x+1)^2(4y^3-6y^2+2y) - x(-x+1) + 1/6$$

$$\sigma_{xy} = x^2(-x+1)^2(12y^2-12y+2) - y^2(-y+1)^2(12x^2-12x+2)$$

$$\sigma_{yy} = -x(-x+1) - 2y^2(2y-2)(4x^3-6x^2+2x) - 4y(-y+1)^2(4x^3-6x^2+2x) + 1/6$$

All the necessary functions to do this benchmark are in `mms/dh.py`:

```
functions for the Donea & Huerta benchmark (dh)

def u_th(x,y):
 return x**2*(1.-x)**2*(2*y-6*y**2+4*y**3)

def v_th(x,y):
 return -y**2*(1.-y)**2*(2*x-6*x**2+4*x**3)

def p_th(x,y):
 return x*(1-x) - 1./6.

def dpdx_th(x,y):
 return 1.-2.*x

def dpdy_th(x,y):
 return 0.

def exx_th(x,y):
 return x**2*(2*x-2)*(4*y**3-6*y**2+2*y)+2*x*(-x+1)**2*(4*y**3-6*y**2+2*y)

def exy_th(x,y):
 return (x**2*(-x+1)**2*(12*y**2-12*y+2)-y**2*(-y+1)**2*(12*x**2-12*x+2))/2

def eyy_th(x,y):
 return -exx_th(x,y)

def bx(x,y):
 return ((12.-24.*y)*x**4+(-24.+48.*y)*x**3*x +
 (-48.*y+72.*y*y-48.*y*y*y+12.)*x**2*x +
 (-2.+24.*y-72.*y*y+48.*y*y*y)*x +
 1.-4.*y+12.*y*y-8.*y*y*y)

def by(x,y):
```

```

return ((8. - 48.*y + 48.*y*y) * x * x * x +
 (-12. + 72.*y - 72.*y*y) * x * x +
 (4. - 24.*y + 48.*y*y - 48.*y*y*y + 24.*y**4) * x -
 12.*y*y + 24.*y*y*y - 12.*y**4)

```

This benchmark is implemented in ASPECT [44] and in STONE 1 and many more. We have

$$\int_0^1 \int_0^1 u^2 dx dy = \int_0^1 \int_0^1 (x^2(1-x)^2(2y-6y^2+4y^3))^2 dx dy = \frac{1}{33075}$$

$$\int_0^1 \int_0^1 v^2 dx dy = \int_0^1 \int_0^1 (-y^2(1-y)^2(2x-6x^2+4x^3))^2 dx dy = \frac{1}{33075}$$

so the root mean square velocity is

$$v_{rms} = \sqrt{\frac{1}{L_x L_y} \int_0^1 \int_0^1 (u^2 + v^2) dx dy} \simeq 0.00777615791$$

We can also look at depth averages. The vertical depth average of the horizontal component of the velocity is given by

$$\begin{aligned} \langle u \rangle(y) &= \frac{1}{L_x} \int_0^{L_x} u(x, y) dx \\ &= \int_0^{L_x} x^2(1-x)^2 2y(1-3y+2y^2) dx \\ &= \left( \int_0^1 x^2(1-x)^2 dx \right) 2y(1-3y+2y^2) \\ &= \frac{1}{30} 2y(1-3y+2y^2) \end{aligned} \tag{12.14}$$

Likewise, the vertical depth average of the vertical component of the velocity is given by

$$\begin{aligned} \langle v \rangle(y) &= \frac{1}{L_x} \int_0^{L_x} v(x, y) dx \\ &= - \int_0^1 y^2(1-y)^2 2x(1-3x+2x^2) dx \\ &= -y^2(1-y)^2 \left( \int_0^1 2x(1-3x+2x^2) dx \right) \\ &= 0 \end{aligned} \tag{12.15}$$

Unfortunately we have seen in Eq.(12.13) that the velocity magnitude is a rather complex function and we won't be able to compute a depth average analytically.

### 12.1.3 Manufactured solution in Dohrmann and Bochev [336] (2004)

mms.dobo.tex

Taken from Dohrmann & Bochev (2004,2006) [336, 101]. This benchmark is also used in Worthen *et al.* [1369] and Lamichhane *et al.* [741]. It is for a unit square with  $\nu = \eta/\rho = 1$  and the smooth exact solution is

$$u(x, y) = x + x^2 - 2xy + x^3 - 3xy^2 + x^2y \tag{12.16}$$

$$v(x, y) = -y - 2xy + y^2 - 3x^2y + y^3 - xy^2 \tag{12.17}$$

$$p(x, y) = xy + x + y + x^3y^2 - 4/3 \tag{12.18}$$

Note that the pressure field is such that  $\int_{\Omega} p \, dV = 0$ . The gradient components of the velocity and pressure fields are given by:

$$\begin{aligned}\frac{\partial u}{\partial x} &= 3x^2 + 2x(y+1) - 3y^2 - 2y + 1 \\ \frac{\partial u}{\partial y} &= x(x - 6y - 2) \\ \frac{\partial v}{\partial x} &= -y(6x + y + 2) \\ \frac{\partial v}{\partial y} &= -3x^2 - 2x(y+1) + 3y^2 + 2y - 1 \\ \frac{\partial p}{\partial x} &= 3x^2y^2 + y + 1 \\ \frac{\partial p}{\partial y} &= 2x^3y + x + 1\end{aligned}$$

so that the strain rate tensor components are

$$\dot{\epsilon}_{xx} = 3x^2 + 2x(y+1) - 3y^2 - 2y + 1 \quad (12.19)$$

$$\dot{\epsilon}_{yy} = -3x^2 - 2x(y+1) + 3y^2 + 2y - 1 \quad (12.20)$$

$$\dot{\epsilon}_{xy} = \frac{1}{2}[x(x - 6y - 2) - y(6x + y + 2)] \quad (12.21)$$

$$= \frac{1}{2}(x^2 - y^2 - 12xy - 2x - 2y) \quad (12.22)$$

We have

$$\frac{\partial \dot{\epsilon}_{xx}}{\partial x} = 2(3x + y + 1) \quad (12.23)$$

$$\frac{\partial \dot{\epsilon}_{xy}}{\partial y} = -6x - y - 1 \quad (12.24)$$

$$\frac{\partial \dot{\epsilon}_{xy}}{\partial x} = -6y + x - 1 \quad (12.25)$$

$$\frac{\partial \dot{\epsilon}_{yy}}{\partial y} = 2(3y - x + 1) \quad (12.26)$$

so that the corresponding body force is given by:

$$\begin{aligned}b_x &= \frac{\partial p}{\partial x} - 2\frac{\partial \dot{\epsilon}_{xx}}{\partial x} - 2\frac{\partial \dot{\epsilon}_{xy}}{\partial y} \\ &= 3x^2y^2 + y + 1 - 4(3x + y + 1) - 2(-6x - y - 1) \\ &= 3x^2y^2 - y - 1\end{aligned} \quad (12.27)$$

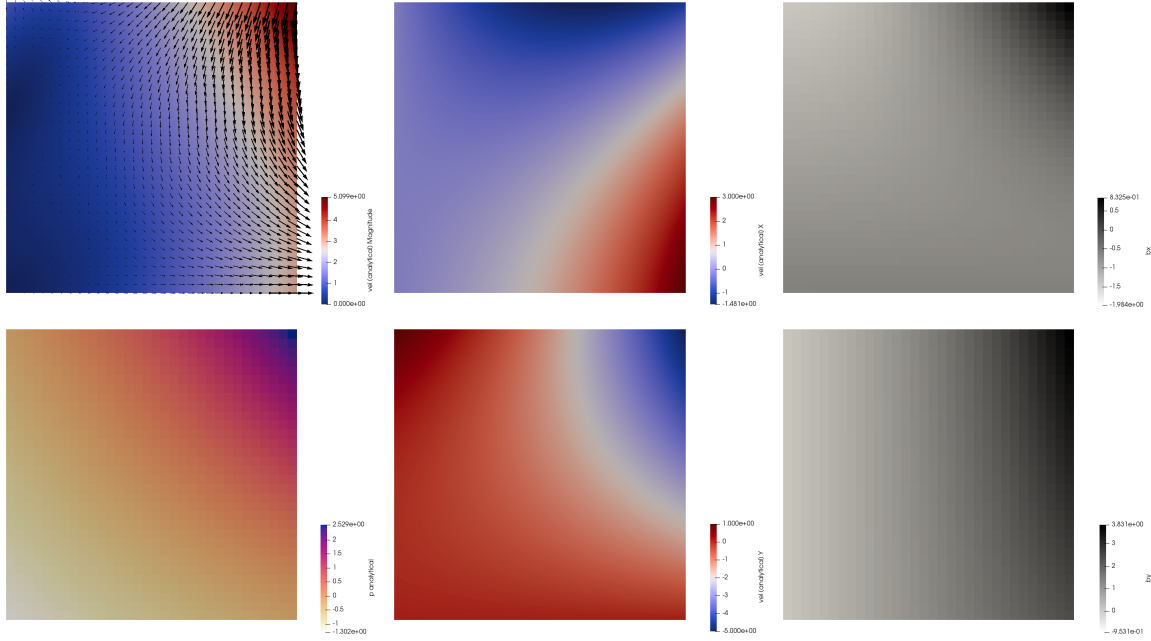
$$\begin{aligned}b_y &= \frac{\partial p}{\partial y} - 2\frac{\partial \dot{\epsilon}_{xy}}{\partial x} - 2\frac{\partial \dot{\epsilon}_{yy}}{\partial y} \\ &= 2x^3y + x + 1 - 2(-6y + x - 1) - 4(3y - x + 1) \\ &= 2x^3y + 3x - 1\end{aligned} \quad (12.28)$$

Finally we have

$$\begin{aligned}\int_0^1 \int_0^1 u^2 dx dy &= \int_0^1 \int_0^1 (x + x^2 - 2xy + x^3 - 3xy^2 + x^2y)^2 dx dy = \frac{401}{504} \\ \int_0^1 \int_0^1 v^2 dx dy &= \int_0^1 \int_0^1 (-y - 2xy + y^2 - 3x^2y + y^3 - xy^2)^2 dx dy = \frac{5911}{2520}\end{aligned}$$

so the root mean square velocity is

$$v_{rms} = \sqrt{\frac{1}{L_x L_y} \int_0^1 \int_0^1 (u^2 + v^2) dx dy} = \sqrt{\frac{401 \cdot 5 + 5911}{2520}} \simeq 1.77236278...$$



#### 12.1.4 Analytical benchmark III - "DB3D"

This benchmark begins by postulating a polynomial solution to the 3D Stokes equation [336]:

$$\vec{v} = \begin{pmatrix} x + x^2 + xy + x^3y \\ y + xy + y^2 + x^2y^2 \\ -2z - 3xz - 3yz - 5x^2yz \end{pmatrix} \quad (12.29)$$

and

$$p = xyz + x^3y^3z - 5/32 \quad (12.30)$$

While it is then trivial to verify that this velocity field is divergence-free (see here under), the corresponding body force of the Stokes equation can be computed by inserting this solution into the momentum equation with a given viscosity  $\eta(x, y, z)$  (constant or position/velocity/strain rate dependent). The domain is a unit cube and velocity boundary conditions simply use Eq. (12.29). Note that the pressure fulfils

$$\int_{\Omega} p(x, y, z) dV = 0.$$

Following [189], the viscosity is given by the smoothly varying function

$$\eta(x, y, z) = \exp(1 - \beta(x(1 - x) + y(1 - y) + z(1 - z))) \quad (12.31)$$

Choosing  $\beta = 0$  yields a constant velocity  $\eta = e^1$  (and greatly simplifies the right-hand side). One can easily show that the ratio of viscosities  $\eta^*$  in the system follows  $\eta^* = \exp(-3\beta/4)$  so that choosing  $\beta = 10$  yields  $\eta^* \simeq 1808$  and  $\beta = 20$  yields  $\eta^* \simeq 3.269 \times 10^6$ .

The exact form of the rhs is carried out in Stone ??.

Let us now compute the root mean square velocity:

$$\int_{\Omega} u^2 dx dy dz = \int_0^{+1} \int_0^{+1} \int_0^{+1} (x + x^2 + xy + x^3 y)^2 dx dy dz = 2867/1260 \quad (12.32)$$

$$\int_{\Omega} v^2 dx dy dz = \int_0^{+1} \int_0^{+1} \int_0^{+1} (y + xy + y^2 + x^2 y^2)^2 dx dy dz = 3947/1800 \quad (12.33)$$

$$\int_{\Omega} w^2 dx dy dz = \int_0^{+1} \int_0^{+1} \int_0^{+1} (-2z - 3xz - 3yz - 5x^2 yz)^2 dx dy dz = 463/36 \quad (12.34)$$

then

$$\mathbf{v}_{rms} = \sqrt{2867/1260 + 3947/1800 + 463/36} \simeq 4.1628459$$

### 12.1.5 Analytical benchmark IV - "Bercovier & Engelman"

From [79]. The two-dimensional domain is a unit square. The body forces are:

$$\begin{aligned} f_x &= 128[x^2(x-1)^2 12(2y-1) + 2(y-1)(2y-1)y(12x^2 - 12x + 2)] \\ f_y &= 128[y^2(y-1)^2 12(2x-1) + 2(x-1)(2x-1)y(12y^2 - 12y + 2)] \end{aligned} \quad (12.35)$$

The solution is

$$\begin{aligned} u &= -256x^2(x-1)^2y(y-1)(2y-1) \\ v &= 256y^2(y-1)^2x(x-1)(2x-1) \\ p &= 0 \end{aligned} \quad (12.36)$$

$$du/dx = 512(1-2x)(-1+x)x(-1+y)y(-1+2y) \quad (12.37)$$

$$du/dy = -256(-1+x)^2x^2(1-6y+6y^2) \quad (12.38)$$

$$dv/dx = 256y^2(y-1)^2x(x-1)(2x-1) \quad (12.39)$$

$$dv/dy = -512(-1+x)x(1-2x)(-1+y)y(-1+2y) \quad (12.40)$$

$$(12.41)$$

and we can easily verify that  $\vec{\nabla} \cdot \vec{\mathbf{v}} = du/dx + dv/dy = 0$ .

CHECK RHS !

Another choice with a non-zero pressure:

$$\begin{aligned} f_x &= 128[x^2(x-1)^2 12(2y-1) + 2(y-1)(2y-1)y(12x^2 - 12x + 2)] + y - 1/2 \\ f_y &= 128[y^2(y-1)^2 12(2x-1) + 2(x-1)(2x-1)y(12y^2 - 12y + 2)] + x - 1/2 \end{aligned} \quad (12.42)$$

The solution is

$$\begin{aligned} u &= -256x^2(x-1)^2y(y-1)(2y-1) \\ v &= 256y^2(y-1)^2x(x-1)(2x-1) \\ p &= (x-1/2)(y-1/2) \end{aligned} \quad (12.43)$$

### 12.1.6 Analytical benchmark VI - "Ilinca & Pelletier"

This is taken from [620].

Let us consider the Poiseuille flow of a Newtonian fluid. The channel has isothermal flat walls located at  $y = \pm h$ . The velocity distribution is parabolic:

$$u = u_0 \left(1 - \frac{y^2}{h^2}\right) \quad v = 0$$

where  $u_0$  is the maximum velocity. The (steady state) temperature field is the solution of the advection-diffusion equation:

$$\rho C_p \vec{v} \cdot \vec{\nabla} T = k \Delta T + \Phi$$

where  $\Phi$  is the dissipation function given by

$$\Phi = \eta \left[ 2 \left( \frac{\partial u}{\partial x} \right)^2 + 2 \left( \frac{\partial v}{\partial y} \right)^2 + \left( \frac{\partial v}{\partial x} + \frac{\partial u}{\partial y} \right)^2 \right] = \eta \left( \frac{\partial u}{\partial y} \right)^2 = 4\eta \frac{u_0^2 y^2}{h^4}$$

We logically assume that  $T = T(y)$  so that  $\partial T / \partial x = 0$  and  $\vec{v} \cdot \vec{\nabla} T = 0$ . We then have to solve:

$$k \frac{\partial^2 T}{\partial y^2} + 4\eta \frac{u_0^2 y^2}{h^4} = 0$$

We can integrate twice and use the boundary conditions  $T(y = \pm h) = T_0$  to arrive at:

$$T(y) = T_0 + \frac{1}{3} \frac{\eta u_0^2}{k} \left[ 1 - \left( \frac{y}{h} \right)^4 \right]$$

with a maximum temperature

$$T_M = T(y = 0) = T_0 + \frac{1}{3} \frac{\eta u_0^2}{k}$$

### 12.1.7 Analytical benchmark VII - "grooves"

mms\_grooves.tex

This benchmark was designed by Dave May. The velocity and pressure fields are given by

$$\begin{aligned} u(x, y) &= x^3 y + x^2 + xy + x \\ v(x, y) &= -\frac{3}{2} x^2 y^2 - 2xy - \frac{1}{2} y^2 - y \\ p(x, y) &= x^2 y^2 + xy + 5 + p_0 \end{aligned} \tag{12.44}$$

where  $p_0$  is a constant to be determined based on the type of pressure normalisation. The viscosity is chosen to be

$$\eta(x, y) = -\sin(p) + 1 + \epsilon = -\sin(x^2 y^2 + xy + 5) + 1 + \epsilon \tag{12.45}$$

where  $\epsilon$  actually controls the viscosity contrast. Note that inserting the polynomial expression of the pressure inside the viscosity expression makes the problem linear. We have

$$\begin{aligned} \dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} &= 3x^2 y + 2x + y + 1 \\ \dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} &= -3x^2 y - 2x - y - 1 \\ \dot{\epsilon}_{xy} = \frac{1}{2} \left( \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) &= \frac{1}{2} (x^3 + x - 3xy^2 - 2y) \end{aligned} \tag{12.46}$$

and we can verify that the velocity field is incompressible since  $\vec{\nabla} \cdot \vec{v} = \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} = 0$ . The pressure gradient is given by

$$\begin{aligned}\frac{\partial p}{\partial x} &= 2xy^2 + y \\ \frac{\partial p}{\partial y} &= 2x^2y + x\end{aligned}$$

The right hand side term of the Stokes equation is such that

$$\begin{aligned}-\frac{\partial p}{\partial x} + \frac{\partial s_{xx}}{\partial x} + \frac{\partial s_{yx}}{\partial y} + f_x &= 0 \\ -\frac{\partial p}{\partial y} + \frac{\partial s_{xy}}{\partial x} + \frac{\partial s_{yy}}{\partial y} + f_y &= 0\end{aligned}\tag{12.47}$$

with

$$\begin{aligned}\frac{\partial s_{xx}}{\partial x} &= \frac{\partial(2\eta\dot{\epsilon}_{xx})}{\partial x} = 2\eta\frac{\partial\dot{\epsilon}_{xx}}{\partial x} + 2\frac{\partial\eta}{\partial x}\dot{\epsilon}_{xx} \\ \frac{\partial s_{zx}}{\partial z} &= \frac{\partial(2\eta\dot{\epsilon}_{zx})}{\partial z} = 2\eta\frac{\partial\dot{\epsilon}_{zx}}{\partial z} + 2\frac{\partial\eta}{\partial z}\dot{\epsilon}_{zx} \\ \frac{\partial s_{xz}}{\partial x} &= \frac{\partial(2\eta\dot{\epsilon}_{xz})}{\partial x} = 2\eta\frac{\partial\dot{\epsilon}_{xz}}{\partial x} + 2\frac{\partial\eta}{\partial x}\dot{\epsilon}_{xz} \\ \frac{\partial s_{zz}}{\partial z} &= \frac{\partial(2\eta\dot{\epsilon}_{zz})}{\partial z} = 2\eta\frac{\partial\dot{\epsilon}_{zz}}{\partial z} + 2\frac{\partial\eta}{\partial z}\dot{\epsilon}_{zz} \\ \frac{\partial\eta}{\partial x} &= -z(2xz + 1)\cos(x^2z^2 + xz + 5) \\ \frac{\partial\eta}{\partial z} &= -x(2xz + 1)\cos(x^2z^2 + xz + 5) \\ \frac{\partial\dot{\epsilon}_{xx}}{\partial x} &= 6xz + 2 \\ \frac{\partial\dot{\epsilon}_{zx}}{\partial z} &= -3xz - 1 \\ \frac{\partial\dot{\epsilon}_{xz}}{\partial x} &= \frac{1}{2}(3x^2 + 1 - 3z^2) \\ \frac{\partial\dot{\epsilon}_{zz}}{\partial z} &= -3x^2 - 1\end{aligned}$$

Velocity boundary conditions are prescribed on all four boundaries so that the pressure is known up to a constant (the pressure solution has a nullspace), and the  $p_0$  constant can be determined by requiring that

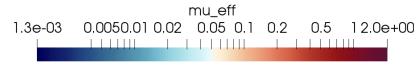
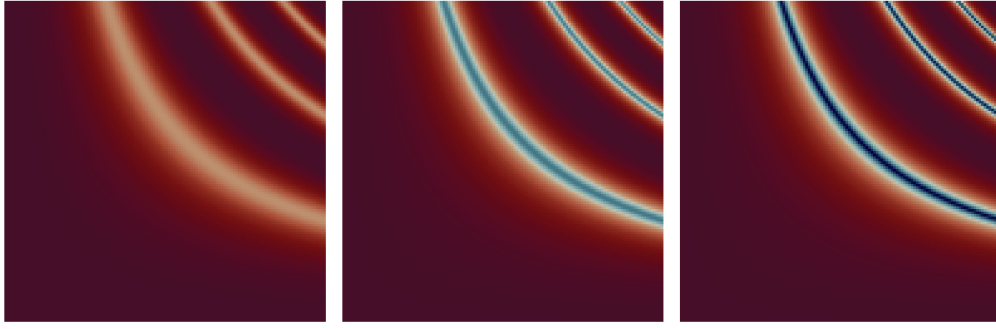
$$\int_0^L \int_0^L p(x, y) dx dy = \int_0^L \int_0^L (x^2y^2 + xy + 5) dx dy + \int_0^L \int_0^L p_0 dx dy = \int_0^L \int_0^L (x^2y^2 + xy + 5) dx dy + p_0 L^2 = 0$$

where  $L$  is the size of the square domain. Then

$$p_0 = -\frac{1}{L^2} \int_0^L \int_0^L (x^2y^2 + xy + 5) dx dy = -\frac{L^4}{9} - \frac{L^2}{4} - 5$$

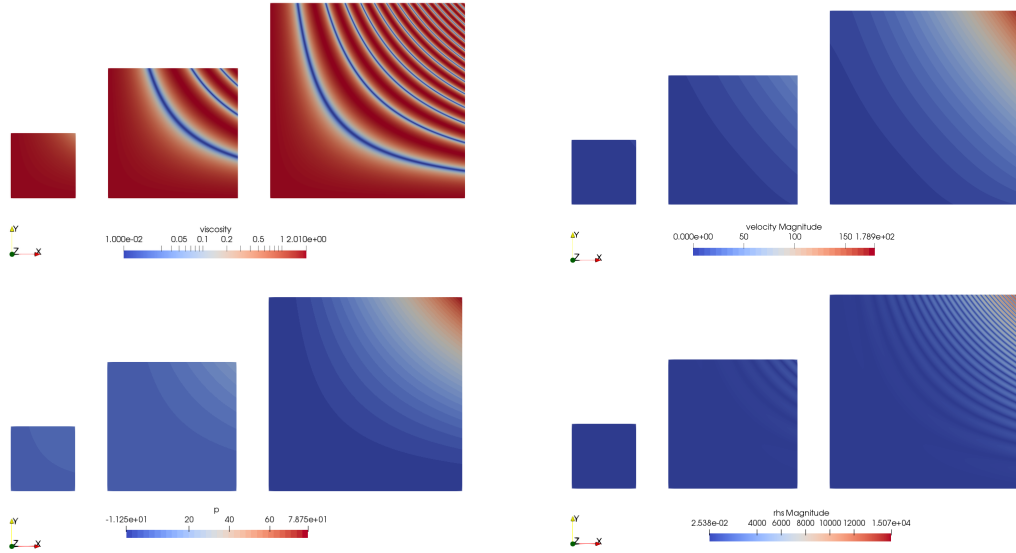
As seen in the following figure, the value of  $\epsilon$  controls the viscosity field amplitude. This is simply explained by the fact that when the sin term of the viscosity takes value 1, the viscosity is then equal to  $\epsilon$ .





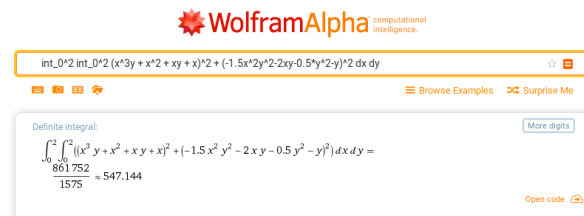
Domain size 2x2 with  $\epsilon = 0.1, 0.01, 0.001$

Another interesting aspect of this benchmark is the fact that increasing the domain size adds complexity to it as it increases the number of low viscosity zones and the spacing between them also decreases:



Three different domain sizes (1x1, 2x2, 3x3) with  $\epsilon = 0.001$ .

Finally, because the analytical expression for both components of the velocity is a polynomial, we can also compute the root mean square velocity exactly. For instance, for a 2x2 domain:

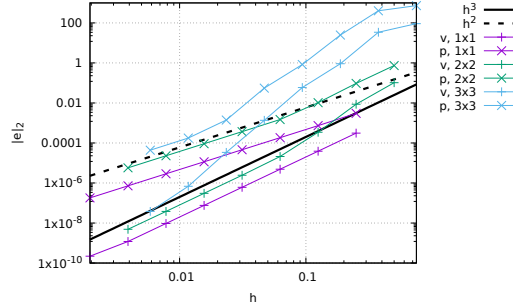


and we end up with (for  $L = 2$ )

$$v_{rms} = \sqrt{\frac{1}{L^2} \frac{861752}{1575}} = \sqrt{\frac{215438}{1575}} \simeq 11.6955560683$$

I have added this benchmark to ASPECT . The velocity and pressure errors (in the  $L_2$  norm) are measured for  $L = 1, 2, 3$ , levels 3 to 9 (resolutions  $8 \times 8$  to  $512 \times 512$ ) and  $\epsilon = 10^{-1}, 10^{-2}, 10^{-3}$ .

The figure below shows the velocity and pressure error convergence as a function of the mesh size for  $\epsilon = 0.1$  (results are identical for the other two  $\epsilon$  values). The expected convergence rates (cubic convergence for velocity and quadratic for pressure) are recovered for the  $1 \times 1$  domain at all resolutions. These rates are recovered for the  $2 \times 2$  domain for resolutions above level 6. We see that the multitude of low viscosity bands in the upper right corner of the  $3 \times 3$  domain will require a refinement level superior to 9 to recover the optimal convergence rates.



Velocity and pressure error convergence as a function of the mesh size  $h$  for 3 domain sizes with  $\epsilon = 0.1$ .

This benchmark is implemented and used in [STONE](#) 112.

### 12.1.8 Analytical benchmark VIII - "Kovasznay"

This flow was published by L.I.G. Kovasznay in 1948 [726]. This paper presents an exact two-dimensional solution of the Navier-Stokes equations with a periodicity in the vertical direction, gives an analytical solution to the steady-state Navier-Stokes equations that is similar which is a flow-field behind a periodic array of cylinders.

$$u(x, y) = 1 - \exp(\lambda x) \cos(2\pi y) \quad (12.48)$$

$$v(x, y) = \frac{\lambda}{2\pi} \exp(\lambda x) \sin(2\pi y) \quad (12.49)$$

$$p(x, y) = \frac{1}{2}(1 - \exp(2\lambda x)) \quad (12.50)$$

$$\lambda = \frac{Re}{2} - \sqrt{\frac{Re^2}{4} + 4\pi^2} \quad (12.51)$$

Following step-55 of deal.II<sup>1</sup> we have to 'cheat' here since we are not solving the non-linear Navier-Stokes equations, but the linear Stokes system without convective term. Therefore, to recreate the exact same solution we move the convective term into the right-hand side.

The analytical solution is prescribed left and right, while free/no (??) slip is prescribed at top and bottom.

Velocity and pressure solution as implemented in step-55:

```
const double pi2 = pi*pi;

u = -exp(x*(-sqrt(25.0 + 4*pi2) + 5.0))*cos(2*y*pi) + 1

v = (1.0L/2.0L)*(-sqrt(25.0 + 4*pi2) + 5.0)*
 exp(x*(-sqrt(25.0 + 4*pi2) + 5.0))*sin(2*y*pi)/pi

p = -1.0L/2.0L*exp(x*(-2*sqrt(25.0 + 4*pi2) + 10.0)) - 2.0*(-6538034.74494422
 + 0.0134758939981709*exp(4*sqrt(25.0 + 4*pi2)))/(-80.0*exp(3*sqrt(25.0 + 4*pi2)))
```

<sup>1</sup>[https://www.dealii.org/current/doxygen/deal.II/step\\_55.html](https://www.dealii.org/current/doxygen/deal.II/step_55.html)

```

+ 16.0*sqrt(25.0 + 4*pi2)*exp(3*sqrt(25.0 + 4*pi2)))
- 1634508.68623606*exp(-3.0*sqrt(25.0 + 4*pi2))/(-10.0 + 2.0*sqrt(25.0 + 4*pi2))
+ (-0.00673794699908547*exp(sqrt(25.0 + 4*pi2))
+ 3269017.37247211*exp(-3*sqrt(25.0 + 4*pi2)))/(-8*sqrt(25.0 + 4*pi2) + 40.0)
+ 0.00336897349954273*exp(1.0*sqrt(25.0 + 4*pi2))/(-10.0 + 2.0*sqrt(25.0 + 4*pi2))

```

while the rhs of the PDE is given by

```

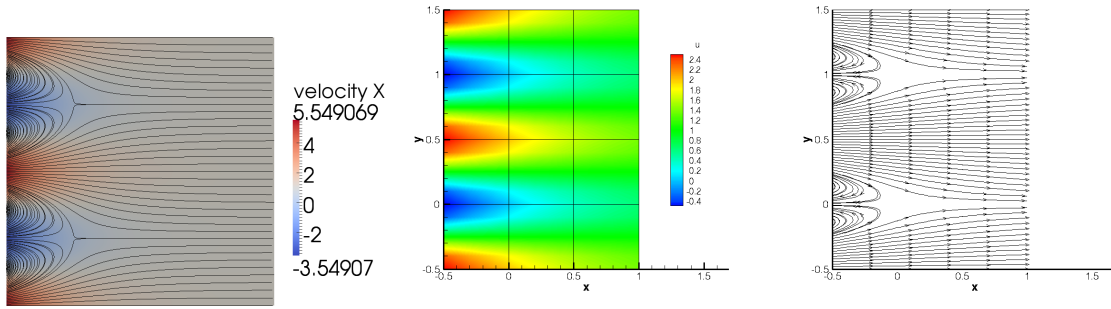
const double pi2 = pi * pi;

values[0] = -1.0L / 2.0L * (-2 * sqrt(25.0 + 4 * pi2) + 10.0) *
 exp(x*(-2*sqrt(25.0 + 4 * pi2) + 10.0)) -
 0.4 * pi2*exp(x * (-sqrt(25.0 + 4 * pi2) + 5.0)) * cos(2 * y * pi) +
 0.1 *pow(-sqrt(25.0 + 4 * pi2) + 5.0, 2) *
 exp(x*(-sqrt(25.0 + 4 * pi2) + 5.0)) * cos(2 * y * pi)

values[1] = 0.2 * pi*(-sqrt(25.0 + 4 * pi2) + 5.0) *
 exp(x*(-sqrt(25.0 + 4 * pi2) + 5.0)) * sin(2 * y * pi) -
 0.05 *pow(-sqrt(25.0 + 4 * pi2) + 5.0, 3) *
 exp(x*(-sqrt(25.0 + 4 * pi2) + 5.0)) * sin(2 * y * pi) / pi

values[2] = 0;

```



Left: solution from Step-55. Right: Solution obtained with NekTar++<sup>2</sup>

This benchmark is carried out in many CFD papers: [269, 118, 936], see also Section 7.4.3 of Hesthaven & Warburton [568].

Find analytical expression for pressure. Compute expression for rhs. Make stone

### 12.1.9 Analytical benchmark IX - "VJ2"

It is presented in [655] and meant to be a peculiar case where the velocity solution is exactly zero. The viscosity is 1, the domain is a unit square, no-slip boundary conditions are prescribed everywhere. The buoyancy force is given by  $\vec{b} = (0, Ra(1 - y + 3y^2))$  where  $Ra > 0$  is a parameter. The flow is incompressible and the analytical pressure solution is given by  $p = Ra(y^3 - y^2/2 + y - 7/12)$ .

### 12.1.10 Manufactured solution in John, Linke, Merdon, Neilan, and Rebholz [655]

mms\_jolm17.tex

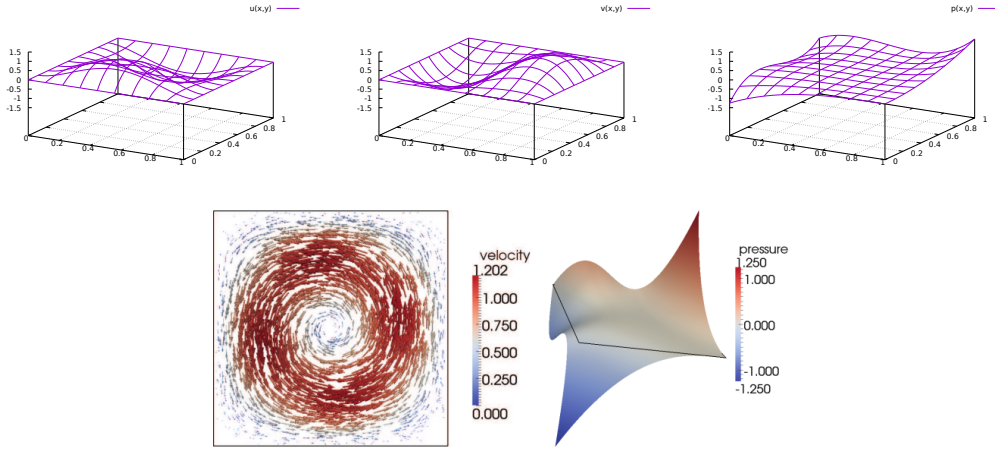
This benchmark comes from John *et al.* [655]. The domain is once again the unit square. The velocity field has the form of a large vortex. Note that velocity field is actually the same velocity field as in the Donea & Huerta benchmark above (albeit multiplied by a factor 100).

<sup>2</sup><http://doc.nektar.info/userguide/4.3.4/user-guidese45.html>

$$u(x, y) = 200x^2(1-x)^2y(1-y)(1-2y) \quad (12.52)$$

$$v(x, y) = -200x(1-x)(1-2x)y^2(1-y)^2 \quad (12.53)$$

$$p(x, y) = 10 \left[ (x - 1/2)^3 y^2 + (1-x)^3 (y - 1/2)^3 \right] \quad (12.54)$$



Taken from John, Linke, Merdon, Neilan, and Rebholz [655] (2017).

$$\dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} = -400(1-x)x(2x-1)(y-1)y(2y-1) \quad (12.55)$$

$$\frac{\partial u}{\partial y} = 200(1-x)^2x^2(6y^2-6y+1) \quad (12.56)$$

$$\frac{\partial v}{\partial x} = -200(6x^2-6x+1)(1-y)^2y^2 \quad (12.57)$$

$$\dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} = 400(x-1)x(2x-1)(1-y)y(2y-1) \quad (12.58)$$

so that

$$\begin{aligned} \dot{\epsilon}_{xy} &= \frac{1}{2} [200(1-x)^2x^2(6y^2-6y+1) - 200(6x^2-6x+1)(1-y)^2y^2] \\ &= 100(1-x)^2x^2(6y^2-6y+1) - 100(6x^2-6x+1)(1-y)^2y^2 \end{aligned} \quad (12.59)$$

Also

$$\begin{aligned} \frac{\partial \dot{\epsilon}_{xx}}{\partial x} &= 400(6x^2-6x+1)y(2y^2-3y+1) \\ \frac{\partial \dot{\epsilon}_{xy}}{\partial x} &= 200(-2x^2(1-x)(6y^2-6y+1) + 2x(1-x)^2(6y^2-6y+1) - 6(2x-1)(1-y)^2y^2) \\ &= 100(-2x^2(1-x)(6y^2-6y+1) + 2x(1-x)^2(6y^2-6y+1) - 6(2x-1)(1-y)^2y^2) \\ \frac{\partial \dot{\epsilon}_{xy}}{\partial y} &= 400(6x^2-6x+1)(1-y)y^2 + 200(1-x)^2x^2(12y-6) - 400(6x^2-6x+1)(1-y)^2y \\ \frac{\partial \dot{\epsilon}_{yy}}{\partial y} &= -400x(2x^2-3x+1)(6y^2-6y+1) \end{aligned} \quad (12.60)$$

$$\frac{\partial p}{\partial x} = 30(x-1/2)^2y^2 - 30(1-x)^2(y-1/2)^3 \quad (12.61)$$

$$\frac{\partial p}{\partial y} = 20(x-1/2)^3y + 30(1-x)^3(y-1/2)^2 \quad (12.62)$$

From  $\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{b} = \vec{0}$  we can obtain the rhs as follows:

$$\begin{aligned}\vec{b} &= -\vec{\nabla} \cdot \boldsymbol{\sigma} \\ &= \vec{\nabla} p - \vec{\nabla} \cdot \mathbf{s} \\ &= \vec{\nabla} p - \vec{\nabla} \cdot (2\eta \dot{\boldsymbol{\epsilon}})\end{aligned}\tag{12.63}$$

Assuming  $\eta = 1$  we arrive at:

$$b_x = \frac{\partial p}{\partial x} - 2 \frac{\partial \dot{\epsilon}_{xx}}{\partial x} - 2 \frac{\partial \dot{\epsilon}_{xy}}{\partial y}\tag{12.64}$$

$$b_y = \frac{\partial p}{\partial y} - 2 \frac{\partial \dot{\epsilon}_{xy}}{\partial x} - 2 \frac{\partial \dot{\epsilon}_{yy}}{\partial y}\tag{12.65}$$

All the necessary functions to do this benchmark are in `mms/jolm17.py`:

```
functions for the John et al (2017) manufactured solution

def u_th(x,y):
 return 200*x**2*(1-x)**2*y*(1-y)*(1-2*y)

def v_th(x,y):
 return -200*x*(1-x)*(1-2*x)*y**2*(1-y)**2

def p_th(x,y):
 return 10*((x-1./2.)**3*y**2+(1-x)**3*(y-1./2.)**3)

def dpdx_th(x,y):
 return 30*(x-1./2.)**2*y**2-30*(1-x)**2*(y-1./2.)**3

def dpdy_th(x,y):
 return 20*(x-1./2.)**3*y + 30*(1-x)**3*(y-1./2.)**2

def exx_th(x,y):
 return -400*(1-x)*x*(2*x-1)*(y-1)*y*(2*y-1)

def exy_th(x,y):
 return 100*(1-x)**2*x**2*(6*y**2-6*y+1)-100*(6*x**2-6*x+1)*(1-y)**2*y**2

def eyy_th(x,y):
 return 400*(x-1)*x*(2*x-1)*(1-y)*y*(2*y-1)

def dexxdx(x,y):
 return 400*(6*x**2-6*x+1)*y*(2*y**2-3*y+1)

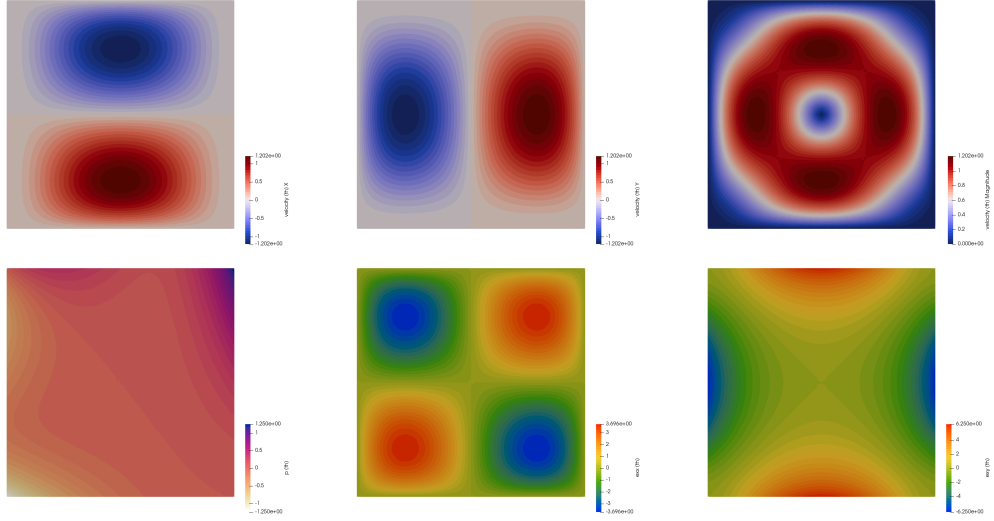
def dexydx(x,y):
 return 100*(-2*x**2*(1-x)*(6*y**2-6*y+1) + 2*x*(1-x)**2*(6*y**2-6*y+1) -6*(2*x-1)
 *(1-y)**2*y**2)

def dexydy(x,y):
 return 200*(6*x**2-6*x+1)*(1-y)*y**2 + 100*(1-x)**2*x**2*(12*y-6) -200*(6*x**2-6*
 x+1)*(1-y)**2*y

def deyydy(x,y):
 return -400*x*(2*x**2-3*x+1)*(6*y**2-6*y+1)

def bx(x,y):
 return dpdx_th(x,y)-2*dexxdx(x,y)-2*dexydy(x,y)

def by(x,y):
 return dpdy_th(x,y)-2*dexydx(x,y)-2*deyydy(x,y)
```



### 12.1.11 Manufactured solution in Lamichhane [741] (2017)

Taken from Lamichhane (2017) [741].

$$u(x, y) = -2x^2y(2y - 1)(x - 1)^2(y - 1) = a(x)b(y) \quad (12.66)$$

$$v(x, y) = 2xy^2(2x - 1)(x - 1)(y - 1)^2 = c(x)d(y) \quad (12.67)$$

$$p(x, y) = x(1 - x)(1 - 2y) \quad (12.68)$$

with

$$a(x) = -2x^2(x - 1)^2 \quad (12.69)$$

$$b(y) = y(2y - 1)(y - 1) \quad (12.70)$$

$$c(x) = x(2x - 1)(x - 1) \quad (12.71)$$

$$d(y) = 2y^2(y - 1)^2 \quad (12.72)$$

and

$$a'(x) = -4x(2x^2 - 3x + 1) \quad (12.73)$$

$$a''(x) = -4(6x^2 - 6x + 1) \quad (12.74)$$

$$b'(y) = 6y^2 - 6y + 1 \quad (12.75)$$

$$b''(y) = 12y - 6 \quad (12.76)$$

$$c'(x) = 6x^2 - 6x + 1 \quad (12.77)$$

$$c''(x) = 12x - 6 \quad (12.78)$$

$$d'(y) = 4y(2y^2 - 3y + 1) \quad (12.79)$$

$$d''(y) = 4(6y^2 - 6y + 1) \quad (12.80)$$

It then follows:

$$L_{xx} = \frac{\partial u}{\partial x} = a'(x)b(y) \quad (12.81)$$

$$L_{xy} = \frac{\partial v}{\partial x} = c'(x)d(y) \quad (12.82)$$

$$L_{yx} = \frac{\partial u}{\partial y} = a(x)b'(y) \quad (12.83)$$

$$L_{yy} = \frac{\partial v}{\partial y} = c(x)d'(y) \quad (12.84)$$

and

$$\dot{\varepsilon}_{xx} = L_{xx} = a'(x)b(y) \quad (12.85)$$

$$\dot{\varepsilon}_{yy} = L_{yy} = c(x)d'(y) \quad (12.86)$$

$$\dot{\varepsilon}_{xy} = \frac{1}{2}(L_{xy} + L_{yx}) \quad (12.87)$$

$$= \frac{1}{2}(a(x)b'(y) + c'(x)d(y)) \quad (12.88)$$

We easily verify that  $\dot{\varepsilon}_{xx} + \dot{\varepsilon}_{yy} = 0$ .

$$\frac{\partial p}{\partial x} = (1 - 2x)(1 - 2y) \quad (12.89)$$

$$\frac{\partial p}{\partial y} = -2x(1 - x) \quad (12.90)$$

We will also need:

$$\frac{\partial \dot{\varepsilon}_{xx}}{\partial x} = a''(x)b(y) \quad (12.91)$$

$$\frac{\partial \dot{\varepsilon}_{xy}}{\partial y} = \frac{1}{2}(a(x)b''(y) + c'(x)d'(y)) \quad (12.92)$$

$$\frac{\partial \dot{\varepsilon}_{xy}}{\partial x} = \frac{1}{2}(a'(x)b'(y) + c''(x)d(y)) \quad (12.93)$$

$$\frac{\partial \dot{\varepsilon}_{yy}}{\partial x} = c(x)d''(y) \quad (12.94)$$

Assuming  $\eta = 1$  we arrive at:

$$b_x = \frac{\partial p}{\partial x} - 2\frac{\partial \dot{\varepsilon}_{xx}}{\partial x} - 2\frac{\partial \dot{\varepsilon}_{xy}}{\partial y} \quad (12.95)$$

$$= (1 - 2x)(1 - 2y) - 2a''(x)b(y) - (a(x)b''(y) + c'(x)d'(y)) \quad (12.96)$$

$$b_y = \frac{\partial p}{\partial y} - 2\frac{\partial \dot{\varepsilon}_{xy}}{\partial x} - 2\frac{\partial \dot{\varepsilon}_{yy}}{\partial y} \quad (12.97)$$

$$= -2x(1 - x) - (a'(x)b'(y) + c''(x)d(y)) - 2c(x)d''(y) \quad (12.98)$$

We have (thank you WolframAlpha<sup>3</sup>):

$$\int_0^1 \int_0^1 u^2 dx dy = \int_0^1 \int_0^1 (-2x^2y(2y - 1)(x - 1)^2(y - 1))^2 dx dy = \frac{1}{33075}$$

$$\int_0^1 \int_0^1 v^2 dx dy = \int_0^1 \int_0^1 (2xy^2(2x - 1)(x - 1)(y - 1)^2)^2 dx dy = \frac{1}{33075}$$

so that

$$v_{rms} = \sqrt{\frac{1}{L_x L_y} \frac{2}{33075}} \simeq 0.00777615791$$

Finally we can verify that the pressure has zero average:

$$\int_0^1 \int_0^1 p(x, y) dx dy = \int_0^1 \int_0^1 x(1 - x)(1 - 2y) dx dy = \int_0^1 x(1 - x) dx \int_0^1 (1 - 2y) dy = 0$$

---

<sup>3</sup><https://www.wolframalpha.com/>

### 12.1.12 Manufactured solution in Mu and Ye [911]

We first consider the Stokes equations with homogeneous boundary condition defined on the unit square. The exact solutions are given by

$$\begin{aligned}u(x, y) &= 2\pi \sin^2(\pi x) \cos(\pi y) \sin(\pi y) \\v(x, y) &= -2\pi \sin(\pi x) \cos(\pi x) \sin^2(\pi y) \\p(x, y) &= \cos(\pi x) \cos(\pi y)\end{aligned}\tag{12.99}$$

Then, since  $\cos' = -\sin$ :

$$\begin{aligned}\partial_x p &= -\pi \sin(\pi x) \cos(\pi y) \\\partial_y p &= -\pi \cos(\pi x) \sin(\pi y)\end{aligned}$$

and

$$\begin{aligned}\partial_x u &= 4\pi^2 \sin(\pi x) \cos(\pi x) \cos(\pi y) \sin(\pi y) \\\partial_y u &= 2\pi^2 \sin^2(\pi x) [\cos^2(\pi y) - \sin^2(\pi y)] \\\partial_x v &= -2\pi^2 [\cos^2(\pi x) - \sin^2(\pi x)] \sin^2(\pi y) \\\partial_y v &= -4\pi^2 \sin(\pi x) \cos(\pi x) \sin(\pi y) \cos(\pi y)\end{aligned}$$

with as expected

$$\partial_x u + \partial_y v = 4\pi^2 \sin(\pi x) \cos(\pi x) \cos(\pi y) \sin(\pi y) - 4\pi^2 \sin(\pi x) \cos(\pi x) \sin(\pi y) \cos(\pi y) = 0$$

The strain rate components are then

$$\begin{aligned}\dot{\epsilon}_{xx} &= \partial_x u \\\dot{\epsilon}_{yy} &= \partial_y v \\\dot{\epsilon}_{xy} &= \frac{1}{2}(\partial_x v + \partial_y u) \\&= \frac{1}{2}(-2\pi^2 [\cos^2(\pi x) - \sin^2(\pi x)] \sin^2(\pi y) + 2\pi^2 \sin^2(\pi x) [\cos^2(\pi y) - \sin^2(\pi y)]) \\&= \pi^2 [\sin^2(\pi x) - \cos^2(\pi x)] \sin^2(\pi y) + \pi^2 \sin^2(\pi x) [\cos^2(\pi y) - \sin^2(\pi y)] \\&= \pi^2 (\sin^2(\pi x) \cos^2(\pi y) - \cos^2(\pi x) \sin^2(\pi y))\end{aligned}$$

We will also need

$$\begin{aligned}\partial_x \dot{\epsilon}_{xx} &= 4\pi^3 [\cos^2(\pi x) - \sin^2(\pi x)] \cos(\pi y) \sin(\pi y) \\\partial_x \dot{\epsilon}_{xy} &= \pi^3 (2 \cos(\pi x) \sin(\pi x) \cos^2(\pi y) + 2 \sin(\pi x) \cos(\pi x) \sin^2(\pi y)) \\&= 2\pi^3 \sin(\pi x) \cos(\pi x) \\\partial_y \dot{\epsilon}_{xy} &= \pi^3 (-2 \sin^2(\pi x) \sin(\pi y) \cos(\pi y) - 2 \cos^2(\pi x) \sin(\pi y) \cos(\pi y)) \\&= -2\pi^3 \sin(\pi y) \cos(\pi y) \\\partial_y \dot{\epsilon}_{yy} &= -4\pi^3 \sin(\pi x) \cos(\pi x) [\cos^2(\pi y) - \sin^2(\pi y)]\end{aligned}$$

Finally we need to compute the root mean square velocity (integrals were obtained with Wolframalpha):

$$\iint u^2 dx dy = 4\pi^2 \underbrace{\int_0^1 \sin^4(\pi x) dx}_{3/8} \cdot \underbrace{\int_0^1 \cos^2(\pi y) \sin^2(\pi y) dy}_{1/8}$$



$$\iint v^2 dx dy = 4\pi^2 \underbrace{\int_0^1 \sin^2(\pi x) \cos^2(\pi x) dx}_{1/8} \cdot \underbrace{\int_0^1 \sin^4(\pi y) dy}_{3/8}$$

so that

$$v_{rms} = \sqrt{\frac{1}{L_x L_y} \iint (u^2 + v^2) dx dy} = \sqrt{\frac{3}{8} \frac{1}{8} + \frac{1}{8} \frac{3}{8}} = \sqrt{\frac{3}{32}} \simeq 0.05412658773$$

Python script:

### 12.1.13 Manufactured solution in Boffi, Cavallini, Gardini, and Gastaldi [110] (2012)

This manufactured solution originates in Boffi, Cavallini, Gardini, and Gastaldi [110] (2012). The velocity field turns out to be identical to the Donea and Huerta [341] manufactured solution. It is based on the stream function

$$\Psi(x, y) = x^2(x-1)^2 y^2(y-1)^2 = f(x)g(y)$$

defined on the unit square. We have

$$\begin{aligned} f'(x) &= 2x(x-1)^2 + 2x^2(x-1) \\ &= 2(x-1)[x(x-1) + x^2] \\ &= 2(x-1)(2x^2 - x) \\ &= 2x(x-1)(2x-1) \\ f''(x) &= 2[(x-1)(2x-1) + x(2x-1) + 2x(x-1)] \\ &= 2(2x^2 - x - 2x + 1 + 2x^2 - x + 2x^2 - 2x) \\ &= 2(6x^2 - 6x + 1) \\ g'(y) &= 2y(y-1)^2 + 2y^2(y-1) \\ &= 2(y-1)[y(y-1) + y^2] \\ &= 2(y-1)(2y^2 - y) \\ &= 2y(y-1)(2y-1) \\ g''(y) &= 2[(y-1)(2y-1) + y(2y-1) + 2y(y-1)] \\ &= 2[2y^2 - 3y + 1 + 2y^2 - y + 2y^2 - 2y] \\ &= 2(6y^2 - 6y + 1) \end{aligned} \tag{12.100}$$

The manufactured 'smooth pressure' solution is then

$$\begin{aligned} u(x, y) &= \partial_y \psi = f(x)g'(y) = x^2(x-1)^2 2y(y-1)(2y-1) \\ v(x, y) &= -\partial_x \psi = -f'(x)g(y) = -2x(x-1)(2x-1)y^2(y-1)^2 \\ p(x, y) &= \frac{1}{2}x^2 - \frac{1}{6} \end{aligned}$$

$$\begin{aligned} \partial_x u(x, y) &= f'(x)g'(y) \\ \partial_y u(x, y) &= f(x)g''(y) \\ \partial_x v(x, y) &= -f''(x)g(y) \\ \partial_y v(x, y) &= -f'(x)g'(y) \end{aligned}$$

$$\begin{aligned}
\dot{\epsilon}_{xx}(x, y) &= f'(x)g'(y) \\
\dot{\epsilon}_{yy}(x, y) &= -f'(x)g'(y) \\
\dot{\epsilon}_{xy}(x, y) &= \frac{1}{2}(f(x)g''(y) - f''(x)g(y))
\end{aligned}$$

$$\begin{aligned}
\frac{\partial p}{\partial x}(x, y) &= x \\
\frac{\partial p}{\partial y}(x, y) &= 0
\end{aligned}$$

$$\begin{aligned}
\partial_x \dot{\epsilon}_{xx}(x, y) &= f''(x)g'(y) \\
\partial_x \dot{\epsilon}_{xy}(x, y) &= \frac{1}{2}(f'(x)g''(y) - f'''(x)g(y)) \\
\partial_y \dot{\epsilon}_{xy}(x, y) &= \frac{1}{2}(f(x)g'''(y) - f''(x)g'(y)) \\
\partial_y \dot{\epsilon}_{yy}(x, y) &= -f'(x)g''(y)
\end{aligned}$$

$$\begin{aligned}
v_{rms} &= \sqrt{\frac{1}{L_x L_y} \iint (u^2 + v^2) dx dy} \\
&= \sqrt{\frac{1}{L_x L_y} \iint (f^2 g'^2 + f'^2 g^2) dx dy} \\
&= \sqrt{\left( \underbrace{\int_0^1 f^2 dx}_{1/630} \underbrace{\int_0^1 g'^2 dy}_{2/105} + \underbrace{\int_0^1 f'^2 dx}_{2/105} \underbrace{\int_0^1 g^2 dy}_{1/630} \right) dx dy} \\
1/33075 * 2 &= \frac{1}{105} \sqrt{\frac{2}{3}} \\
&\simeq 0.007776157913597390787927885951...
\end{aligned} \tag{12.101}$$

The authors also define a non-smooth pressure case:

$$p(x, y) = \begin{cases} y(1-y) \exp(x-1/2)^2 + 1/2 & \text{for } x \geq 1/2 \\ y(1-y) \exp(x-1/2)^2 - 1/2 & \text{for } x < 1/2 \end{cases}$$

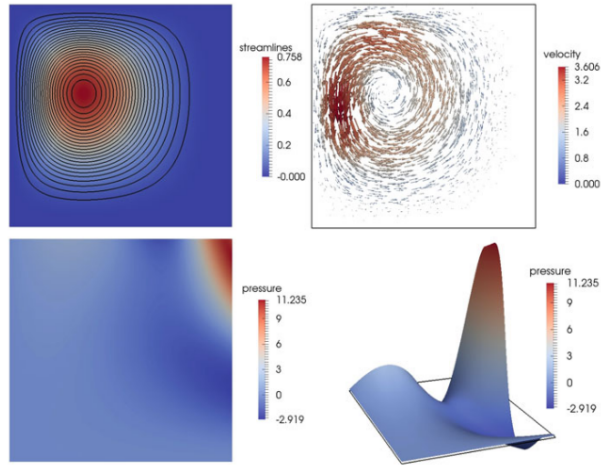
with

$$\begin{aligned}
\frac{\partial p}{\partial x}(x, y) &= 2(x-1/2)y(1-y) \exp(x-1/2)^2 \\
\frac{\partial p}{\partial y}(x, y) &= (1-2y) \exp(x-1/2)^2
\end{aligned}$$

However it is not clear to me how to implement this since only the gradient of the pressure appears in the equations.

### 12.1.14 Manufactured solution in John [650] (book)

This manufactured solution originates in appendix D.1 of John [650] (book).



**Fig. D.1** Example D.3. Stream function (*top left*) velocity (*top right*) and pressure (*bottom*). These plots are based on results obtained with numerical simulations

Taken from [650].

The stream function is given by

$$\Phi(x, y) = 1000x^2(1 - x)^4y^3(1 - y)^2 = 1000f(x)g(y)$$

with

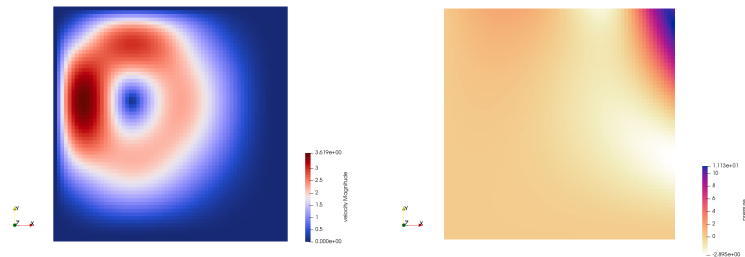
$$\begin{aligned}
f(x) &= x^2(1-x)^4 \\
f'(x) &= 2x(1-x)^4 - 4x^2(1-x)^3 \\
&= [2x(1-x) - 4x^2](1-x)^3 \\
&= (2x - 6x^2)(1-x)^3 \\
&= 2x(1-3x)(1-x)^3 \\
f''(x) &= 2(1-3x)(1-x)^3 - 6x(1-x)^3 - 6x(1-3x)(1-x)^2 \\
&= 2[(1-3x)(1-x) - 3x(1-x) - 3x(1-3x)](1-x)^2 \\
&= 2[1-4x+3x^2-3x+3x^2-3x+9x^2](1-x)^2 \\
&= 2(1-10x+15x^2)(1-x)^2 \\
f'''(x) &= 2[(-10+30x)(1-x)^2 - 2(1-10x+15x^2)(1-x)] \\
&= 2[-10+40x-30x^2-2+20x-30x^2](1-x) \\
&= 2(-12+60x-60x^2)(1-x) \\
&= 24(-1+5x-5x^2)(1-x) \\
&= 24(-1+5x-5x^2+x-5x^2+5x^3) \\
&= 24(-1+6x-10x^2+5x^3)
\end{aligned}$$

$$\begin{aligned}
g(y) &= y^3(1-y)^2 \\
g'(y) &= 3y^2(1-y)^2 - 2y^3(1-y) \\
&= [3y^2(1-y) - 2y^3](1-y) \\
&= y^2(3-3y-2y)(1-y) \\
&= y^2(3-5y)(1-y) \\
g''(y) &= 2y(3-5y)(1-y) - 5y^2(1-y) - y^2(3-5y) \\
&= 2y(3-8y+5y^2) - 5y^2 + 5y^3 - 3y^2 + 5y^3 \\
&= 6y - 16y^2 + 10y^3 - 8y^2 + 10y^3 \\
&= 2y(3-12y+10y^2) \\
g'''(y) &= 6 - 48y + 60y^2
\end{aligned}$$

$$u(x, y) = \partial_y \Phi = 1000f(x)g'(y) = 1000x^2(1-x)^4y^2(3-5y)(1-y) \quad (12.102)$$

$$v(x, y) = -\partial_x \Phi = -1000f'(x)g(y) = -10002x(1-3x)(1-x)^3y^3(1-y)^2 \quad (12.103)$$

$$p(x, y) = \pi^2[xy^3 \cos(2\pi x^2 y) - x^2 y \sin(2\pi xy)] + 1/8 \quad (12.104)$$



$$\begin{aligned}
\partial_x u(x, y) &= 1000 f'(x) g'(x) \\
\partial_y u(x, y) &= 1000 f(x) g''(y) \\
\partial_x v(x, y) &= -1000 f''(x) g(y) \\
\partial_y v(x, y) &= -1000 f'(x) g'(y) \\
\dot{\epsilon}_{xx}(x, y) &= 1000 f' g' \\
\dot{\epsilon}_{yy}(x, y) &= -1000 f' g' \\
\dot{\epsilon}_{xy}(x, y) &= 500(f g'' - f'' g) \\
\partial_x \dot{\epsilon}_{xx}(x, y) &= 1000 f'' g' \\
\partial_x \dot{\epsilon}_{xy}(x, y) &= 500(f' g'' - f''' g) \\
\partial_y \dot{\epsilon}_{xy}(x, y) &= 500(f g''' - f'' g') \\
\partial_y \dot{\epsilon}_{yy}(x, y) &= -1000 f' g''
\end{aligned}$$

Of course we have  $\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} = 0$ .

$$\begin{aligned}
\frac{\partial p}{\partial x}(x, y) &= \pi^2 [y^3 \cos(2\pi x^2 y) - 4\pi x^2 y^4 \sin(2\pi x^2 y) - 2xy \sin(2\pi xy) - 2\pi x^2 y^2 \cos(2\pi xy)] \\
\frac{\partial p}{\partial y}(x, y) &= \pi^2 [3xy^2 \cos(2\pi x^2 y) - 2\pi x^3 y^3 \sin(2\pi x^2 y) - x^2 \sin(2\pi xy) - 2\pi x^3 y \cos(2\pi xy)]
\end{aligned}$$

$$\begin{aligned}
v_{rms} &= \sqrt{\frac{1}{L_x L_y} \iint (u^2 + v^2) dx dy} \\
&= \sqrt{\int_0^1 \int_0^1 u^2 dx dy + \int_0^1 \int_0^1 v^2 dx dy} \\
&= 1000 \sqrt{\int_0^1 \int_0^1 (f g')^2 dx dy + \int_0^1 \int_0^1 (-f' g)^2 dx dy} \\
&= 1000 \sqrt{\underbrace{\int_0^1 f^2 dx}_{\frac{1}{6435}} \underbrace{\int_0^1 g'^2 dy}_{\frac{2}{315}} + \underbrace{\int_0^1 f'^2 dx}_{\frac{2}{693}} \underbrace{\int_0^1 g^2 dy}_{\frac{1}{2310}}} \\
&= 1000 \sqrt{\frac{1}{5 \cdot 9 \cdot 13 \cdot 11} \frac{2}{5 \cdot 9 \cdot 7} + \frac{1}{9 \cdot 11 \cdot 7} \frac{1}{5 \cdot 11 \cdot 3 \cdot 7}} \\
&\simeq 1.4953325891041323968540981
\end{aligned} \tag{12.105}$$

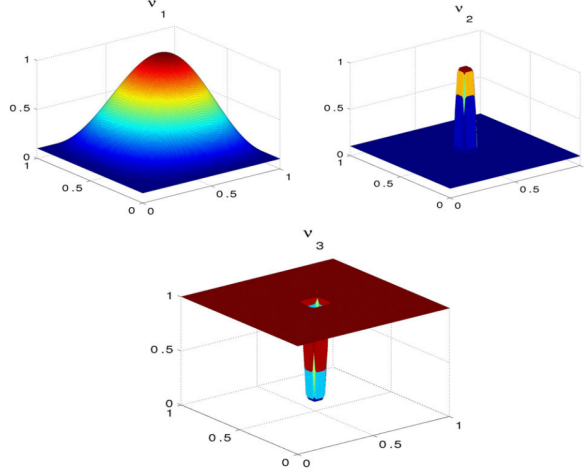
### 12.1.15 Manufactured solution in John, Kaiser, and Novo [652]

This benchmark is identical to the one presented in Section 12.1.14 but it is not isoviscous. There are three different viscosity fields:

$$\eta_1(x, y) = \eta_{min} + (\eta_{max} - \eta_{min})x^2(1-x)y^2(1-y)\frac{721}{16} \quad (12.106)$$

$$\eta_2(x, y) = \eta_{min} + (\eta_{max} - \eta_{min}) \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}] \quad (12.107)$$

$$\eta_3(x, y) = \eta_{min} + (\eta_{max} - \eta_{min}) [1 - \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}]] \quad (12.108)$$



Taken from John, Kaiser, and Novo [652] (2016).

We have

$$b_x = \frac{\partial p}{\partial x} - 2\eta \frac{\partial \dot{\epsilon}_{xx}}{\partial x} - 2\eta \frac{\partial \dot{\epsilon}_{xy}}{\partial y} - 2\frac{\partial \eta}{\partial x} \dot{\epsilon}_{xx} - 2\frac{\partial \eta}{\partial y} \dot{\epsilon}_{xy} \quad (12.109)$$

$$b_y = \frac{\partial p}{\partial y} - 2\eta \frac{\partial \dot{\epsilon}_{xy}}{\partial x} - 2\eta \frac{\partial \dot{\epsilon}_{yy}}{\partial y} - 2\frac{\partial \eta}{\partial x} \dot{\epsilon}_{xy} - 2\frac{\partial \eta}{\partial y} \dot{\epsilon}_{yy} \quad (12.110)$$

The three first terms on each line have already been obtained so we must focus on the last two:

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_1}{\partial x} = \frac{721}{16} [2x(1-x) - x^2]y^2(1-y) \quad (12.111)$$

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_1}{\partial y} = \frac{721}{16} x^2(1-x)[2y(1-y) - y^2] \quad (12.112)$$

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_2}{\partial x} = -10^{14}(x-0.5)^9 \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}] \quad (12.113)$$

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_2}{\partial y} = -10^{14}(y-0.5)^9 \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}] \quad (12.114)$$

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_3}{\partial x} = 10^{14}(x-0.5)^9 \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}] \quad (12.115)$$

$$\frac{1}{(\eta_{max} - \eta_{min})} \frac{\partial \eta_3}{\partial y} = 10^{14}(y-0.5)^9 \exp[-10^{13}(x-0.5)^{10} + (y-0.5)^{10}] \quad (12.116)$$

with

$$u(x, y) = 1000x^2(1-x)^4y^2(3-5y)(1-y) \quad (12.117)$$

$$v(x, y) = -10002x(1-3x)(1-x)^3y^3(1-y)^2 \quad (12.118)$$

$$p(x, y) = \pi^2[xy^3 \cos(2\pi x^2y) - x^2y \sin(2\pi xy)] + 1/8 \quad (12.119)$$

Note that between John's book and the papers [652] and [653] there seems to be a small difference:  $xy^2$  vs.  $xy^3$ . Also the papers have  $-1/8$  while it should be  $+1/8$  (thank you Wolfram Alpha).

### 12.1.16 Manufactured solution in John [649] (1998) on a disc

Let  $\Omega$  be the disc with the center  $(0,0)$  and the radius 1 which has a crack along the  $x$ -axis between the points  $(0,0)$  and  $(1,0)$ .

$$\vec{v}(x, y) = \frac{3}{2}\sqrt{r} \left( \cos \frac{\theta}{2} - \cos \frac{3\theta}{2}, 3 \sin \frac{\theta}{2} - \sin \frac{3\theta}{2} \right)$$

$$p(x, y) = -\frac{6}{\sqrt{r}} \cos \frac{\theta}{2}$$

Viscosity is 1. Note that the solution has a singularity in the origin.

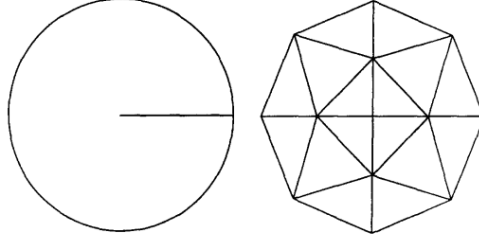


Fig. 4. Domain and coarsest grid (level 0) for Example 11.

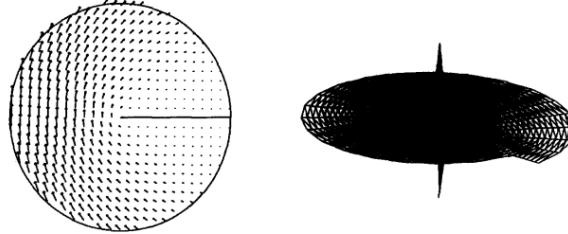


Fig. 5. Solution (velocity and pressure) of Example 11.

Taken from John [649] (1998).

### 12.1.17 Annulus with kinematical b.c. - pure rotation

mms\_annulus.tex

The domain is a hollow cylinder or inner radius  $R_i$  and outside radius  $R_o = 1$ . Boundary conditions are prescribed both on the inside and the outside with  $\vec{v} = (u, v) = (-y, x)$ , or in polar coordinates  $\vec{v} = r\vec{e}_\theta$ .

The gravity is radial and is set to

$$g_x = -x/r \quad g_z = -y/r$$

where  $r = \sqrt{x^2 + z^2}$ , which in polar coordinates is  $\vec{g} = -\vec{e}_r$ . The viscosity is also set to 1, and the density is given by

$$\rho(r) = r^n$$

where  $n$  is a positive or nul integer. The pressure is set to zero at the outer boundary.

The gradient operator in polar coordinates writes:

$$\vec{\nabla} = \frac{\partial}{\partial r} \vec{e}_r + \frac{1}{r} \frac{\partial}{\partial \theta} \vec{e}_\theta$$

and the Laplacian operator:

$$\Delta = \frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2}$$

Note that in our case we need to take the Laplacian of a vector, and unfortunately the Laplacian of a vector is not the Laplacian of the vector's coordinates in polar coordinates (unlike cartesian coordinates). The Laplacian of a vector is given by<sup>4</sup>

$$\nabla^2 \vec{A} = \nabla(\nabla \cdot \vec{A}) - \nabla \times (\nabla \times \vec{A}) = \begin{pmatrix} \frac{\partial^2 A_r}{\partial r^2} + \frac{1}{r} \frac{\partial A_r}{\partial r} - \frac{1}{r^2} A_r + \frac{1}{r^2} \frac{\partial^2 A_r}{\partial \theta^2} - \frac{2}{r^2} \frac{\partial A_\theta}{\partial \theta} \\ \frac{\partial^2 A_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial A_\theta}{\partial r} - \frac{1}{r^2} A_\theta + \frac{1}{r^2} \frac{\partial^2 A_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial A_r}{\partial \theta} \end{pmatrix} = \begin{pmatrix} \Delta A_r \\ \Delta A_\theta \end{pmatrix}$$

The Stokes equation writes:

$$-\vec{\nabla} p + \eta \Delta \vec{v} + \rho \vec{g} = \vec{0}$$

The velocity solution is expected to be  $\vec{v} = r \vec{e}_\theta$ . The Stokes equation in polar coordinates then writes:

$$\begin{aligned} -\frac{\partial p}{\partial r} + \Delta v_r + \rho(r)(-1) &= 0 \\ -\frac{1}{r} \frac{\partial p}{\partial \theta} + \Delta v_\theta &= 0 \end{aligned}$$

Since  $\Delta v_\theta = 0$ , then  $\frac{\partial p}{\partial \theta} = 0$  and then the pressure is independent of  $\theta$ , which is what we expect since the density distribution is radial. We then focus on the first equation, and since  $v_r = 0$ , we then obtain:

$$\frac{\partial p}{\partial r} = -\rho(r)$$

- If  $\rho(r) = 1$ , then

$$\frac{\partial p}{\partial r} = -1$$

yields  $p(r) = -r + C$  where  $C$  is a constant determined by means of b.c. ( $p(r = 1) = 0$ ) so finally

$$\boxed{p(r) = 1 - r}$$

- If  $\rho(r) = r$ , then

$$\frac{\partial p}{\partial r} = -r$$

so that  $p(r) = -\frac{1}{2}r^2 + C$  and likewise

$$\boxed{p(r) = \frac{1}{2}(1 - r^2)}$$

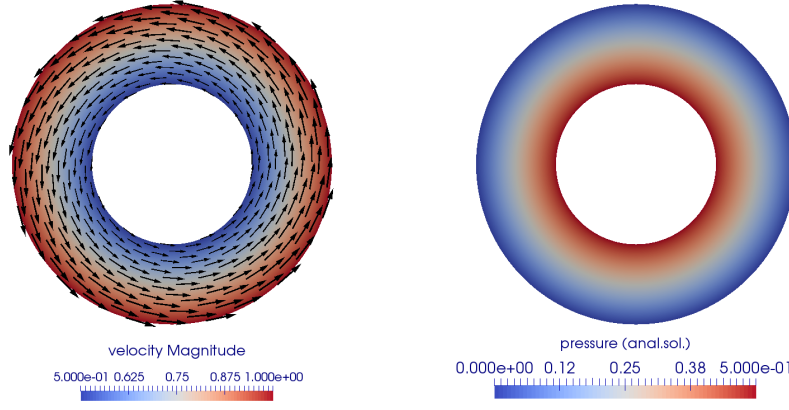
In general, by taking  $\rho(r) = r^n$  with  $n = 0, 1, \dots$  one arrives to a pressure field given by

$$\boxed{p(r) = \frac{1}{n+1}(1 - r^{n+1})}$$

---

<sup>4</sup>[https://en.wikipedia.org/wiki/Vector\\_Laplacian](https://en.wikipedia.org/wiki/Vector_Laplacian)





This benchmark is of course very simple and the fact that the solution is independent of  $\theta$  renders it not so useful. It has succesfully been implemented in ELEFANT .

### 12.1.18 Viscous beam under extension

The domain is a Cartesian box of size  $L_x \times L_y$ . Velocity  $-u_0$  is applied on the left boundary and velocity  $+u_0$  is applied on the right boundary. Bottom and top boundaries are left free. If no vertical velocity is prescribed anywhere there is an obvious nullspace in the solution which is problematic (numerically of course, but also because the solution is then not unique). One might want to set  $v = 0$  at  $y = L_y/2$  on each side for example. The solution to this problem (incompressible Stokes equations) is given by

$$u(x, y) = 2u_0(x/L_x - 1/2) \quad (12.120)$$

$$v(x, y) = -2u_0L_y/L_x(y/L_y - 1/2) \quad (12.121)$$

in the absence of gravity. The strain rate tensor is then:

$$\dot{\epsilon} = \begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{xy} & \dot{\epsilon}_{yy} \end{pmatrix} = \begin{pmatrix} 2u_0/L_x & 0 \\ 0 & -2u_0/L_x \end{pmatrix}$$

and we see that the flow is indeed incompressible as the trace of the strain rate tensor is zero.

The momentum equation is

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta\dot{\epsilon}) = \rho\vec{g}$$

where the viscosity  $\eta$  is constant in space. If gravity is set to zero, we obtain:

$$-\frac{\partial p}{\partial x} = 0 \quad (12.122)$$

$$-\frac{\partial p}{\partial y} = 0 \quad (12.123)$$

since the strain rate is constant in space and the divergence operator applied to it returns the zero tensor. We there fore can conclude that pressure should be constant.

Since the top and bottom boundaries are free, we have  $\sigma \cdot \vec{n} = \vec{0}$  on these. The stress tensor is given by  $\sigma = -\mathbf{1} + 2\eta\dot{\epsilon}$  and the normal on the top is  $\vec{n} = (0, +1)$  so that on the top boundary we have

$$-p + 2\eta\dot{\epsilon}_{yy} = 0$$

or,

$$p = 2\eta\dot{\epsilon}_{yy}$$

Note that using the bottom boundary with  $\vec{n} = (0, -1)$  yields the same result.

### 12.1.19 Channel flow with Herschel-Bulkley rheology

We start from the following formulation for the Herschel-Bulkley rheology:

$$\eta_{HB} = \begin{cases} \eta_0 & \dot{\epsilon}_e \leq \dot{\epsilon}_0 \\ K\dot{\epsilon}_e^{n-1} + \frac{\tau_0}{\dot{\epsilon}_e} & \dot{\epsilon}_e \geq \dot{\epsilon}_0 \end{cases}$$

and the limiting viscosity  $\eta_0$  is such that

$$\eta_0 = K\dot{\epsilon}_0^{n-1} + \frac{\tau_0}{\dot{\epsilon}_0}$$

We consider a two-dimensional channel in the  $x, y$  plane. The walls are at  $y = 0$  and  $y = H$  with no-slip boundary conditions. In the absence of gravity, the Stokes equation simplify to

$$-\frac{\partial p}{\partial x} + \frac{\partial}{\partial y}(2\eta_{HB}\dot{\epsilon}_{xy}) = 0 \quad \text{and} \quad \dot{\epsilon}_{xy} = \frac{1}{2} \frac{\partial u}{\partial y} \quad (12.124)$$

where we assume the velocity  $\vec{v} = (u(y), 0)$ . It then follows that

$$\dot{\epsilon}_e = \sqrt{\mathcal{I}_2(\dot{\epsilon})} = \sqrt{\frac{1}{2}\dot{\epsilon} : \dot{\epsilon}} = \sqrt{\frac{1}{2}[(\dot{\epsilon}_{xx})^2 + (\dot{\epsilon}_{yy})^2 + (\dot{\epsilon}_{zz})^2] + (\dot{\epsilon}_{xy})^2 + (\dot{\epsilon}_{xz})^2 + (\dot{\epsilon}_{yz})^2} = \sqrt{\dot{\epsilon}_{xy}^2} = \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|$$

In the case of a Newtonian fluid, the analytical solution is known and the velocity profile is a parabola with zero velocity on the walls and maximum velocity in the middle. Although the rheology of the fluid is non-linear we assume that a similar velocity profile is expected (although not described by a parabola). We then expect three zones (and we assume that the fluid flows from left to right):

- In the middle, where it is expected that  $\frac{\partial u}{\partial y} = 0$  (at least in one point) because of symmetry. We also therefore expect  $\dot{\epsilon}_e \leq \dot{\epsilon}_0$  in this region so that  $\eta_{HB} = \eta_0$ . How thick this region is will be determined later.

Eq. (12.124) must then be solved

$$\frac{\partial p}{\partial x} = \frac{\partial}{\partial y} \left( 2\eta_{HB} \frac{1}{2} \frac{\partial u}{\partial y} \right) = \eta_0 \frac{\partial^2 u}{\partial y^2} \quad (12.125)$$

Let us call  $\Pi = \frac{\partial p}{\partial x} < 0$ , then we must solve:

$$\frac{\partial^2 u}{\partial y^2} = \frac{\Pi}{\eta_0}$$

The solution is then of the form

$$u(y)|_{mid} = \frac{1}{2} \frac{\Pi}{\eta_0} y^2 + 2ay + b$$

and

$$\dot{\epsilon}_{xy}|_{mid} = \frac{1}{2} \frac{\Pi}{\eta_0} y + a$$

We will determine  $a$  and  $b$  later.

- Near the bottom wall, with  $\frac{\partial u}{\partial y} > 0$  so that  $\dot{\varepsilon}_e = +\frac{1}{2} \left( \frac{\partial u}{\partial y} \right)$  and  $\dot{\varepsilon}_e \geq \dot{\varepsilon}_0$ . We solve Eq. (12.124) again, this time with the non-linear formulation of the viscosity:

$$\begin{aligned}
\frac{\partial p}{\partial x} &= \frac{\partial}{\partial y} \left( 2\eta_{HB} \frac{1}{2} \frac{\partial u}{\partial y} \right) \\
&= 2 \frac{\partial}{\partial y} \left[ \left( K \dot{\varepsilon}_e^{n-1} + \frac{\tau_0}{\dot{\varepsilon}_e} \right) \frac{1}{2} \frac{\partial u}{\partial y} \right] \\
&= 2 \frac{\partial}{\partial y} \left[ \left( K \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|^{n-1} + \tau_0 \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|^{-1} \right) \frac{1}{2} \frac{\partial u}{\partial y} \right] \\
&= 2 \frac{\partial}{\partial y} \left[ K \left( \frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 \right]
\end{aligned} \tag{12.126}$$

We then must solve:

$$\begin{aligned}
\frac{\partial}{\partial y} \left[ K \left( \frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 \right] &= \frac{\Pi}{2} \\
K \left( \frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 &= \frac{\Pi}{2} y + c \\
\left( \frac{1}{2} \frac{\partial u}{\partial y} \right)^n &= \frac{1}{K} \left( \frac{\Pi}{2} y + c - \tau_0 \right)
\end{aligned}$$

or,

$$\boxed{\dot{\varepsilon}_{xy}|_{bot} = \frac{1}{2} \frac{\partial u}{\partial y} = \left( \frac{1}{K} (\Pi_2 y + c - \tau_0) \right)^{1/n}}$$

so

$$\boxed{u(y)|_{bot} = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left( \frac{1}{K} (\Pi_2 y + c - \tau_0) \right)^{1+1/n} + d}$$

where  $\Pi_2 = \Pi/2$

- Near the top wall, with  $\frac{\partial u}{\partial y} < 0$  so that  $\dot{\varepsilon}_e = -\frac{1}{2} \left( \frac{\partial u}{\partial y} \right)$  and  $\dot{\varepsilon}_e \geq \dot{\varepsilon}_0$ . We solve yet again Eq. (12.124):

$$\begin{aligned}
\frac{\partial p}{\partial x} &= \frac{\partial}{\partial y} \left( 2\eta_{HB} \frac{1}{2} \frac{\partial u}{\partial y} \right) \\
&= 2 \frac{\partial}{\partial y} \left[ \left( K \dot{\varepsilon}_e^{n-1} + \frac{\tau_0}{\dot{\varepsilon}_e} \right) \frac{1}{2} \frac{\partial u}{\partial y} \right] \\
&= 2 \frac{\partial}{\partial y} \left[ \left( K \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|^{n-1} + \tau_0 \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|^{-1} \right) \frac{1}{2} \frac{\partial u}{\partial y} \right] \\
&= -2 \frac{\partial}{\partial y} \left[ \left( K \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right)^{n-1} + \tau_0 \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right)^{-1} \right) \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right) \right] \\
&= -2 \frac{\partial}{\partial y} \left[ K \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 \right]
\end{aligned} \tag{12.127}$$

We then must solve:

$$\begin{aligned}
-\frac{\partial}{\partial y} \left[ K \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 \right] &= \frac{\Pi}{2} \\
K \left( -\frac{1}{2} \frac{\partial u}{\partial y} \right)^n + \tau_0 &= -\frac{\Pi}{2} y + e
\end{aligned}$$

$$\left(-\frac{1}{2}\frac{\partial u}{\partial y}\right)^n = \frac{1}{K}(-\frac{\Pi}{2}y + e - \tau_0)$$

which yields

$$\dot{\varepsilon}_{xy}|_{top} = -\left(\frac{1}{K}(-\Pi_2 y + e - \tau_0)\right)^{1/n}$$

$$u(y)|_{top} = 2\frac{n}{n+1}\frac{K}{\Pi_2}\left(\frac{1}{K}(-\Pi_2 y + e - \tau_0)\right)^{1+1/n} + f$$

We have 6 integration constants  $a, b, c, d, e, f$  and 6 additional constraints from continuity or boundary conditions:

$$(1) \quad u(0) = 0 \text{ boundary condition} \quad (12.128)$$

$$(2) \quad u(H) = 0 \text{ boundary condition} \quad (12.129)$$

$$(3) \quad u(y_1) \text{ must be continuous} \quad (12.130)$$

$$(4) \quad u(y_2) \text{ must be continuous} \quad (12.131)$$

$$(5) \quad \dot{\epsilon}_{xy}(y_1) \text{ must be continuous} \quad (12.132)$$

$$(6) \quad \dot{\epsilon}_{xy}(y_2) \text{ must be continuous} \quad (12.133)$$

**Using symmetry to compute  $a$**  Because of symmetry, we expect  $y_1 = H/2 - \delta$  and  $y_2 = H/2 + \delta$  with  $\delta \neq 0$  (i.e.  $y_1 \neq y_2$ ) and we expect  $u(y_1) = u(y_2)$  so that

$$u(y_1)|_{mid} = \frac{1}{2} \frac{\Pi}{\eta_0} y_1^2 + 2ay_1 + b = \frac{1}{2} \frac{\Pi}{\eta_0} y_2^2 + 2ay_2 + b = u(y_2)|_{mid}$$

or,

$$\frac{1}{2} \frac{\Pi}{\eta_0} (y_1^2 - y_2^2) + 2a(y_1 - y_2) = 0$$

$$\frac{1}{2} \frac{\Pi}{\eta_0} (y_1 - y_2)(y_1 + y_2) + 2a(y_1 - y_2) = 0$$

$$\frac{1}{2} \frac{\Pi}{\eta_0} (y_1 + y_2) + 2a = 0$$

$$\frac{1}{2} \frac{\Pi}{\eta_0} H + 2a = 0$$

and finally we obtain  $a$ :

$$a = -\frac{1}{4} \frac{\Pi}{\eta_0} H$$

Note that we could have obtained the same thing by enforcing that the strain rate at  $y_1$  and  $y_2$  are the opposite of one another. It then follows:

$$u(y)|_{mid} = \frac{1}{2} \frac{\Pi}{\eta_0} y^2 - 2 \frac{1}{4} \frac{\Pi}{\eta_0} Hy + b = \frac{1}{2} \frac{\Pi}{\eta_0} (y^2 - yH) + b = \frac{\Pi_2}{\eta_0} (y^2 - yH) + b$$

and

$$\dot{\epsilon}_{xy}|_{mid} = \frac{1}{2} \frac{\Pi}{\eta_0} y - \frac{1}{4} \frac{\Pi}{\eta_0} H = \frac{1}{2} \frac{\Pi}{\eta_0} (y - \frac{H}{2}) = \frac{\Pi_2}{\eta_0} (y - \frac{H}{2})$$

Because of the parabola-like flow profile, we expect the strain rate to be zero in the middle  $y = H/2$ , and positive for  $z_1 < y < H/2$  and negative for  $H/2 < y < z_2$ , which is indeed what we recover ( $\Pi < 0$ ).

**Using bottom boundary condition to obtain  $d$**

$$u(y=0)|_{bot} = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left( \frac{1}{K} (c - \tau_0) \right)^{1+1/n} + d = 0$$

so

$$d = -2 \frac{n}{n+1} \frac{K}{\Pi_2} \left( \frac{1}{K} (c - \tau_0) \right)^{1+1/n}$$

and then

$$u(y)|_{bot} = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{1}{K} (\Pi y + c - \tau_0) \right)^{1+1/n} - \left( \frac{1}{K} (c - \tau_0) \right)^{1+1/n} \right]$$

Using top boundary condition to obtain  $f$

$$u(y = H)|_{top} = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left( \frac{1}{K} (-\Pi_2 H + e - \tau_0) \right)^{1+1/n} + f = 0$$

so

$$f = -2 \frac{n}{n+1} \frac{K}{\Pi_2} \left( \frac{1}{K} (-\Pi_2 H + e - \tau_0) \right)^{1+1/n}$$

and then

$$\boxed{u(y)|_{top} = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{1}{K} (-\Pi_2 y + e - \tau_0) \right)^{1+1/n} - \left( \frac{1}{K} (-\Pi_2 H + e - \tau_0) \right)^{1+1/n} \right]}$$

**computing  $\delta$**  The coordinates of the transitions  $y_1$  and  $y_2$  are the location where the strain rate  $\dot{\epsilon}_e$  reaches  $\dot{\epsilon}_0$ . In other words:

$$\dot{\epsilon}_e|_{mid}(y = y_1) = \dot{\epsilon}_{xy}|_{mid}(y = y_1) = \frac{1}{2} \frac{\Pi}{\eta_0} \left( y_1 - \frac{H}{2} \right) = \frac{1}{2} \frac{\Pi}{\eta_0} \left( \frac{H}{2} - \delta - \frac{H}{2} \right) = -\frac{1}{2} \frac{\Pi}{\eta_0} \delta = \dot{\epsilon}_0$$

or,

$$\delta = -\frac{2\dot{\epsilon}_0\eta_0}{\Pi}$$

Since  $\Pi < 0$  it adds up and  $\delta > 0$ . We can also write

$$\boxed{\delta = \frac{2\dot{\epsilon}_0\eta_0}{|\Pi|}}$$

and we will use throughout what follows:

$$\dot{\epsilon}_0 = -\frac{1}{2} \frac{\Pi}{\eta_0} \delta$$

**Using strain rate continuity at  $y_1$  to compute  $c$**

$$\left( \frac{1}{K} (\Pi_2 y_1 + c - \tau_0) \right)^{1/n} = \frac{1}{2} \frac{\Pi}{\eta_0} \left( y_1 - \frac{H}{2} \right) = -\frac{1}{2} \frac{\Pi}{\eta_0} \delta$$

$$\Pi_2 y_1 + c - \tau_0 = K \left( -\frac{1}{2} \frac{\Pi}{\eta_0} \delta \right)^n$$

$$c = K \left( -\frac{1}{2} \frac{\Pi}{\eta_0} \delta \right)^n + \tau_0 - \Pi_2 y_1$$

$$\boxed{c = K \dot{\epsilon}_0^n + \tau_0 - \Pi_2 y_1}$$

$$\begin{aligned}
u(y)|_{bot} &= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{1}{K} (\Pi_2 y + c - \tau_0) \right)^{1/n+1} - \left( \frac{1}{K} (c - \tau_0) \right)^{1/n+1} \right] \\
&= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{1}{K} (\Pi_2 (y - y_1) + K \dot{\varepsilon}_0^n) \right)^{1/n+1} - \left( \frac{1}{K} (K \dot{\varepsilon}_0^n - \Pi_2 y_1) \right)^{1/n+1} \right] \\
&= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{\Pi_2}{K} (y - y_1) + \dot{\varepsilon}_0^n \right)^{1/n+1} - \left( \dot{\varepsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] \\
\dot{\varepsilon}_{xy}|_{bot} &= \left( \frac{1}{K} (\Pi_2 y + c - \tau_0) \right)^{1/n} \\
&= \left( \frac{\Pi_2}{K} (y - y_1) + \dot{\varepsilon}_0^n \right)^{1/n}
\end{aligned}$$

**Using strain rate continuity at  $y_2$  to compute  $e$**

$$\begin{aligned}
-\left( \frac{1}{K} (-\Pi_2 y_2 + e - \tau_0) \right)^{1/n} &= \frac{1}{2} \frac{\Pi}{\eta_0} \left( y_2 - \frac{H}{2} \right) = \frac{1}{2} \frac{\Pi}{\eta_0} \delta \\
-\Pi_2 y_2 + e - \tau_0 &= K \left( -\frac{1}{2} \frac{\Pi}{\eta_0} \delta \right)^n \\
e &= K \left( -\frac{1}{2} \frac{\Pi}{\eta_0} \delta \right)^n + \tau_0 + \Pi_2 y_2 \\
\boxed{e} &= K \dot{\varepsilon}_0^n + \tau_0 + \Pi_2 y_2
\end{aligned}$$

$$\begin{aligned}
u(y)|_{top} &= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{1}{K} (-\Pi_2 y + e - \tau_0) \right)^{1/n+1} - \left( \frac{1}{K} (-\Pi_2 H + e - \tau_0) \right)^{1/n+1} \right] \\
&= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( -\frac{\Pi_2}{K} (y - y_2) + \dot{\varepsilon}_0^n \right)^{1/n+1} - \left( -\frac{\Pi_2}{K} (H - y_2) + \dot{\varepsilon}_0^n \right)^{1/n+1} \right] \\
\dot{\varepsilon}_{xy}|_{top} &= -\left( \frac{1}{K} (-\Pi_2 y + e - \tau_0) \right)^{1/n} \\
&= -\left( -\frac{\Pi_2}{K} (y - y_2) + \dot{\varepsilon}_0^n \right)^{1/n}
\end{aligned}$$

**Using velocity continuity to compute  $b$**  We use  $u(y_1)|_{bot} = u(y_1)|_{mid}$ :

$$\begin{aligned}
2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{\Pi_2}{K} (y_1 - y_1) + \dot{\varepsilon}_0^n \right)^{1/n+1} - \left( \dot{\varepsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] &= \frac{\Pi_2}{\eta_0} (y_1^2 - y_1 H) + b \\
2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \dot{\varepsilon}_0^{n+1} - \left( \dot{\varepsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] &= \frac{\Pi_2}{\eta_0} y_1 (y_1 - H) + b
\end{aligned}$$

so

$$b = 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \dot{\varepsilon}_0^{n+1} - \left( \dot{\varepsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] - \frac{\Pi_2}{\eta_0} y_1 (y_1 - H)$$

**Using velocity continuity to compute  $b$  (again?)** This time we use  $u(y_2)|_{top} = u(y_2)|_{mid}$ :

$$2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( -\frac{\Pi_2}{K} (y_2 - y_2) + \varepsilon_0^n \right)^{1/n+1} - \left( -\frac{\Pi_2}{K} (H - y_2) + \varepsilon_0^n \right)^{1/n+1} \right] = \frac{\Pi_2}{\eta_0} (y_2^2 - y_2 H) + b$$

$$2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \varepsilon_0^{n+1} - \left( -\frac{\Pi_2}{K} (H - y_2) + \varepsilon_0^n \right)^{1/n+1} \right] = \frac{\Pi_2}{\eta_0} (y_2^2 - y_2 H) + b$$

and since  $H - y_2 = H - H/2 - \delta = H/2 - \delta = y_1$  and

$$y_2^2 - y_2 H = y_2 (y_2 - H) = (H/2 + \delta)(-y_1) = (-H/2 - \delta)y_1 = (-H + H/2 - \delta)y_1 = (-H + y_1)y_1$$

so that we indeed recover the same  $b$  value as above.



To summarize:

$$\begin{aligned}
u(y)|_{bot} &= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( \frac{\Pi_2}{K} (y - y_1) + \dot{\epsilon}_0^n \right)^{1/n+1} - \left( \dot{\epsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] \\
u(y)|_{mid} &= \frac{\Pi_2}{\eta_0} (y^2 - y) + 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ +\dot{\epsilon}_0^{n+1} - \left( \dot{\epsilon}_0^n - \frac{\Pi_2}{K} y_1 \right)^{1/n+1} \right] - \frac{\Pi_2}{\eta_0} y_1 (y_1 - H) \\
u(y)|_{top} &= 2 \frac{n}{n+1} \frac{K}{\Pi_2} \left[ \left( -\frac{\Pi_2}{K} (y + y_2) + \dot{\epsilon}_0^n \right)^{\frac{1}{n}+1} - \left( -\frac{\Pi_2}{K} (H + y_2) + \dot{\epsilon}_0^n \right)^{\frac{1}{n}+1} \right] \\
\dot{\epsilon}_{xy}|_{bot} &= \left( \frac{\Pi_2}{K} (y - y_1) + \dot{\epsilon}_0^n \right)^{1/n} \\
\dot{\epsilon}_{xy}|_{mid} &= \frac{\Pi_2}{\eta_0} \left( y - \frac{H}{2} \right) \\
\dot{\epsilon}_{xy}|_{top} &= - \left( -\frac{\Pi_2}{K} (y - y_2) + \dot{\epsilon}_0^n \right)^{1/n}
\end{aligned}$$

Rather interestingly we find that  $\tau_0$  does not directly enter the equations above. This can be explained as follows: since we have the relationship

$$\eta_0 = K \dot{\epsilon}_0^{n-1} + \frac{\tau_0}{\dot{\epsilon}_0}$$

the parameters  $\eta_0$ ,  $\dot{\epsilon}_0$ ,  $\tau_0$  and  $K$  cannot be all chosen freely. The viscosity  $\eta_0$  is reached when the strain rate becomes smaller than  $\dot{\epsilon}_0$ , so these two parameters have a physical meaning. We set  $\eta_0 = 10^{25}$  and  $\dot{\epsilon}_0 = 10^{-17}$ . When/if  $K$  is zero, then  $\tau_0$  can be interpreted as a yield value for a rigid plastic material so we arbitrarily set it to  $\tau_0 = 10^7$ . Having fixed these parameters we can compute

$$K = \frac{\eta_0 \dot{\epsilon}_0 - \tau_0}{\dot{\epsilon}_0^n}$$

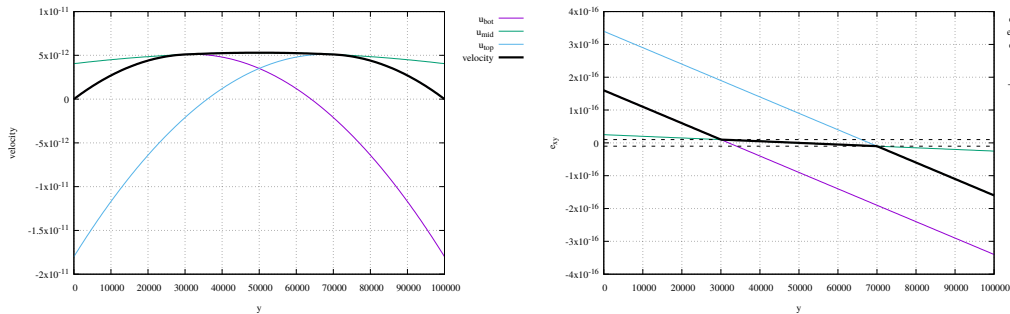
The data used to produce all the following plots is generated by the python program and a gnuplot script to be found in `images/mms/channel_hb/`.

**Let's start simple:**  $n = 1$  In this case the viscosity is given by

$$\eta_{HB} = \begin{cases} \eta_0 & \dot{\epsilon}_e \leq \dot{\epsilon}_0 \\ K + \frac{\tau_0}{\dot{\epsilon}_e} & \dot{\epsilon}_e \geq \dot{\epsilon}_0 \end{cases}$$

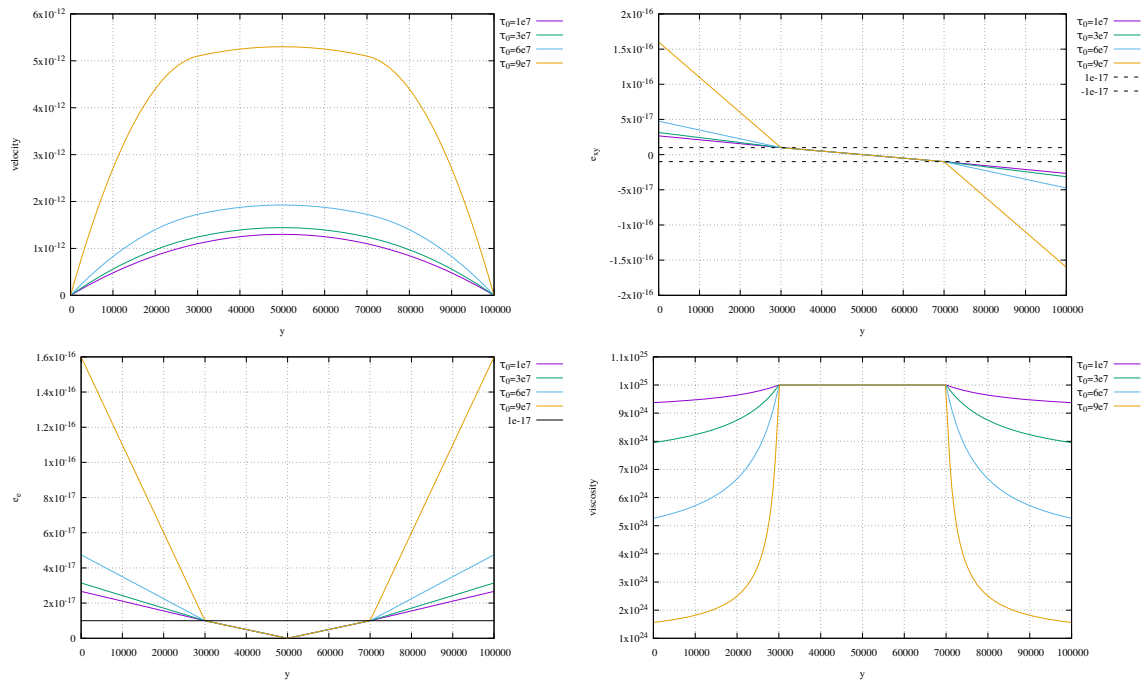
Since  $\dot{\epsilon}_e = \left| \frac{1}{2} \frac{\partial u}{\partial y} \right|$  then

$$\eta_{HB} = \begin{cases} \eta_0 & \dot{\epsilon}_e \leq \dot{\epsilon}_0 \\ K + \frac{2\tau_0}{\left| \frac{\partial u}{\partial y} \right|} & \dot{\epsilon}_e \geq \dot{\epsilon}_0 \end{cases}$$



Obtained for  $n = 1$  and  $\tau_0 = 9e7$ . The black lines are the resulting velocity and strain rate profiles obtained by joining the bottom, middle and top functions.

In the following I explore the effect of the  $\tau_0$  value ( $K$  is calculated correspondingly as we have seen before).



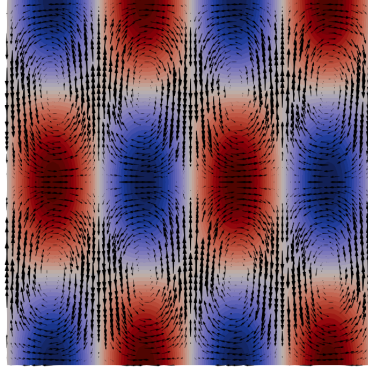
### 12.1.20 Flow in a square using Stream Functions

I wish to arrive at an analytical formulation for a 2D incompressible flow in the square domain  $[-1 : 1] \times [-1 : 1]$ . The fluid has constant viscosity  $\eta = 1$  and is subject to free slip boundary conditions on all sides. For reasons that will become clear in what follows I postulate the following stream function:

$$\Psi(x, y) = \sin(m\pi x) \sin(n\pi y) \quad (12.134)$$

We have the velocity being defined as:

$$\vec{v} = (u, v) = \left( \frac{\partial \Psi}{\partial y}, -\frac{\partial \Psi}{\partial x} \right) = (n\pi \sin(m\pi x) \cos(n\pi y), -m\pi \cos(m\pi x) \sin(n\pi y)) \quad (12.135)$$



Velocity field for  $(m, n) = (2, 1)$

The strain rate components are then:

$$\dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} = mn\pi^2 \cos(m\pi x) \cos(n\pi y) \quad (12.136)$$

$$\dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} = -mn\pi^2 \cos(m\pi x) \cos(n\pi y) \quad (12.137)$$

$$2\dot{\epsilon}_{xy} = \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \quad (12.138)$$

$$= \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \quad (12.139)$$

$$= -n^2\pi^2\Psi + m^2\pi^2\Psi \quad (12.140)$$

$$= (m^2 - n^2)\pi^2 \sin(m\pi x) \sin(n\pi y) \quad (12.141)$$

Note that if  $m = n$  the last term is identically zero, which is not desirable (flow is too 'simple') so in what follows we will prefer  $m \neq n$ .

It is also easy to verify that  $u = 0$  on the sides and  $v = 0$  at the top and bottom and that the term  $\dot{\epsilon}_{xy}$  is nul on all four sides, thereby guaranteeing free slip.

Our choice of stream function yields:

$$\nabla^4 \Psi = \frac{\partial^4 \Psi}{\partial x^4} + \frac{\partial^4 \Psi}{\partial y^4} + 2\frac{\partial^2 \Psi}{\partial x^2 \partial y^2} = \pi^4(m^4\Psi + n^4\Psi + 2m^2n^2\Psi) = (m^4 + n^4 + 2m^2n^2)\pi^4\Psi$$

Let us recall Eq. (??):

$$\vec{\nabla}^4 \Psi = -\frac{\partial \rho g_y}{\partial x} + \frac{\partial \rho g_x}{\partial y} \quad (12.142)$$

We assume  $g_x = 0$  and  $g_y = -1$  so that we simply have

$$\frac{\partial \rho}{\partial x} = (m^4 + n^4 + 2m^2n^2)\pi^4 \Psi = (m^4 + n^4 + 2m^2n^2)\pi^4 \sin(m\pi x) \sin(n\pi y) \quad (12.143)$$

so that (assuming the integration constant to be zero):

$$\rho(x, y) = -\frac{m^4 + n^4 + 2m^2n^2}{m}\pi^3 \cos(m\pi x) \sin(n\pi y)$$

The  $x$ -component of the momentum equation is (since  $g_x = 0$ ):

$$-\frac{\partial p}{\partial x} + \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = -\frac{\partial p}{\partial x} - m^2 n \pi^3 \sin(m\pi x) \cos(n\pi y) - n^3 \pi^3 \sin(m\pi x) \cos(n\pi y) = 0$$

so

$$\frac{\partial p}{\partial x} = -(m^2 n + n^3)\pi^3 \sin(m\pi x) \cos(n\pi y)$$

and the pressure field is then (once again neglecting the integration constant):

$$p(x, y) = \frac{m^2 n + n^3}{m}\pi^2 \cos(\pi x) \cos(\pi y)$$

Note that we then have the interesting property that the pressure average over the domain is zero, i.e.  $\int p dV = 0$ .

### 12.1.21 One-dimensional advection-diffusion equation

Let us start with the 1D steady advection-diffusion equation:

$$\rho C_p u \frac{dT}{dx} - k \frac{d^2 T}{dx^2} = f \quad \text{in } [0, L_x] \quad (12.144)$$

with the boundary conditions  $T(x=0) = 0$  and  $T(x=L_x) = 0$ .

The solution to this problem is:

$$T(x) = \frac{f}{u} \left( x - \frac{1 - \exp \frac{ux}{\kappa}}{1 - \exp \frac{u}{\kappa}} \right)$$

or

$$T(x) = \frac{f}{u} \left( x - \frac{1 - \exp \frac{2Pe x}{h}}{1 - \exp \frac{2Pe}{h}} \right)$$

where

$$Pe = \frac{uh}{2\kappa} = \frac{uh\rho C_p}{2k}$$

### 12.1.22 Annulus with kinematical b.c. - shear flow

Let us consider an annulus domain of inner radius  $R_1$  and outer radius  $R_2$ . Boundary conditions are  $\vec{v}(r=R_1) = \mathbf{v}_1 \vec{e}_\theta$  and  $\vec{v}(r=R_2) = \mathbf{v}_2 \vec{e}_\theta$ . Density is assumed to be zero in the domain.

As seen in Section 12.1.17, the Laplacian of a vector  $\vec{A}$  is given by<sup>5</sup>

$$\nabla^2 \vec{A} = \nabla(\nabla \cdot \vec{A}) - \nabla \times (\nabla \times \vec{A}) = \begin{pmatrix} \frac{\partial^2 A_r}{\partial r^2} + \frac{1}{r} \frac{\partial A_r}{\partial r} - \frac{1}{r^2} A_r + \frac{1}{r^2} \frac{\partial^2 A_r}{\partial \theta^2} - \frac{2}{r^2} \frac{\partial A_\theta}{\partial \theta} \\ \frac{\partial^2 A_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial A_\theta}{\partial r} - \frac{1}{r^2} A_\theta + \frac{1}{r^2} \frac{\partial^2 A_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial A_r}{\partial \theta} \end{pmatrix}$$

<sup>5</sup>[https://en.wikipedia.org/wiki/Vector\\_Laplacian](https://en.wikipedia.org/wiki/Vector_Laplacian)

Given the symmetry of the problem and the boundary conditions we know that the solution is as follows:

$$\vec{v}(r, \theta) = v_\theta(r) \vec{e}_\theta$$

Using this velocity field, we can now obtain the pressure field by solving the Stokes equation

$$-\vec{\nabla} p + \eta \vec{\nabla}^2 \vec{v} = \vec{0}$$

since density is zero. Because of symmetry we also expect the pressure to be a function of  $r$  only, i.e.  $p(r)$ . The Stokes equation in polar coordinates then writes:

$$-\frac{\partial p}{\partial r} = 0 \quad (12.145)$$

$$-\frac{1}{r} \underbrace{\frac{\partial p}{\partial \theta}}_{=0} + \frac{\partial^2 v_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial v_\theta}{\partial r} - \frac{1}{r^2} v_\theta = 0 \quad (12.146)$$

so that the pressure is a constant which we arbitrarily set to zero. We assume  $v_\theta(r) = r^\alpha$ , and the second equation above becomes:

$$\alpha(\alpha - 1)r^{\alpha-2} + \alpha r^{\alpha-2} - r^{\alpha-2} = 0$$

reducing to  $\alpha^2 - 1 = 0$ , i.e.  $\alpha = \pm 1$ , since the above equation must be valid for any value of  $r$ . The generic solution then can be written as

$$v_\theta(r) = Ar + \frac{B}{r}$$

Using the b.c. :

$$v_\theta(R_1) = AR_1 + \frac{B}{R_1} = v_1$$

$$v_\theta(R_2) = AR_2 + \frac{B}{R_2} = v_2$$

or,

$$A + \frac{B}{R_1^2} = \frac{v_1}{R_1}$$

$$A + \frac{B}{R_2^2} = \frac{v_2}{R_2}$$

so

$$B = \frac{\frac{v_1}{R_1} - \frac{v_2}{R_2}}{\frac{1}{R_1^2} - \frac{1}{R_2^2}} = R_1 R_2 \frac{v_1 R_2 - v_2 R_1}{R_2^2 - R_1^2}$$

and

$$A = \frac{v_2 R_2 - v_1 R_1}{R_2^2 - R_1^2}$$

$$v_\theta(r) = \frac{v_2 R_2 - v_1 R_1}{R_2^2 - R_1^2} r + \frac{R_1 R_2}{r} \frac{v_1 R_2 - v_2 R_1}{R_2^2 - R_1^2}$$

We can verify that the flow is indeed incompressible:

$$\vec{\nabla} \cdot \vec{v} = \frac{1}{r} \frac{\partial(r v_r)}{\partial r} + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} = 0$$

since  $v_r = 0$  and  $v_\theta$  does not depend on  $\theta$ .

Note that we could have used a non-zero density: as long as it does not depend on  $\theta$  and the gravity points towards the center, it allows for a decoupling of the equations, thereby only contributing to a lithostatic pressure field.

### 12.1.23 Generic framework for 3D solution in Cartesian coordinates

mms\_generic3D.tex

We postulate

$$u(x, y, z) = f(x)g'(y)h'(z) \quad (12.147)$$

$$v(x, y, z) = f'(x)g(y)h'(z) \quad (12.148)$$

$$w(x, y, z) = -2f'(x)g'(y)h(z) \quad (12.149)$$

so that the flow is indeed incompressible:

$$\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} = f'(x)g'(y)h'(z) + f'(x)g'(y)h'(z) - 2f'(x)g'(y)h'(z) = 0$$

The velocity gradient  $\mathbf{L}(\vec{v})$  is then given by

$$\mathbf{L}(\vec{v}) = \begin{pmatrix} f'g'h' & f''gh' & -2f''g'h \\ fg''h' & f'g'h' & -2f'g''h \\ fg'h'' & f'gh'' & -2f'g'h' \end{pmatrix}$$

and the strain rate tensor by:

$$\dot{\epsilon}(\vec{v}) = \frac{1}{2}(\mathbf{L}(\vec{v}) + \mathbf{L}(\vec{v})^T) = \frac{1}{2} \begin{pmatrix} 2f'g'h' & (f''g + fg'')h' & fg'h'' - 2f''g'h \\ (f''g + fg'')h' & 2f'g'h' & f'gh'' - 2f'g''h \\ fg'h'' - 2f''g'h & f'gh'' - 2f'g''h & -4f'g'h' \end{pmatrix}$$

We assume for simplicity that  $\eta = 1$  so

$$\begin{aligned} \vec{\nabla} \cdot 2\eta\dot{\epsilon}(\vec{v}) &= \vec{\nabla} \cdot \begin{pmatrix} 2f'g'h' & (f''g + fg'')h' & fg'h'' - 2f''g'h \\ (f''g + fg'')h' & 2f'g'h' & f'gh'' - 2f'g''h \\ fg'h'' - 2f''g'h & f'gh'' - 2f'g''h & -4f'g'h' \end{pmatrix} \\ &= \begin{pmatrix} 2f''g'h' + (f''g' + fg''')h' + fg'h''' - 2f''g'h' \\ (f'''g + f'g'')h' + 2f'g''h' + f'gh''' - 2f'g''h' \\ f'g'h'' - 2f'''g'h + f'g'h'' - 2f'g'''h - 4f'g'h'' \end{pmatrix} \\ &= \begin{pmatrix} f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - 2f'g'h'' \end{pmatrix} \end{aligned} \quad (12.150)$$

The Stokes equation then writes

$$\begin{pmatrix} -\partial_x p \\ -\partial_y p \\ -\partial_z p \end{pmatrix} + \begin{pmatrix} f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - 2f'g'h'' \end{pmatrix} + \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} = \vec{0}$$

or,

$$\begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} = \begin{pmatrix} \partial_x p \\ \partial_y p \\ \partial_z p \end{pmatrix} - \begin{pmatrix} f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - 2f'g'h'' \end{pmatrix}$$

**First application** Let us assume

$$f(x) = x(1-x) \quad g(y) = y(1-y) \quad h(z) = z(1-z)$$

then

$$\begin{aligned} f'(x) &= 1-2x & g'(y) &= 1-2y & h'(z) &= 1-2z \\ f''(x) &= -2 & g''(y) &= -2 & h''(z) &= -2 \\ f'''(x) &= 0 & g'''(y) &= 0 & h'''(z) &= 0 \end{aligned}$$

Which ensures that there is no flow through the boundaries of the unit cube. Then

$$\begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} = \begin{pmatrix} \partial_x p \\ \partial_y p \\ \partial_z p \end{pmatrix} - \begin{pmatrix} f''g'h' \\ f'g''h' \\ -2f'g'h'' \end{pmatrix}$$

with

$$u(x, y, z) = x(1-x)(1-2y)(1-2z) \quad (12.151)$$

$$v(x, y, z) = (1-2x)y(1-y)(1-2z) \quad (12.152)$$

$$w(x, y, z) = -2(1-2x)(1-2y)z(1-z) \quad (12.153)$$

Finally we postulate

$$p(x, y, z) = -f'g'h' = (2x-1)(2y-1)(2z-1)$$

with  $\int_{\Omega} p(x, y, z) dx dy dz = 0$ , so

$$\begin{pmatrix} \partial_x p \\ \partial_y p \\ \partial_z p \end{pmatrix} = \begin{pmatrix} -f''g'h' \\ -f'g''h' \\ -f'g'h'' \end{pmatrix} = \begin{pmatrix} 2(2y-1)(2z-1) \\ 2(2x-1)(2z-1) \\ 2(2x-1)(2y-1) \end{pmatrix}$$

Then we find that

$$\vec{b} = \begin{pmatrix} -f''g'h' \\ -f'g''h' \\ -f'g'h'' \end{pmatrix} - \begin{pmatrix} f''g'h' \\ f'g''h' \\ -2f'g'h'' \end{pmatrix} = \begin{pmatrix} -2f''g'h' \\ -2f'g''h' \\ f'g'h'' \end{pmatrix} = \begin{pmatrix} 4(2y-1)(2z-1) \\ 4(2x-1)(2z-1) \\ -2(2x-1)(2y-1) \end{pmatrix}$$

The root mean square velocity is given by

$$\int_{\Omega} u^2 dV = 1/270 \quad (12.154)$$

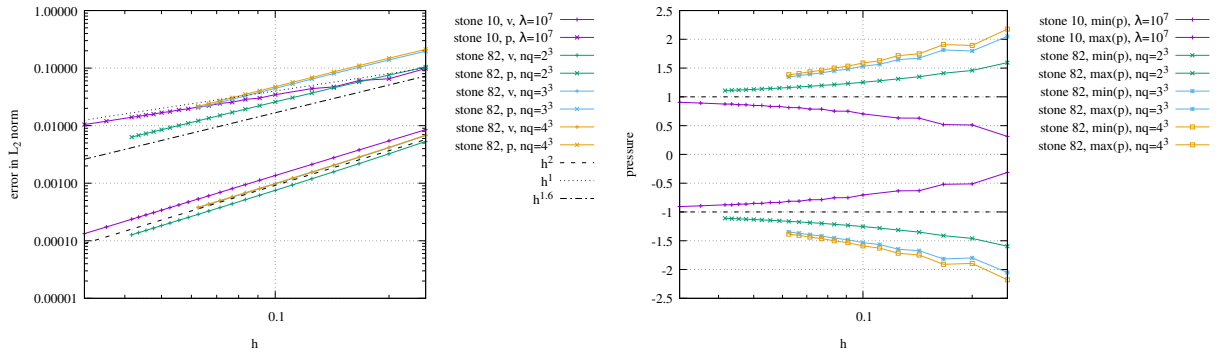
$$\int_{\Omega} v^2 dV = 1/270 \quad (12.155)$$

$$\int_{\Omega} w^2 dV = 2/135 \quad (12.156)$$

and then

$$v_{rms} = \sqrt{\frac{6}{270}} \simeq 0.1490712$$

This benchmark is implemented in [STONE](#) 10 ( $Q_1 \times P_0$ ), [STONE](#) 75 (MINI-1 bubble) and [STONE](#) 82 (MINI-2 bubbles). The code is available in `/mms/mms3D.py`.



Error convergence and pressure statistics



**Second application** Once again, no flow through the boundaries of the unit cube:

$$f(x) = x^2(1-x)^2 \quad g(y) = y^2(1-y)^2 \quad h(z) = z^2(1-z)^2$$

then

$$\begin{aligned} f'(x) &= 2x(2x^2 - 3x + 1) & g'(y) &= 2y(2y^2 - 3y + 1) & h'(z) &= 2z(2z^2 - 3z + 1) \\ f''(x) &= 2(6x^2 - 6x + 1) & g''(y) &= 2(6y^2 - 6y + 1) & h''(z) &= 2(6z^2 - 6z + 1) \\ f'''(x) &= 24x - 12 & g'''(y) &= 24y - 12 & h'''(z) &= 24z - 12 \end{aligned}$$

The velocity field is then given by

$$\begin{aligned} u(x, y, z) &= f(x)g'(y)h'(z) = 4x^2(1-x)^2y(2y^2 - 3y + 1)z(2z^2 - 3z + 1) \\ v(x, y, z) &= f'(x)g(y)h'(z) = 4x(2x^2 - 3x + 1)y^2(1-y)^2z(2z^2 - 3z + 1) \\ w(x, y, z) &= -2f'(x)g'(y)h(z) = -8x(2x^2 - 3x + 1)y(2y^2 - 3y + 1)z^2(1-z)^2 \end{aligned} \quad (12.157)$$

We choose

$$p(x, y, z) = -f'g'h' = -2x(2x^2 - 3x + 1)2y(2y^2 - 3y + 1)2z(2z^2 - 3z + 1)$$

The rhs vector is then

$$\begin{aligned} \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} &= \begin{pmatrix} \partial_x p \\ \partial_y p \\ \partial_z p \end{pmatrix} - \begin{pmatrix} f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - 2f'g'h'' \end{pmatrix} \\ &= - \begin{pmatrix} f''g'h' \\ f'g''h' \\ f'g'h'' \end{pmatrix} - \begin{pmatrix} f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - 2f'g'h'' \end{pmatrix} \\ &= - \begin{pmatrix} 2f''g'h' + fg'''h' + fg'h''' \\ f'''gh' + 2f'g''h' + f'gh''' \\ -2f'''g'h - 2f'g'''h - f'g'h'' \end{pmatrix} \end{aligned} \quad (12.158)$$

Let us look at the root mean square velocity:

$$\begin{aligned} \mathbf{v}_{rms}^2 &= \iiint (u^2 + v^2 + w^2) dx dy dz \\ &= \iiint u^2 dx dy dz + \iiint v^2 dx dy dz + \iiint w^2 dx dy dz \\ &= \iiint (f(x)g'(y)h'(z))^2 dx dy dz + \iiint (f'(x)g(y)h'(z))^2 dx dy dz + \iiint (-2f'(x)g'(y)h(z))^2 dx dy dz \\ &= \left( \int_0^1 f(x)^2 dx \right) \left( \int_0^1 g'(y)^2 dy \right) \left( \int_0^1 h'(z)^2 dz \right) \\ &+ \left( \int_0^1 f'(x)^2 dx \right) \left( \int_0^1 g(y)^2 dy \right) \left( \int_0^1 h'(z)^2 dz \right) \\ &+ 4 \left( \int_0^1 f'(x)^2 dx \right) \left( \int_0^1 g'(y)^2 dy \right) \left( \int_0^1 h(z)^2 dz \right) \end{aligned} \quad (12.159)$$

Since

$$\int_0^1 f(x)^2 dx = \int_0^1 g(y)^2 dy = \int_0^1 h(z)^2 dz$$

and

$$\int_0^1 g'(y)^2 dy = \int_0^1 f'(x)^2 dx = \int_0^1 h'(z)^2 dz$$

we only have 2 integrals to compute and using WolframAlpha we find:

$$\begin{aligned} \left( \int_0^1 f(x)^2 dx \right) &= \frac{1}{630} \\ \left( \int_0^1 f'(x)^2 dx \right) &= \frac{2}{105} \end{aligned}$$

so that

$$\begin{aligned} \mathbf{v}_{rms}^2 &= \frac{1}{630} \frac{2}{105} \frac{2}{105} + \frac{2}{105} \frac{1}{630} \frac{2}{105} + 4 \frac{1}{630} \frac{2}{105} \frac{1}{630} \\ &= 6 \frac{1}{630} \frac{2}{105} \frac{2}{105} \\ &= \frac{1}{105} \frac{2}{105} \frac{2}{105} \end{aligned} \tag{12.160}$$

In the end:

$$\mathbf{v}_{rms} \simeq 0.00185885728$$

This benchmark is carried out in Stone 10.

### 12.1.24 2D Analytical benchmark XII

This is presented in Soulaïmani *et al.* (1987) [1181]. The velocity field is given by

$$\vec{v}(x, y) = (x^3, -3x^2y)$$

and the pressure is

$$p(x, y) = x^3 + y^3 - 1/2$$

so that, assuming that the viscosity is 1, the body force is:

$$\vec{b} = (-6x + 3x^2, 6y + 3y^2)$$

Note that I have added the  $-1/2$  term to the pressure so that  $\int \int p dx dy = 0$ . The root mean square velocity over a unit square is

$$v_{rms} = \sqrt{\int_0^1 \int_0^1 (u^2 + v^2) dx dy} = \sqrt{\int_0^1 \int_0^1 (x^6 + 9x^4y^2) dx dy} = \sqrt{\frac{1}{7} + 9\frac{1}{5}\frac{1}{3}} = \sqrt{\frac{26}{35}} \simeq 0.861892$$

The strain rate tensor terms are

$$\begin{aligned}\dot{\epsilon}_{xx} &= 3x^2 \\ \dot{\epsilon}_{yy} &= -3x^2 \\ \dot{\epsilon}_{xy} &= -3xy\end{aligned}$$

Another one mentioned in the paper:

$$\vec{v}(x, y) = (x^2, -2xy) \quad p(x, y) = 0 \quad \vec{b} = (-2, 0)$$

### 12.1.25 2D analytical benchmark from Burman & Hansbo (2006)

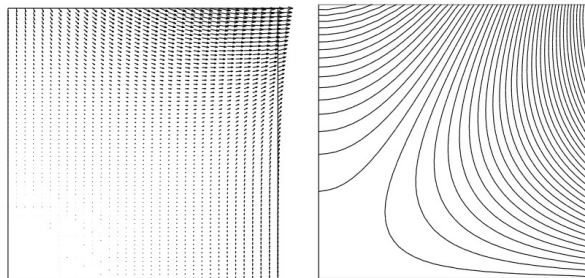
This is presented in Burman and Hansbo [181] (2006) and Burman and Hansbo [182] (2007) and apparently originates in Norburn and Silvester [945] (1998).

The velocity and pressure fields are given in the unit square by

$$\begin{aligned}u(x, y) &= 20xy^3 \\ v(x, y) &= 5x^4 - 5y^4 \\ p(x, y) &= 60x^2y - 20y^3 - 5\end{aligned} \tag{12.161}$$

with

$$\begin{aligned}\frac{\partial u}{\partial x} &= 20y^3 \\ \frac{\partial u}{\partial y} &= 60xy^2 \\ \frac{\partial v}{\partial x} &= 20x^3 \\ \frac{\partial v}{\partial y} &= -20y^3\end{aligned} \tag{12.162}$$



Taken from Burman and Hansbo [181]. Left is velocity, right is pressure.

The flow is incompressible:

$$\operatorname{div}(\vec{v}) = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = 0$$

Then the strain rate tensor is given by

$$\dot{\epsilon}(\vec{v}) = \begin{pmatrix} 20y^3 & 30xy^2 + 10x^3 \\ 30xy^2 + 10x^3 & -20y^3 \end{pmatrix}$$

Assuming the viscosity  $\eta = 1$ , then the full stress tensor is given by

$$\begin{aligned} \sigma &= -p\mathbf{1} + 2\eta\dot{\epsilon}(\vec{v}) \\ &= \begin{pmatrix} -60x^2y + 20y^3 + 5 + 40y^3 & 60xy^2 + 20x^3 \\ 60xy^2 + 20x^3 & -60x^2y + 20y^3 + 5 - 40y^3 \end{pmatrix} \\ &= \begin{pmatrix} -60x^2y + 60y^3 + 5 & 60xy^2 + 20x^3 \\ 60xy^2 + 20x^3 & -60x^2y - 20y^3 + 5 \end{pmatrix} \end{aligned} \quad (12.163)$$

finally

$$\vec{b} = -\vec{\nabla} \cdot \sigma = \begin{pmatrix} -120xy + 120xy \\ 60y^2 + 60x^2 - 60x^2 - 60y^2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

This is particularly convenient...

It is implemented in [STONE](#) 14, 18 and 115.

### 12.1.26 2D analytical benchmark from Cioncolini & Boffi (2019)

This is Test #1 presented in Cioncolini and Boffi [258] (2019). The domain is a unit square and the solution is given by

$$\begin{aligned} u(x, y) &= x^2(1-x)^2 2y(1-y)(2y-1) \\ v(x, y) &= y^2(1-y)^2 2x(1-x)(1-2x) \\ p(x, y) &= x(1-x)(1-y) - 1/12 \\ b_x &= -\eta \{ 4y(1-y)(2y-1)[(1-2x)^2 - 2x(1-x)] + 12x^2(1-x)^2(1-2y) \} + (1-2x)(1-2y) \\ b_y &= \end{aligned} \quad (12.165)$$

with no-slip boundary conditions on all sides. Note that the velocity solution is close to the mms of Section ??

### 12.1.27 3D Analytical benchmark XIII

This is presented in Soulaïmani *et al.* (1987) [1181].

$$\vec{b} = (1, 1, 1) \quad \vec{v}(x, y, z) = (y, z, x) \quad p(x, y, z) = x + y + z - 1/2$$

### 12.1.28 2D Analytical benchmark XIV

It originates in Section 6 of John *et al.* (2017) [655]. The velocity is given by

$$\begin{aligned} v_x &= 200x^2(1-x)^2y(1-y)(1-2y) = 100a(x)a'(y) \\ v_y &= -200x(1-x)(1-2x)y^2(1-y)^2 = -100a'(x)a(y) \end{aligned}$$

with

$$\begin{aligned}
a(x) &= x^2(1-x)^2 \\
a'(x) &= 2x(1-x)^2 - 2x^2(1-x) = 2x(1-x)(1-2x) \\
a''(x) &= 2(1-6x+6x^2) \\
a'''(x) &= 24x-12
\end{aligned}$$

We can now compute the components of the strain rate tensor:

$$\begin{aligned}
\dot{\epsilon}_{xx} &= \frac{\partial \mathbf{v}_x}{\partial x} = \frac{\partial(100a(x)a'(y))}{\partial x} = 100a'(x)a'(y) \\
\dot{\epsilon}_{yy} &= \frac{\partial \mathbf{v}_y}{\partial y} = \frac{\partial(-100a'(x)a(y))}{\partial y} = -100a'(x)a'(y) \\
\dot{\epsilon}_{xy} = \dot{\epsilon}_{yx} &= \frac{1}{2} \left( \frac{\partial \mathbf{v}_x}{\partial y} + \frac{\partial \mathbf{v}_y}{\partial x} \right) = \frac{1}{2} (100a(x)a''(y) - 100a''(x)a(y)) = 50 (a(x)a''(y) - a''(x)a(y))
\end{aligned}$$

The momentum conservation equation is given by

$$\begin{aligned}
-\partial_x p + \partial_x(2\eta\dot{\epsilon}_{xx}) + \partial_y(2\eta\dot{\epsilon}_{xy}) + b_x &= 0 \\
-\partial_y p + \partial_x(2\eta\dot{\epsilon}_{xy}) + \partial_y(2\eta\dot{\epsilon}_{yy}) + b_y &= 0
\end{aligned}$$

Then

$$\begin{aligned}
b_x &= \partial_x p - \partial_x(2\eta\dot{\epsilon}_{xx}) - \partial_y(2\eta\dot{\epsilon}_{xy}) \\
&= \partial_x p - \partial_x[2\eta 100a'(x)a'(y)] - \partial_y[2\eta 50 (a(x)a''(y) - a''(x)a(y))] \\
&= \partial_x p - 200\eta a''(x)a'(y) - 100\eta[a(x)a'''(y) - a''(x)a'(y)] \\
&= \partial_x p - 100\eta a''(x)a'(y) - 100\eta a(x)a'''(y) \\
&= \partial_x p - 100\eta[a''(x)a'(y) + a(x)a'''(y)] \\
b_y &= \partial_y p - \partial_x(2\eta\dot{\epsilon}_{xy}) - \partial_y(2\eta\dot{\epsilon}_{yy}) \\
&= \partial_y p - \partial_x[2\eta 50 (a(x)a''(y) - a''(x)a(y))] + \partial_y 2\eta 100a'(x)a'(y) \\
&= \partial_y p - 100\eta(a'(x)a''(y) - a'''(x)a(y)) + 200\eta a'(x)a''(y) \\
&= \partial_y p + 100\eta a'(x)a''(y) + 100\eta a'''(x)a(y) \\
&= \partial_y p + 100\eta[a'(x)a''(y) + a'''(x)a(y)]
\end{aligned}$$

with

$$\begin{aligned}
p(x, y) &= 10 \left[ \left( x - \frac{1}{2} \right)^3 y^2 + (1-x)^3 \left( y - \frac{1}{2} \right)^3 \right] \\
\frac{\partial p}{\partial x} &= 10 \left[ 3 \left( x - \frac{1}{2} \right)^2 y^2 - 3(1-x)^2 \left( y - \frac{1}{2} \right)^3 \right] \\
\frac{\partial p}{\partial y} &= 10 \left[ \left( x - \frac{1}{2} \right)^3 2y + (1-x)^3 3 \left( y - \frac{1}{2} \right)^2 \right]
\end{aligned}$$

See [STONE](#) 104.

### 12.1.29 Poisson equation on 3D shell

This benchmark is presented in Phillips *et al.* [996]. Inner radius is 1, outer radius is 3. The right hand side term of the Poisson equation is given by:

$$\begin{aligned}
 f(r, \theta, \phi) = & \frac{\sin^2 \theta}{R^2} [(\cos \phi - \sin \phi)(20 \sin^2 \theta - 15) - \sin 2\phi(10 \sin^2 \theta - 6)] \\
 & \times \left[ \left( \frac{r}{R_{inner}} \right)^2 - 1 \right] \left[ \left( \frac{r}{R_{outer}} \right)^2 - 1 \right] \\
 & + \sin^4 \theta \left[ \cos \phi - \sin \phi - \frac{1}{2} \sin 2\phi \right] \left[ \frac{20r^2}{R_{inner}^2 R_{outer}^2} - 6 \left( \frac{1}{R_{inner}^2} + \frac{1}{R_{outer}^2} \right) \right] \quad (12.166)
 \end{aligned}$$

Note: what is  $R$  here??

The solution is then

$$T(r, \theta, \phi) = \sin^4 \theta \left[ \cos \phi - \sin \phi - \frac{1}{2} \sin 2\phi \right] \left[ \left( \frac{r}{R_{inner}} \right)^2 - 1 \right] \left[ \left( \frac{r}{R_{outer}} \right)^2 - 1 \right]$$

This expression is used to generate the Dirichlet boundary conditions on the inner and outer surfaces.

### 12.1.30 SolCx

The SolCx benchmark is intended to test the accuracy of the solution to a problem that has a large jump in the viscosity along a line through the domain. Such situations are common in geophysics: for example, the viscosity in a cold, subducting slab is much larger than in the surrounding, relatively hot mantle material.

The SolCx benchmark computes the Stokes flow field of a fluid driven by spatial density variations, subject to a spatially variable viscosity. Specifically, the domain is  $\Omega = [0, 1]^2$ , gravity is  $\vec{g} = (0, -1)^T$  and the density is given by

$$\rho(x, y) = \sin(\pi y) \cos(\pi x) \quad (12.167)$$

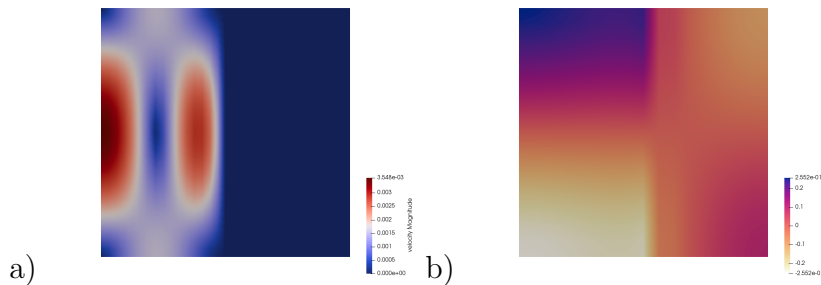
Boundary conditions are free slip on all of the sides of the domain and the temperature plays no role in this benchmark. The viscosity is prescribed as follows:

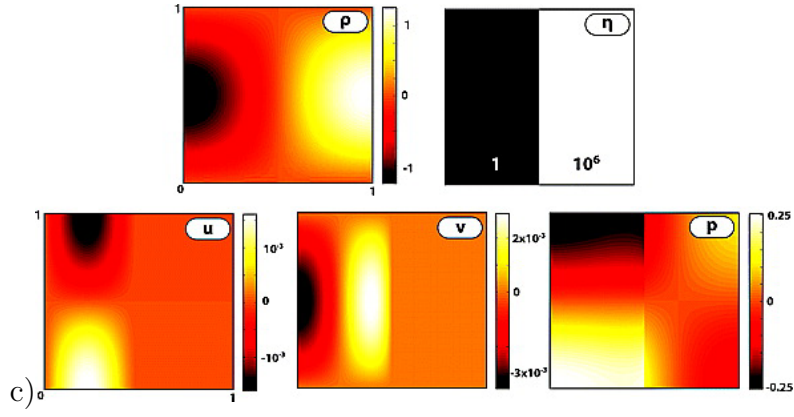
$$\eta(x, y) = \begin{cases} 1 & \text{for } x < 0.5 \\ 10^6 & \text{for } x > 0.5 \end{cases} \quad (12.168)$$

Note the strongly discontinuous viscosity field yields a stagnant flow in the right half of the domain and thereby yields a pressure discontinuity along the interface.

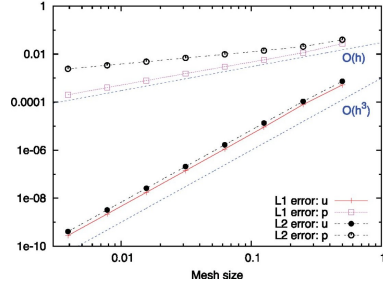
The SolCx benchmark was previously showcased in Duretz *et al.* (2011) [352] and its analytic solution is given in Zhong (1996) [1416]. It has been carried out in Kronbichler *et al.* (2012) [732] and Gerya *et al.* (2013) [452], and is also found in the ASPECT manual [44].

Note that the source code which evaluates the velocity and pressure fields for both SolCx and SolKz is distributed as part of the open source package Underworld (Moresi *et al.* , 2007 [901], <http://underworldproject.org>). I have translated this code to python.





a,b) obtained with ASPECT . c) taken from Duretz *et al.* (2011) [352].



|       | $Q_2^d \times Q_1$    |                        | $Q_3^d \times Q_2$    |                        | $Q_3^d \times P_{-1}$ , even mesh |                      | $Q_2^d \times P_{-1}$ , odd mesh |                        |
|-------|-----------------------|------------------------|-----------------------|------------------------|-----------------------------------|----------------------|----------------------------------|------------------------|
| $h$   | $\ e_u\ _{L_2}$       | $\ e_p\ _{L_2}$        | $\ e_u\ _{L_2}$       | $\ e_p\ _{L_2}$        | $\ e_u\ _{L_2}$                   | $\ e_p\ _{L_2}$      | $\ e_u\ _{L_2}$                  | $\ e_p\ _{L_2}$        |
| 1/8   | $1.3 \times 10^{-5}$  | $1.4 \times 10^{-2}$   | $6.3 \times 10^{-7}$  | $8.8 \times 10^{-3}$   | $1.3 \times 10^{-5}$              | $1.5 \times 10^{-3}$ | $6.5 \times 10^{-4}$             | $1.1 \times 10^{-2}$   |
| 1/16  | $1.7 \times 10^{-6}$  | $9.8 \times 10^{-3}$   | $4.0 \times 10^{-8}$  | $6.2 \times 10^{-3}$   | $1.7 \times 10^{-6}$              | $3.7 \times 10^{-4}$ | $3.6 \times 10^{-4}$             | $6.8 \times 10^{-3}$   |
| 1/32  | $2.1 \times 10^{-7}$  | $6.9 \times 10^{-3}$   | $2.6 \times 10^{-9}$  | $4.4 \times 10^{-3}$   | $2.2 \times 10^{-7}$              | $9.2 \times 10^{-5}$ | $1.9 \times 10^{-4}$             | $4.5 \times 10^{-3}$   |
| 1/64  | $2.6 \times 10^{-8}$  | $4.9 \times 10^{-3}$   | $1.7 \times 10^{-10}$ | $3.1 \times 10^{-3}$   | $2.6 \times 10^{-8}$              | $2.3 \times 10^{-5}$ | $9.8 \times 10^{-5}$             | $3.3 \times 10^{-3}$   |
| 1/128 | $3.3 \times 10^{-9}$  | $3.4 \times 10^{-3}$   | $2.0 \times 10^{-11}$ | $2.2 \times 10^{-3}$   | $3.2 \times 10^{-9}$              | $5.7 \times 10^{-6}$ | $5.0 \times 10^{-5}$             | $2.1 \times 10^{-3}$   |
| 1/256 | $4.1 \times 10^{-10}$ | $2.4 \times 10^{-3}$   | $1.7 \times 10^{-11}$ | $1.5 \times 10^{-3}$   | $4.1 \times 10^{-10}$             | $1.4 \times 10^{-6}$ | $2.5 \times 10^{-5}$             | $1.5 \times 10^{-3}$   |
|       | $\mathcal{O}(h^3)$    | $\mathcal{O}(h^{1/2})$ | $\mathcal{O}(h^4)$    | $\mathcal{O}(h^{1/2})$ | $\mathcal{O}(h^3)$                | $\mathcal{O}(h^2)$   | $\mathcal{O}(h)$                 | $\mathcal{O}(h^{1/2})$ |

Taken from Kronbichler *et al.* (2012) [732]. Velocity and pressure errors  $e_u$ ,  $e_p$  and convergence rates for different choices of the Stokes finite element spaces, using globally refined meshes. For ‘odd’ meshes, the numbers shown are the average errors from nearby meshes (e.g. for  $h = 1/64$ , the average of the errors on  $63 \times 63$  and  $65 \times 65$  meshes).

## Relevant Literature:

- D.A. May and L. Moresi. “Preconditioned iterative methods for Stokes flow problems arising in computational geodynamics”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 33–47. DOI: 10.1016/j.pepi.2008.07.036
- M Velić, L. Moresi, D. May, and M. Knepley. “A Family of Numerically Stable Analytic Solutions for Geodynamic Code Verification”. In: ()
- Albert de Montserrat, Jason P Morgan, and Jörg Hasenclever. “LaCoDe: a Lagrangian two-dimensional thermo-mechanical code for large-strain compressible visco-elastic geodynamical modeling”. In: *Tectonophysics* 767 (2019), p. 228173. DOI: 10.1016/j.tecto.2019.228173
- Y.A. Mishin, O.V. Vasilyev, and T.V. Gerya. “A Wavelet-Based Adaptive Finite Element Method for the Stokes Problems”. In: *Fluids* 7 (2022), p. 221. DOI: 10.3390/fluids7070221
- Ruben Sevilla and Thibault Duretz. “A face-centered finite volume method for high-contrast Stokes interface problems”. In: *International Journal for Numerical Methods in Engineering* 124 (2023), pp. 3709–3732. DOI: 10.1002/nme.7294

STONE 5, 77

## 12.1.31 SolKz

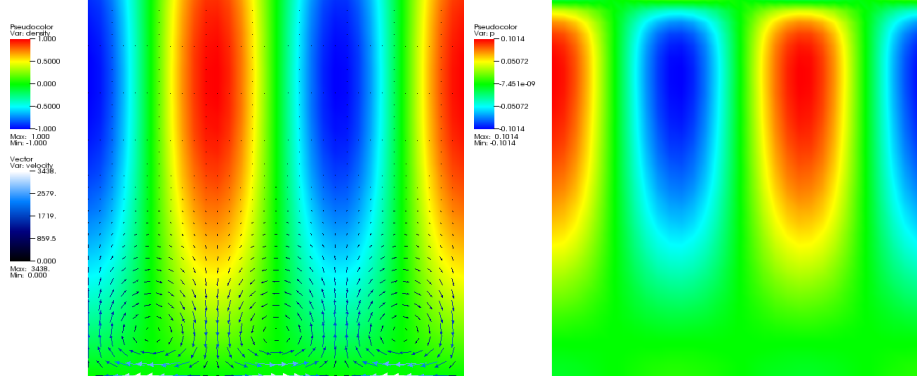
The SolKz benchmark is similar to the SolCx benchmark: the viscosity is a function of the space coordinates too and is given by

$$\eta(y) = \exp(2By) \quad \text{with} \quad B = 13.8155$$

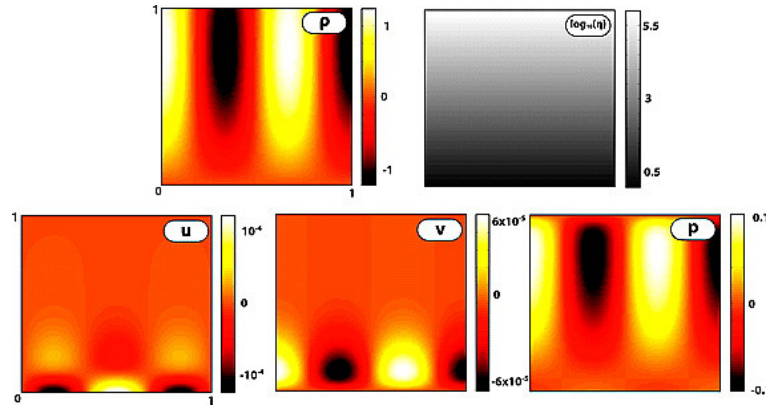
It is not a discontinuous function but grows exponentially with the vertical coordinate so that its overall variation is again  $10^6$ . The forcing is again chosen by imposing a spatially variable density variation as follows:

$$\rho(x, y) = \sin(2y) \cos(3\pi x)$$


Free slip boundary conditions are imposed on all sides of the domain. This benchmark too is presented in [1416] and is studied in [352] and [452].



Taken from ASPECT manual [44].



Taken from Duretz *et al.* (2011) [352].

 Relevant Literature: [902], [846], [1316], [322], [1062] [STONE 06](#)

### 12.1.32 SolVi

SolVi is another very common benchmark carried out in the computational geodynamics literature.

This inclusion benchmark solves a problem with a discontinuous viscosity, which is chosen in such a way that the discontinuity is a circle. Given the regular nature of the used by a majority of codes, this ensures that the discontinuity in the viscosity never aligns to cell boundaries. This in turns leads to almost discontinuous pressures along the interface which are difficult to represent accurately.

Schmid & Podlachikov (2003) [1128]. derived a simple analytic solution for the pressure and velocity fields for such a circular inclusion under simple shear.

A characteristic of the analytical solution is that the pressure is zero inside the inclusion, while outside it follows the relation

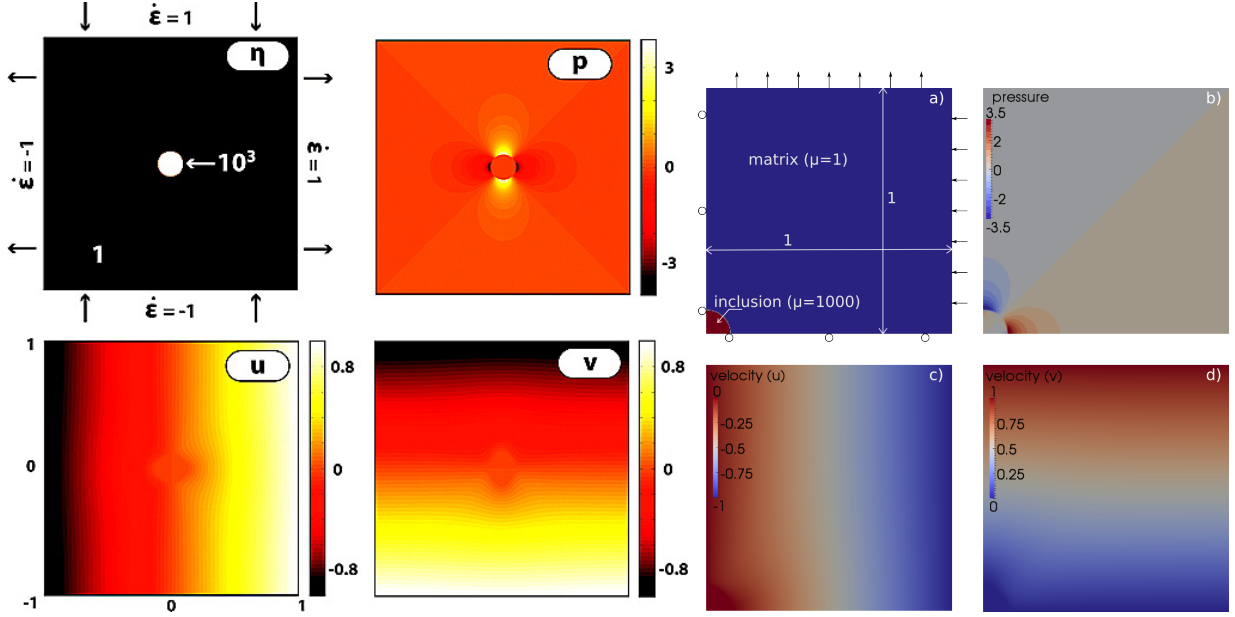
$$p_m = 4\epsilon \frac{\eta_m(\eta_i - \eta_m)}{\eta_i + \eta_m} \frac{r_i^2}{r^2} \cos(2\theta)$$

where  $\eta_i$  is the viscosity of the inclusion (often taken to be 1000) and  $\eta_m$  is the viscosity of the background media (often taken to be 1).


One important observation with this benchmark is the fact that the velocity is not zero even far away from the inclusion, so that the analytical solution must be imposed on the sides. Also, because



of symmetry, it is often run on the top quadrant  $x > 0, y > 0$  with free slip imposed on the left and bottom boundaries.



Left: taken from Duretz *et al.* (2011) [352].

 **Relevant Literature:** Kaus and Podlachikov [682], Maierová [825], Deubelbeiss and Kaus [330], Beuchert and Podlachikov [86], Suckale, Nave, and Hager [1218], von Tschanner and Schmalholz [1328], de Montserrat, Morgan, and Hasenclever [322], Bangerth *et al.* [44], Liu and Tu [799] (2002); Kronbichler, Heister, and Bangerth [732] (2012); Gerya, May, and Duretz [452] (2013); Sevilla and Duretz [1150] (2023). [STONE](#) 07,

### 12.1.33 Simple shear heating

The domain is a  $L_x \times L_y$  Cartesian box. The velocity field  $\vec{v} = (L_y - y)y\vec{e}_x$  is prescribed on all boundaries, or simply prescribed everywhere in the domain. Temperature is set to  $T = 0$  everywhere in the domain at  $t = 0$ .

As we have seen in Section 2.6, the shear heating,  $\Phi$  is expressed as:

$$\Phi = 2\eta\dot{\epsilon}^d(\vec{v}) : \dot{\epsilon}^d(\vec{v}) \quad (12.169)$$

We have

$$\begin{aligned} \dot{\epsilon}_{xx}(\vec{v}) &= \partial_x u = 0 \\ \dot{\epsilon}_{yy}(\vec{v}) &= \partial_y v = 0 \\ \dot{\epsilon}_{xy}(\vec{v}) &= \frac{1}{2}(\partial_x v + \partial_y u) = \frac{1}{2}(L_y - 2y) \end{aligned} \quad (12.170)$$

We see that the flow is incompressible ( $\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} = 0$ ) so that  $\dot{\epsilon}^d(\vec{v}) = \dot{\epsilon}(\vec{v})$  and

$$\Phi(x, y) = 2\eta[\dot{\epsilon}_{xx}(\vec{v})^2 + \dot{\epsilon}_{yy}(\vec{v})^2 + 2\dot{\epsilon}_{xy}(\vec{v})^2] = 2\eta 2\frac{1}{4}(L_y - 2y)^2 = \eta(L_y - 2y)^2$$

We set  $\eta = 1$  and  $L_y = 1$  ( $L_x$  is actually irrelevant) so that  $\Phi(x, y) = (1 - 2y)^2$ .

The energy equation is given by

$$\rho C_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = k\Delta T + \Phi, \quad (12.171)$$

assuming that there is no conduction (i.e.  $k = 0$ ) and since  $\vec{\nabla}T \propto \vec{e}_y$  while  $\vec{v} \propto \vec{e}_x$  then it simplifies to

$$\frac{\partial T}{\partial t} = \Phi$$

where we have taken  $\rho = 1$  and  $C_p = 1$  for convenience. Since  $T(t = 0) = 0$ , we then have

$$T(t) = \Phi t$$

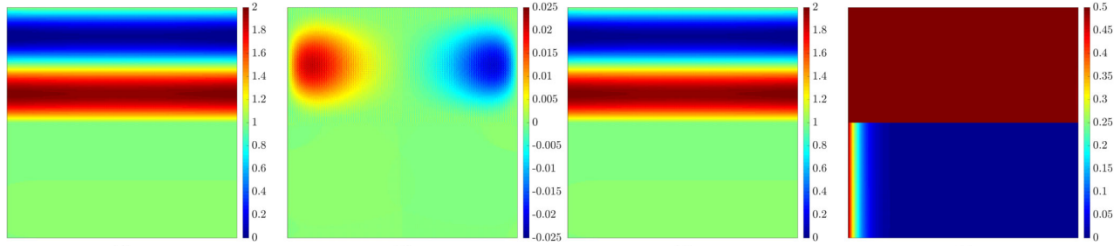
### 12.1.34 2D solution with nontrivial interface jump

benchmark\_jump2D.tex

This benchmark is featured in Sevilla and Duretz [1150] (2023) but originates in Wang and Khoo [1336] (2013). The domain is  $\Omega = [0, 2] \times [-0.5, 1.5]$  and the viscosity is given by  $\eta(x, y) = \eta_1$  if  $y \leq 0.5$  and  $\eta(x, y) = \eta_2$  if  $y > 0.5$ . The analytical solution is given by

$$\vec{v} = (1 - \exp(\lambda) \sin(2\pi y), 0) \quad p = \frac{1}{2} \exp(2\lambda x) \quad \text{with} \quad \lambda = \frac{1}{2\eta} - \sqrt{\frac{1}{4\eta^2} + 4\pi^2} \quad (12.172)$$

The pressure and the gradient tensor exhibit a discontinuity across the interface. Dirichlet boundary conditions are imposed in the whole boundary.



Taken from [1150]. From left to right:  $u$ ,  $v$ ,  $|\vec{v}|$ ,  $p$ , for  $\eta_1 = 1$  and  $\eta = 10^{-4}$ .

### 12.1.35 3D solution with nontrivial interface jump

This benchmark is taken from Sevilla and Duretz [1150] but is taken from Kirchhart, Gross, and Reusken [709] (2016). It is a three dimensional problem in  $[-1, 1]^3$  and the viscosity is given by  $\eta(x, y) = \eta_1$  if  $r \leq r_I$  and  $\eta(x, y) = \eta_2$  if  $r > r_I$  where  $r = |\vec{v}|_2$  and  $r_I = 2/3$ . The analytical solution is given by

$$\vec{v} = \alpha(r) \exp(-r^2)(-y, x, 0) \quad p = x^3 + \lambda(r) \quad (12.173)$$

where

$$\alpha(r) = 1/\eta_1 \quad \text{if } r \leq r_I \quad (12.174)$$

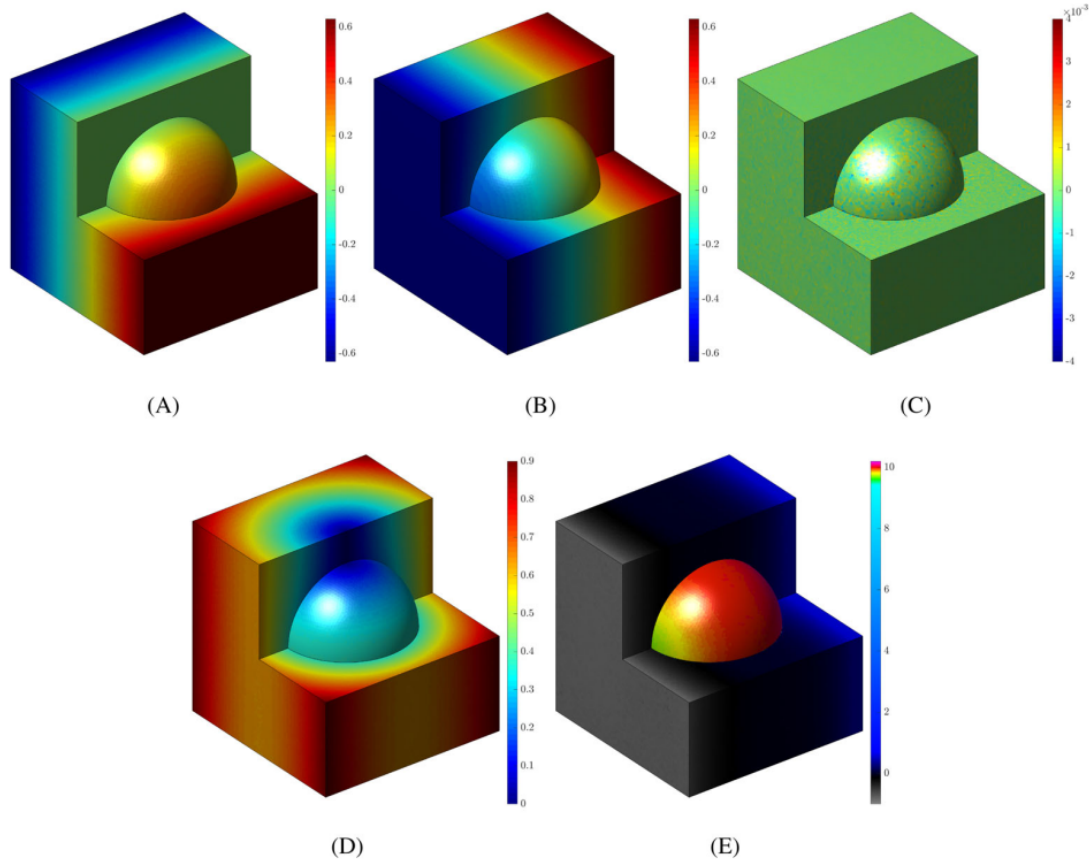
$$= 1/\eta_2 + (1/\eta_1 - 1/\eta_2) \exp(r^2 - r_I^2) \quad \text{if } r > r_I \quad (12.175)$$

and

$$\lambda(r) = 10 \quad \text{if } r \leq r_I \quad (12.176)$$

$$= 0 \quad \text{if } r > r_I \quad (12.177)$$

Dirichlet boundary conditions, corresponding to the analytical solution, are imposed in whole boundary of  $\Omega$ .



Taken from [1150]. From left to right:  $u$ ,  $v$ ,  $w$ ,  $|\vec{v}|$ ,  $p$ , for  $\eta_1 = 1$  and  $\eta = 10^2$ .

## 12.2 Geodynamical benchmarks

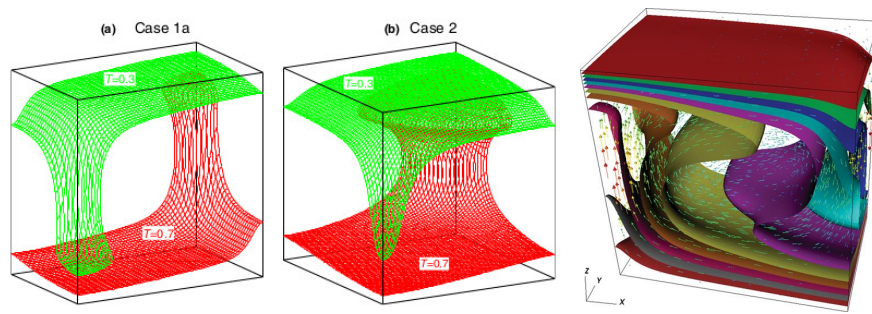
geodynamics\_benchmarks.tex

Some published numerical experiments have over time become benchmarks for other codes, while some others showcased comparisons between codes. Here is a short list of 'famous' benchmarks' in the computational geodynamics community.

- the 'plastic brick' (See section 12.2.3):  
Lemiale, Mühlhaus, Moresi, and Stafford [764], Kaus [679], Quinteros, Ramos, and Jacovkis [1029], Mishin [878], Mühlhaus, Shi, Olsen-Kettle, and Moresi [917], Maierová [825], Spiegelman, May, and Wilson [1187] (2016), Kaus et al. [683] (2016), Glerum, Thieulot, Fraters, Blom, and Spakman [467] (2018), Fraters, Bangerth, Thieulot, Glerum, and Spakman [415] (2019), Mishin, Vasilyev, and Gerya [879] (2022).
- indenter, punch problem (see Section 12.2.8):  
Vilotte *et al.* [1322, 1323, 1324], Hubert-Ferrari *et al.* (2003) [602], Fournier *et al.* (2004) [407], Thieulot *et al.* (2008) [1261], Gerbault (2012) [447], Glerum *et al.* (2018) [467], Stone 8.
- 2D Rayleigh-Benard convection (see Section 12.2.26).
- 2D Rayleigh-Benard convection, lateral heating, 30+ codes: Vahl Davis and Jones [1300] (1983).
- 2D Rayleigh-Benard convection with nonlinear rheology: Tosi et al. [1276] (2015), ASPECT manual [44], Trim, Butler, and Spiteri [1282] (2021), [STONE](#) 28, Davies, Kramer, Ghelichkhan, and Gibson [306] (2022), Sime and Wilson [1169] (2020), Candioti, Schmalholz, and Duretz [205] (2020).

- 2D Rayleigh-Benard laminar plumes, comparison of laboratory and numerical modeling : Vatteville *et al.* (2009) [1314]
- 2D Cartesian flow with extremely temperature-dependent viscosity: Moresi & Solomatov (1995) [903], Trim *et al.* (2021) [1282]
- 2D Rayleigh-Taylor convection/instability: Prosperetti (1981) [1019], Travis *et al.* (1990) [1278], Weinberg & Schmeling (1992) [1346], Poliakov & Podlachikov (1992) [1008], Ogawa (1993) [953], Conrad & Molnar (1997) [277], van Keken *et al.* (1997) [1309], de Smet *et al.* (2000) [323], Soboutia *et al.* (2001) [1176], Babeyko *et al.* (2002) [35], Tackley & King (2003) [1229], Bourgouin *et al.* (2006) [124], Davies, Davies, Hassan, Morgan, and Nithiarasu [307] (2007), Battaglia, Storti, and D’Elia [55] (2008), Deubelbeiss and Kaus [330] (2008), Quinteros *et al.* (2009) [1029], Samuel & Evonuk (2010) [1103], Suckale *et al.* (2010) [1218], Leng & Zhong (2011) [769], Mishin (2011) [878], Logg *et al.* (2012) [806], Maierova (2012) [825], Vynnytska *et al.* (2013) [1333], Choi *et al.* (2013) [238], Robey & Puckett (2019) [1079], Robey (2019) [1078], Fuchs & Schmeling [421], de Montserrat, Morgan, and Hasenclever [322] (2019), Louis-Napoléon, Gerbault, Bonometti, Thieulot, Martin, and Vanderhaeghe [811] (2020), Schuh-Senlis, Thieulot, Cupillard, and Caumon [1144] (2020), Mishin, Vasilyev, and Gerya [879] (2022), Burcet, Oliveira, Afonso, and Zlotnik [175] (2024), ASPECT manual [44].
- 3D Rayleigh-Taylor instability: Furuichi *et al.* (2008) [429], von Tscharner & Schmalholtz (2015) [1328]
- subduction problems: Spiegelman and Katz [1188], Schmeling *et al.* [1124], van Keken *et al.* [1311], Cerpa, Hassani, Gerbault, and Prévost [216], Glerum, Thieulot, Fraters, Blom, and Spakman [467], OzBench *et al.* [968], Sime and Wilson [1169].
- Benchmark of 3D numerical models of subduction against a laboratory experiment: Meriaux *et al.* (2018) [864]
- numerical sandbox [161, 161, 825, 164, 467]
- the Stokes sphere: Gale manual [744], ASPECT manual [44], in visco-plastic fluid: Liu *et al.* [794], Deglo de Besses *et al.* [85]. Finite deformation in and around a fluid sphere [1127, 291].
- the sinking block (sinker) Thieulot [1258] (2011), Cerpa, Hassani, Gerbault, and Prévost [216] (2014), Gerya [455] (2010), Gerya and Yuen [453] (2003), May and Moresi [846] (2008), Mishin [878] (2011), Furuichi, May, and Tackley [431] (2011), Maierová [825] (2012), Schuh-Senlis, Thieulot, Cupillard, and Caumon [1144] (2020), Mishin, Vasilyev, and Gerya [879] (2022). (see Section 12.2.6)
- multiple sinkers [848, 845, 261]
- Thin layer entrainment (see Section 12.2.28)
- 1D compression [899]
- 2D compressible Stokes flow problem [627, 1234, 770, 700, 802]
- 3D convection at infinite Prandtl number with modest viscosity variation: Busse *et al.* [193] (1994), Trompert and Hansen [1283] (1998), Kameyama, Kageyama, and Sato [668] (2005), O’Neill, Moresi, Müller, Albert, and Dufour [949] (2006), Kronbichler, Heister, and Bangerth

[732] (2012), Trim, Butler, and Spiteri [1282] (2021), Davies, Kramer, Ghelichkhan, and Gibson [306] (2022).



Left: Taken from Kameyama *et al.* (2005). a) Isothermal surfaces obtained for the benchmark calculations of stationary convections in Busse *et al.* . (1993). (a) Case 1a is for constant viscosity, while (b) Case 2 is for modestly temperature-dependent viscosity whose viscosity contrast is 20. The calculations were carried out with (a)  $64 \times 32 \times 64$  and (b)  $64 \times 64 \times 64$  mesh divisions. Right: Taken from Kronbichler *et al.* (2012).

- Numerical simulations of three-dimensional thermal convection in a fluid with strongly temperature-dependent viscosity: Ogawa *et al.* [954, 668]
- Free surface evolution: Cramer *et al.* (2012) [285], ASPECT manual [44], Schuh-Senlis *et al.* (2020) [1144]
- Love's problem: Becker & Bevis (2004) [64]
- Poiseuille flow: [400, 423, 1231] (see Section 12.2.1)
- Couette flow with temperature dependent viscosity [367, 322]
- Couette flow with shear heating [367]
- Poiseuille-Couette flow [421]
- Lid driven cavity [684, 239, 1130, 405, 459, 724, 116, 1375, 157, 381]
- Lid driven cavity with analytical solution (see Section 12.2.9)
- Lid driven cavity with nonlinear rheology [80, 1221]
- Wannier flow [1342, 1384, 216]
- bending of elastic plate/beam [216, 114, 1328, 367, 322, 899, 799]
- flexure of finite length elastic plate [238]
- thermal diffusion of half-cooling space (see Section 12.2.12)
- thermal diffusion of Gaussian distribution (see compgeo notes, elephant manual)
- stress build-up in Maxwell visco-elastic material [450, 238, 367, 322]
- plastic oedometer test [238]
- channel flow (nonlinear) [453, 825, 415, 455, 367]
- relaxation of sinusoidal interface [285, 1085]
- single layer visco-elastic folding [1117, 1328]

- Three-dimensional folding of an embedded viscous layer in pure shear [398]
- dam-break problem [888, 42, 788, 757, 584, 17, 494, 574, 55]
- hot blob problem [188, 431] (see Section 12.2.7)
- viscous(-elastic) flow around a cylinder in a channel (see Section 12.2.10)
- Sinking cylinder (2D Stokes sphere): appendix A of [114], [1340].
- Infinite plate with a circular hole [1386, 1032]
- Semi-infinite elastic half plane with a circular hole [1320]
- Slope stability for elasto-plastic materials [1032]
- Time-dependent flow in an annulus [440] (see Section 12.2.4)
- Convection in 2D-box [440] (see Section 12.2.5)
- Onset of convection [44]
- Polydiapirism [1346, 44]
- Slab detachment benchmark (see Section 12.2.14)
- 3D Hollow sphere Stokes flow benchmark:  
Thieulot (2017) [1256], Horbach *et al.* (2020) [591], Kramer *et al.* [730]
- Axisymmetric hollow sphere compressible Stokes flow benchmark:  
Machetel & yuen (1989) [819]
- Annulus benchmark [44], [1002]
- Viscosity grooves benchmark [44]
- Latent heat benchmark [44]
- Layered flow with viscosity contrast [44] (see Section 12.2.15)
- Brittle thrust wedges benchmark [164, 44]
- mantle convection in 3D spherical shell: Ratcliff, Schubert, and Zebib [1046] (1996), Iwase (1996) [628], Zhong *et al.* (2000) [1414], Yoshida & Kageyama [1387], Stemmer *et al.* (2006) [1205], Choblet *et al.* (2007) [236], Zhong *et al.* (2008) [1412], Kameyama *et al.* (2008) [665], Wright *et al.* (2010) [1370], Davies *et al.* (2013) [308], Burstedde *et al.* (2013) [189], Arrial *et al.* (2014) [29], Liu & King (2019) [801], Trim *et al.* (2021) [1282]
- Heat flow around a cylinder (see Section 12.2.11)
- Laplace equation on a semi infinite plate (see Section 12.2.13)
- 2D Stokes flow over cavity: Popov & Makeev (2014) [1014].
- fractal networks of shear bands: Poliakov & Herrmann (1994) [1009]
- Square plate with a crack subjected to a horizontal tensile traction [799]
- Analytical solution for solitary porosity waves: Connolly & Podlachikov (2015) [276]



- Analytical solution for solitary wave of magma: Dannberg & Heister (2016) [302] and refs therein
- Stokes flow caused by the motion of a rigid sphere close to a viscous interface: Danov *et al.* (1998) [304]
- Deformation caused by a closed vertical volcanic pipe [113]
- Mantle convection with reversing mobile plates [718]
- A comparison of mantle convection models featuring plates [1197]
- Uniform strip load on elastic material (see Section ??)
- Linear Stability Analysis for Thermal Convections in Spherical Shells [1392]
- Channel flow: Mancktelow (2008) [832]
- Viscous half-space loading [551]
- Nakiboglu and Lambeck (1982) has cylindrical load on variety of rheologies [925]
- Jull and McKenzie (1996) [659] parabolic load on viscoelastic half-space (and melt fractions)
- Squeezing flow between moving parallel plates [512]
- generalized half-plane and half-space Cerruti [947, 1417]
- Analytical Solutions of Displacements Produced by spherically-shaped Internal Overpressure [448]
- Sagging viscous bridge [1216]
- Deformation around a terminating fault in a viscous medium [47]
- Stress distribution in elastic sphere under equal and opposite loads [1210]
- Flow of a Power Law Fluid Through a Tube - page 87 of Macosko [821]

### 12.2.1 Poiseuille flow

We consider a two-dimensional channel in the  $x, y$  plane. The walls are at  $y = 0$  and  $y = H$  with no-slip boundary conditions. In the absence of gravity, the Stokes equation simplify to

$$-\frac{\partial p}{\partial x} + \frac{\partial}{\partial y}(2\eta_0 \dot{\epsilon}_{xy}) = 0 \quad \text{and} \quad \dot{\epsilon}_{xy} = \frac{1}{2} \frac{\partial u}{\partial y} \quad (12.178)$$

where we assume the velocity  $\vec{v} = (u(y), 0)$ . In the case of a Newtonian fluid, the analytical solution is known and the velocity profile is a parabola with zero velocity on the walls and maximum velocity in the middle.

Eq. (12.178) must then be solved

$$\frac{\partial p}{\partial x} = \frac{\partial}{\partial y} \left( 2\eta_0 \frac{1}{2} \frac{\partial u}{\partial y} \right) = \eta_0 \frac{\partial^2 u}{\partial y^2} \quad (12.179)$$

We pose  $\Pi = \frac{\partial p}{\partial x} < 0$ , i.e. there is more pressure applied to the left than to the right of the channel. We then must solve:

$$\frac{\partial^2 u}{\partial y^2} = \frac{\Pi}{\eta_0}$$

The solution is then of the form

$$u(y) = \frac{1}{2} \frac{\Pi}{\eta_0} y^2 + 2ay + b$$

and

$$\dot{\epsilon}_{xy} = \frac{1}{2} \frac{\Pi}{\eta_0} y + a$$

We will now determine  $a$  and  $b$ .

The velocity must be zero at  $y = 0$  and  $y = H$  so

$$u(y = 0) = b = 0$$

and

$$u(y = H) = \frac{1}{2} \frac{\Pi}{\eta_0} H^2 + 2aH = 0$$

or,

$$2a = -\frac{1}{2} \frac{\Pi}{\eta_0} H$$

so

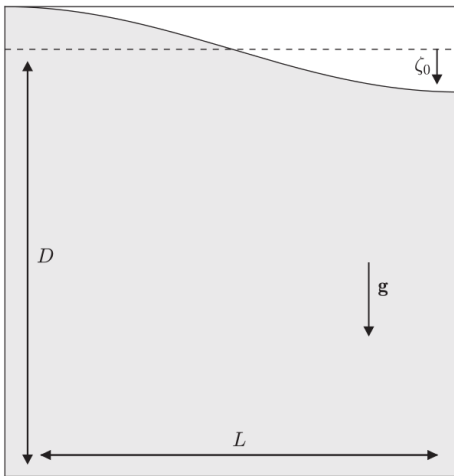
$$u(y) = \frac{1}{2} \frac{\Pi}{\eta_0} (y^2 - yH) \quad (12.180)$$

and

$$\dot{\epsilon}_{xy} = \frac{1}{2} \frac{\Pi}{\eta_0} \left( y - \frac{H}{2} \right) \quad (12.181)$$

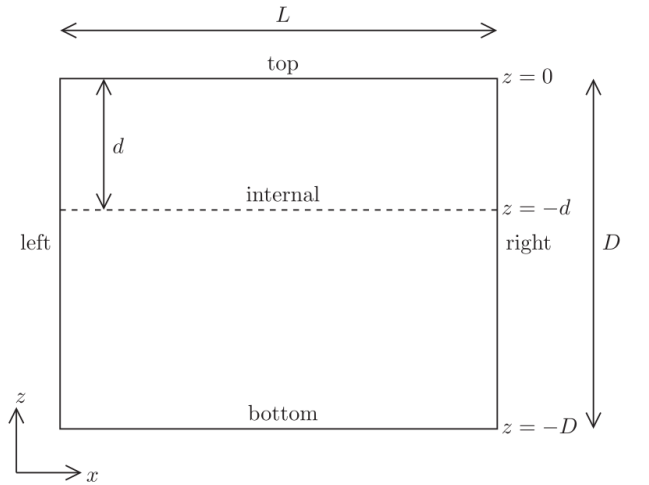
## 12.2.2 Relaxation of sinusoidal topography

Following Kramer *et al.* [731, Section 3.1.1] and [1085] the benchmark consists of the relaxation of surface topography in a two-dimensional Cartesian box with an isoviscous fluid. Free slip boundary conditions are imposed on the sides and bottom of the domain. The setup is as follows:



Taken from [1085]. Setup for the free surface relaxation benchmark.

For the tests  $\rho = \eta = g = L = D = 1$  and  $\xi_0 = 0.005$ .



Taken from [731].  $D = 3 \cdot 10^6, \eta = 10^{21}, \rho = 4500, g = 10, \xi_0 = 10^3 \text{ m}$ ,

and  $L = D/4, D/2, D, 2D, 4D$ .




and the infinitesimal sinusoidal perturbations to the free surface is given by

$$\xi(x, t = 0) = \xi_0 \cos\left(\frac{2\pi nx}{L}\right)$$

where  $n$  is a wavenumber which is an integer multiple of  $1/2$  (taken to be  $1/2$  exactly in both cases).

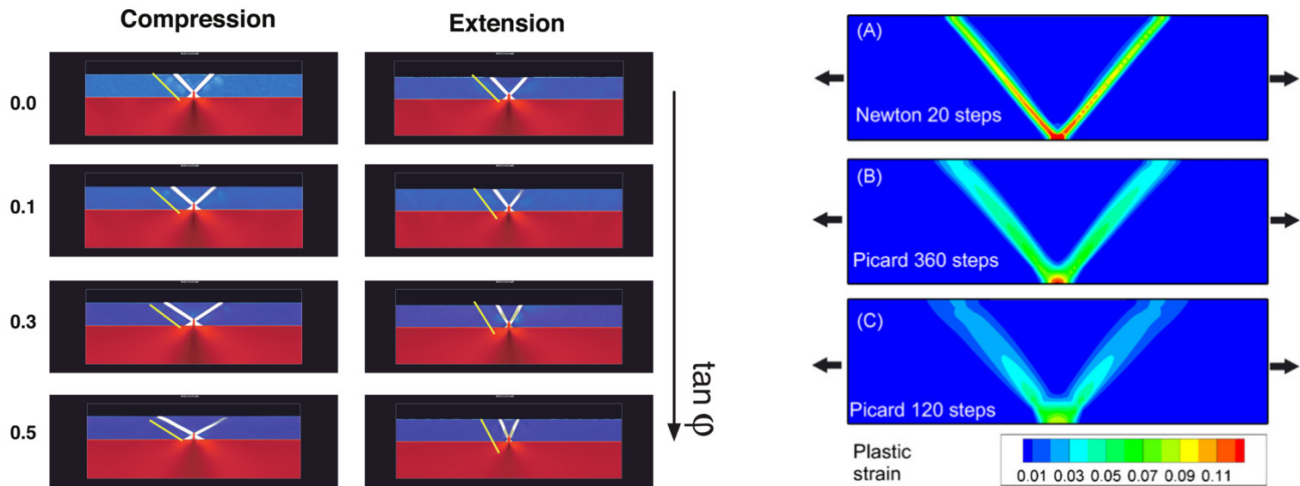
### 12.2.3 the plastic brick

 **Relevant Literature** Hansen [530], Moresi, Mühlhaus, Lemiale, and May [896], Lemiale, Mühlhaus, Moresi, and Stafford [764], [679, 367, 1029, 878, 825, 1187, 467, 415, 44] Davies, Kramer, Ghelichkhan, and Gibson [306] (2022), Mishin, Vasilyev, and Gerya [879] (2022).

Pretty much all of the brick-type (elasto-)visco-plastic experiments in the literature introduce a weak seed at the bottom of the domain to seed deformation (the shear bands will ultimately stem from it). Dimensioned and dimensionless experiments have been carried out, with or without elastic behaviour, with or without adaptive mesh refinement, with first order and second order quadrilateral elements or Taylor-Hood triangles, with or without Newton algorithm, in extension and compression, with or without time-stepping, with or without viscous lower layer.

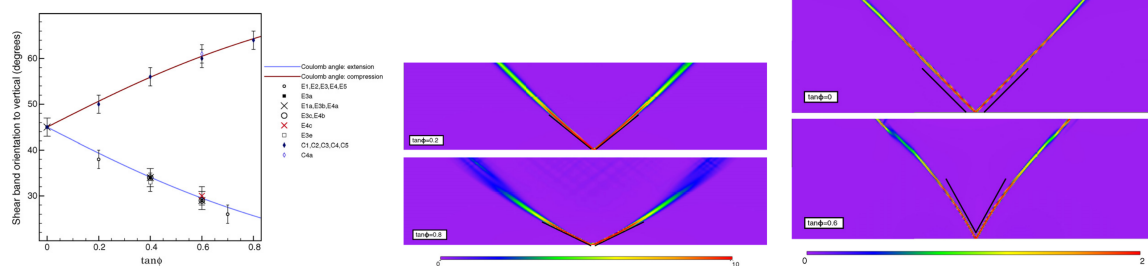


Moresi & Mühlhaus, 2006 [895]

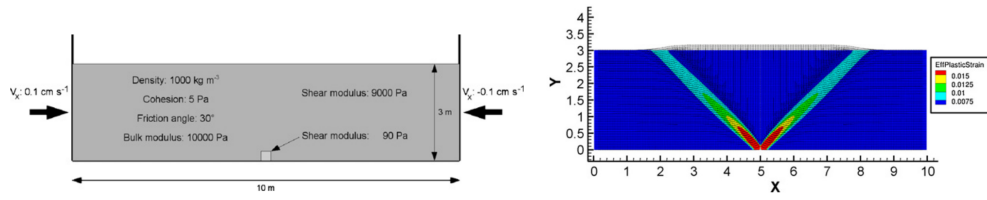


Moresi *et al.* , 2007 [896]

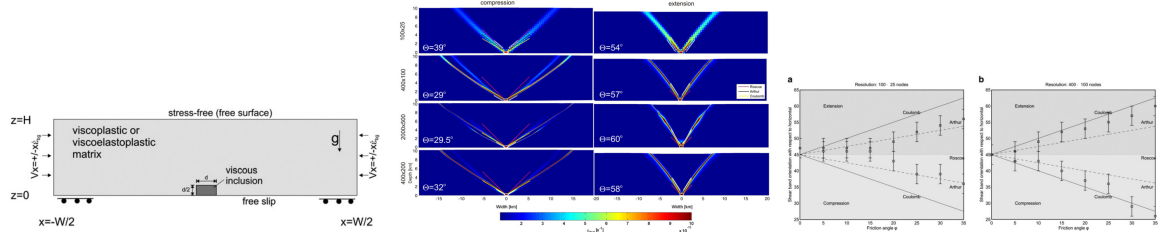
Popov *et al.* , 2008 [1011]



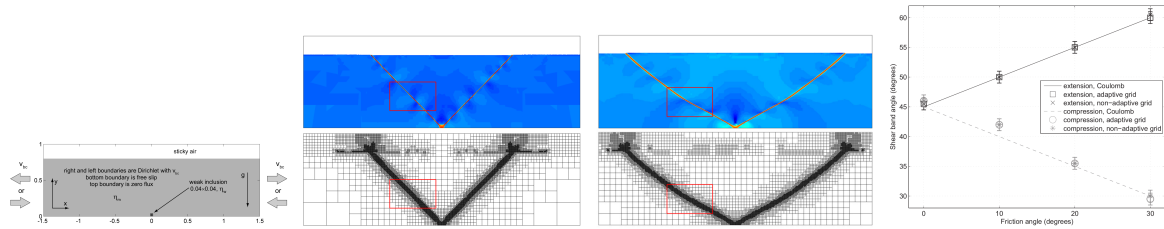
Lemiale *et al.* , 2008 [764]



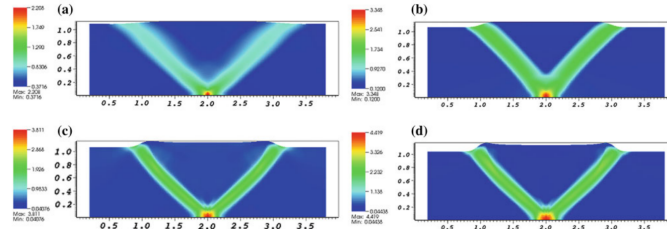
Quinteros *et al.* , 2009 [1029]



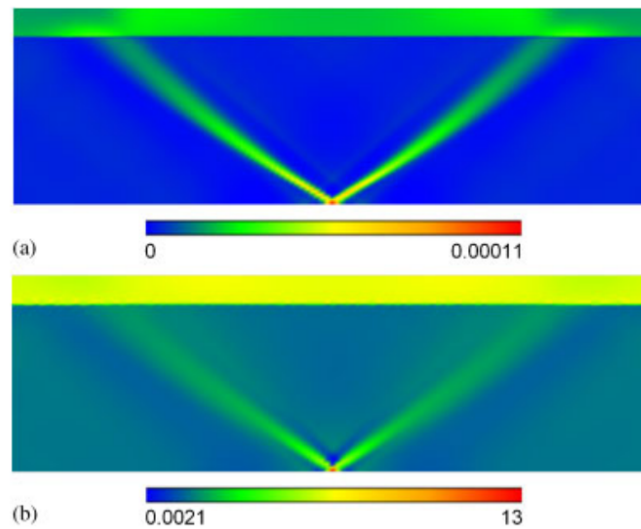
Kaus, 2010 [679]



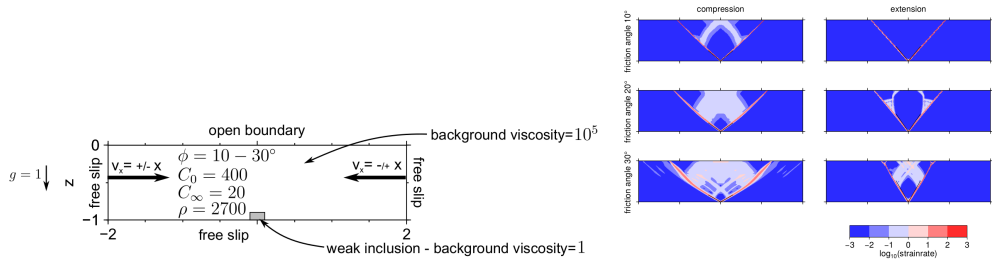
Mishin, phd thesis, 2011 [878]



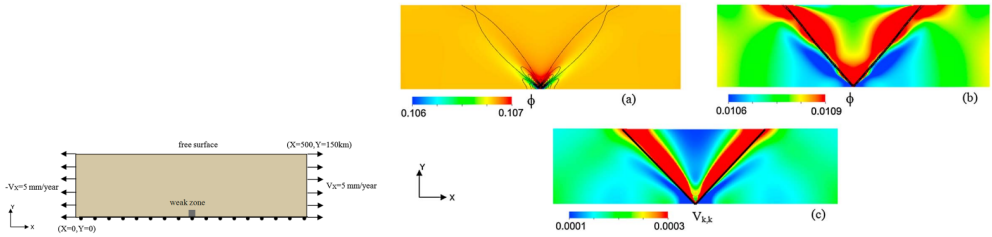
Mühlhaus *et al.* , 2011 [917].



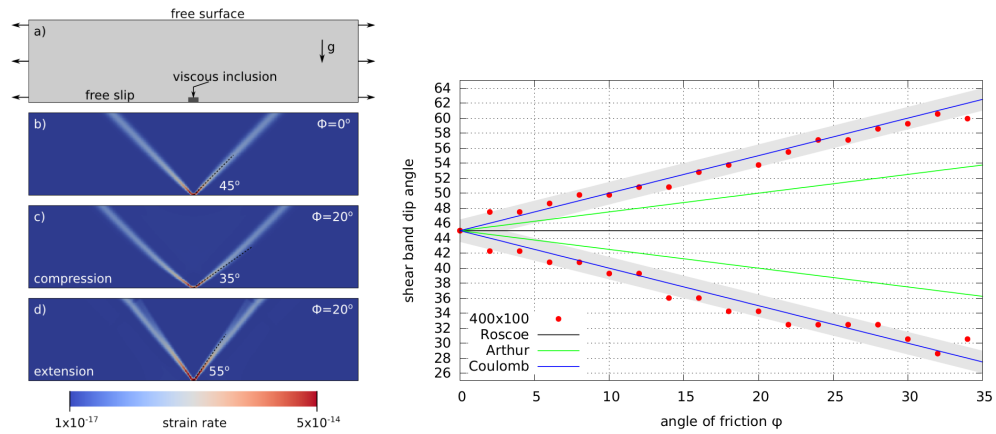
Lemiale *et al.* , 2011 [765].



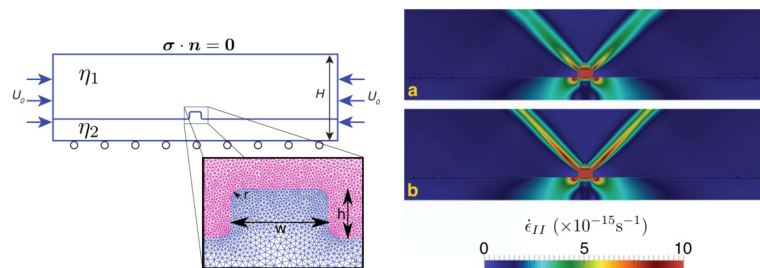
Maierova, phd thesis, 2012 [825]



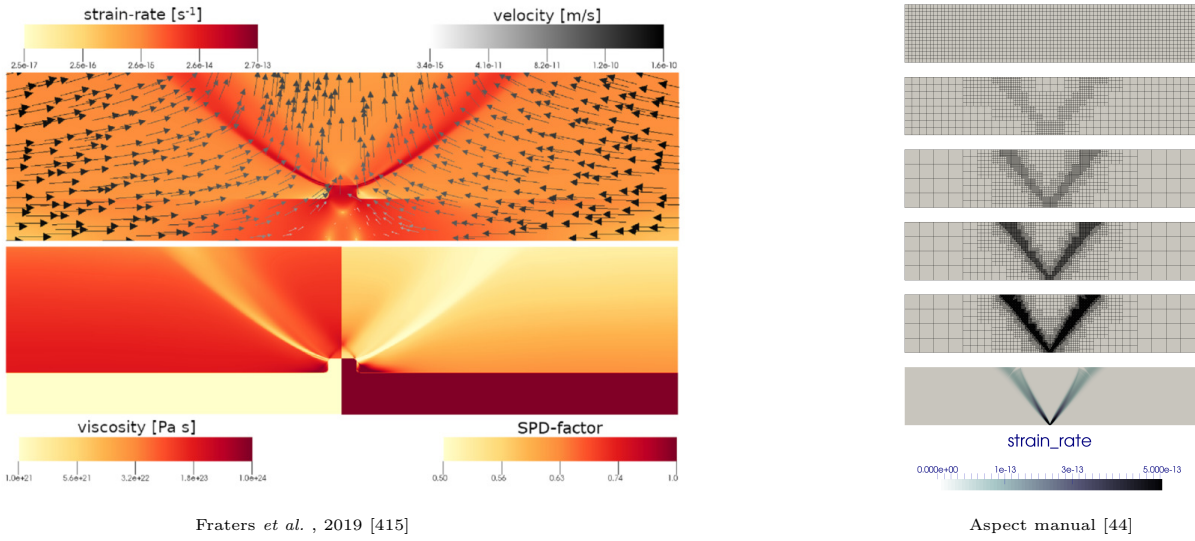
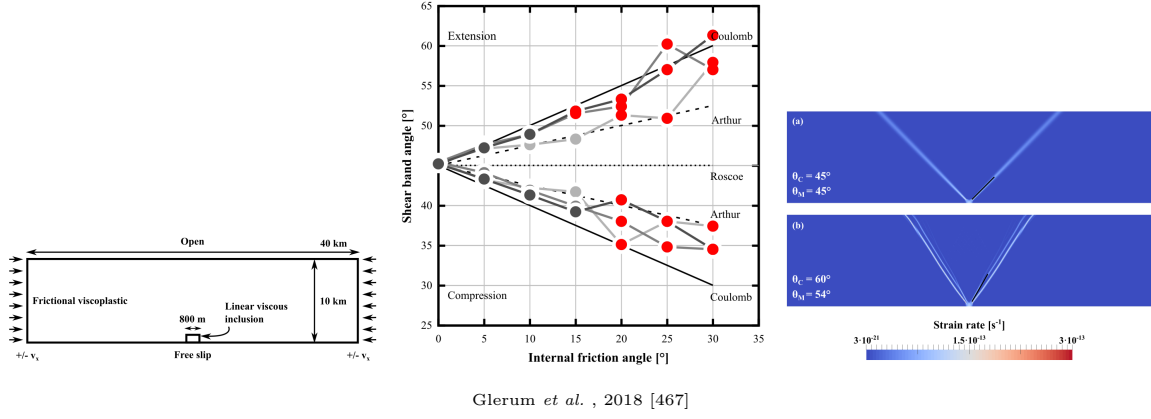
Mohajeri *et al.*, 2013 [887].



Thieulot, 2014 [1257].



Spiegelman *et al.*, 2016 [1187]



## 12.2.4 Time-dependent benchmark in an annulus

This benchmark is presented in Gassmöller *et al.* [440]. The domain is a 2D annulus with inner and outer radii  $R_1 = 1$  and  $R_2 = 2$ , respectively. In this situation, the incompressible isothermal Stokes equations and their solution can be expressed in a cylindrical coordinate system in terms of the radius  $r$  and the azimuthal angle  $\theta$ . The viscosity is set to  $\eta = 1$ , and the density is given by

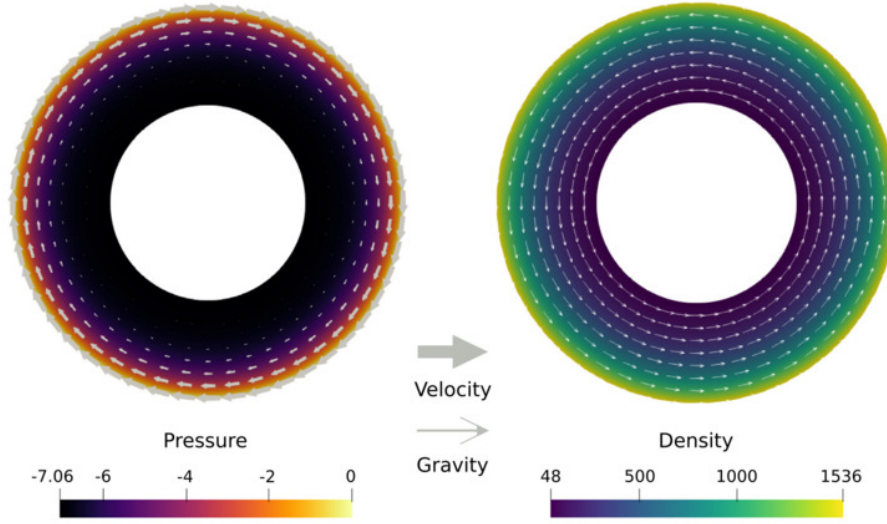
$$\rho(r, \theta) = 48r^5 \quad (12.182)$$

The gravity vector is set to

$$\vec{g}(r, \theta) = \frac{r^3}{384} \vec{e}_r + \vec{e}_\theta \quad (12.183)$$

Note that this gravity vector is not the gradient of a gravity potential and consequently not physical. The Stokes system can then be solved using a separation of variables approach and yields

$$\vec{v} = -r^7 \vec{e}_\theta \quad p(r, \theta) = \frac{r^9}{72} - \frac{512}{72} \quad (12.184)$$



Taken from [440]

Rather importantly, this benchmark was arrived at by means of a stream function (see Section ??)  $\psi(r, \theta) = F(r)G(\theta)$  with  $F(r) = r^8/8$  and  $G(\theta) = 1$ .

## 12.2.5 Convection in 2D-box

We start from the following stream function (see Section ??):

$$\psi(x, y) = \frac{1}{\pi} \sin \pi x \sin \pi y \quad (12.185)$$

which yields:

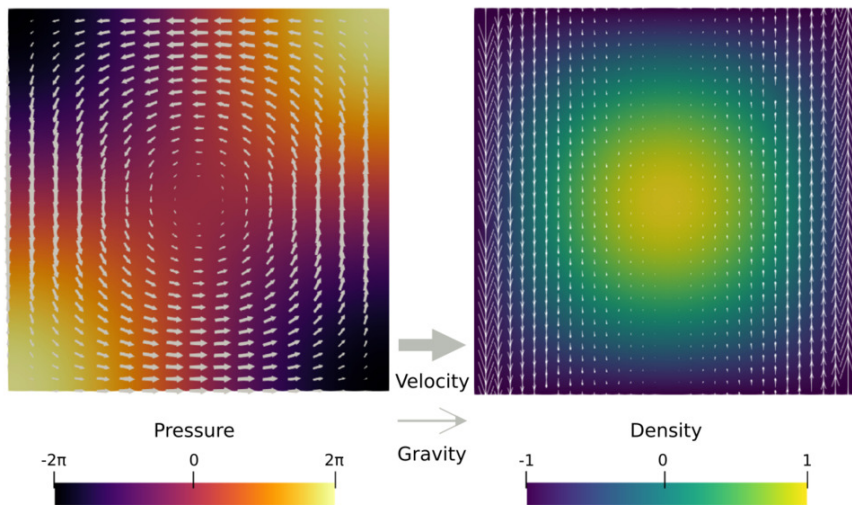
$$\begin{aligned} u(x, y) &= \frac{\partial \psi}{\partial y} = \sin \pi x \cos \pi y \\ v(x, y) &= -\frac{\partial \psi}{\partial x} = -\cos \pi x \sin \pi y \end{aligned} \quad (12.186)$$

The pressure field is

$$p(x, y) = 2\pi \cos(\pi x) \cos(\pi y) \quad (12.187)$$

with

$$\rho(x, y) = \sin(\pi x) \sin(\pi y) \quad g_y = -4\pi^2 \frac{\cos(\pi x)}{\sin(\pi x)} \quad (12.188)$$



Taken from [440]

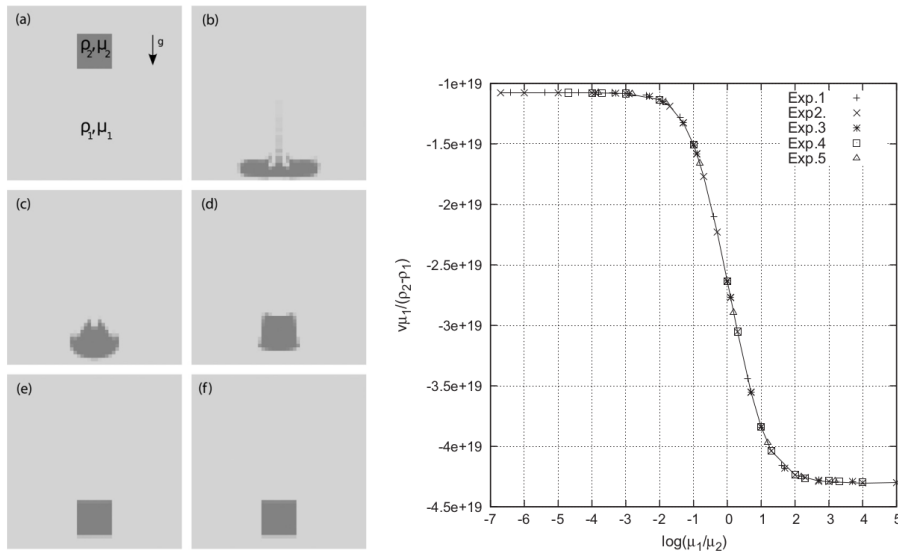
$$\begin{aligned}
v_{rms} &= \sqrt{\frac{1}{L_x L_y} \int_0^1 \int_0^1 (u^2 + v^2) dx dy} \\
&= \sqrt{\int_0^1 \int_0^1 (\sin^2(\pi x) \cos^2(\pi y) + \cos^2(\pi x) \sin^2(\pi y)) dx dy} \\
&= \sqrt{\int_0^1 \sin^2(\pi x) dx \cdot \int_0^1 \cos^2(\pi y) dy + \int_0^1 \cos^2(\pi x) dx \cdot \int_0^1 \sin^2(\pi y) dy} \\
&= \sqrt{\frac{1}{2} \frac{1}{2} + \frac{1}{2} \frac{1}{2}} \\
&= \frac{\sqrt{2}}{2} \\
&\simeq 0.70711...
\end{aligned} \tag{12.189}$$

### 12.2.6 The sinker problem

This experiment is not a benchmark *stricto sensu* since there is no analytical solution. However, it is widely used in the technical literature because of its simple setup and since it allows to test solving strategies. Also, it can conveniently be carried out in both two and three dimensions.

**In two dimensions** The time dependent version of the experiment is for instance to be found in Gerya [455] and the same is repeated in Thieulot [1258].

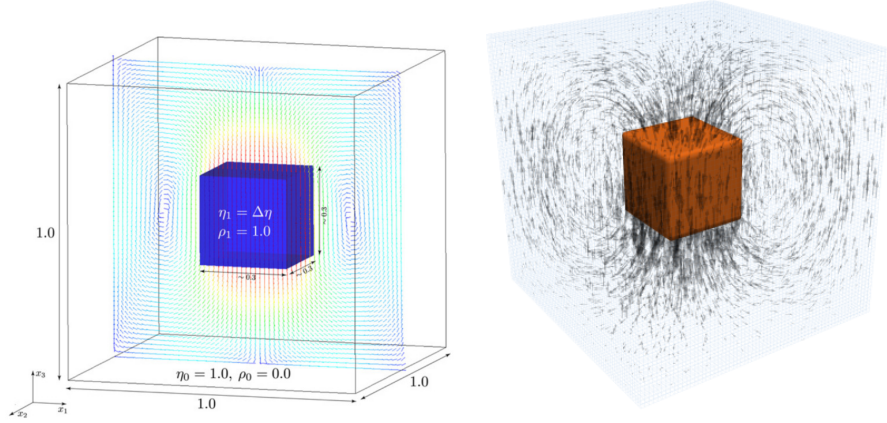
This simple benchmark provides challenging numerical experiments dealing with large viscosity variations within the simulation domain. It consists of a bulk of fluid 1 ( $\eta_1, \rho_1$ ) in which a block of fluid 2 ( $\eta_2, \rho_2$ ) falls under its own weight. The domain is a square of size  $L_x = L_y = 500\text{km}$  and the block is initially centred at point  $(x = 250\text{km}, y = 400\text{km})$  with size  $100 \times 100\text{km}$ . Free slip boundary conditions are imposed on all sides of the domain. In [1258] five experiments have been conducted:  $\eta_1 = 10^{20}\text{Pa.s}$ ,  $\rho_2 = 3220\text{kg m}^{-3}$ ;  $\eta_1 = 10^{21}\text{Pa.s}$ ,  $\rho_2 = 3300\text{kg m}^{-3}$ ;  $\eta_1 = 10^{22}\text{Pa.s}$ ,  $\rho_2 = 6600\text{kg m}^{-3}$ ;  $\eta_1 = 10^{23}\text{Pa.s}$ ,  $\rho_2 = 3300\text{kg m}^{-3}$ ;  $\eta_1 = 10^{24}\text{Pa.s}$ ,  $\rho_2 = 9900\text{kg m}^{-3}$ ; while in all experiments the density of the surrounding fluid is  $\rho_1 = 3200\text{kg m}^{-3}$  and the viscosity of the block is varied between  $10^{19}$  and  $5 \cdot 10^{27}\text{Pa.s}$ .





Left:  $\eta_1 = 10^{21}$  Pa s,  $\rho_2 = 3300 \text{ kg m}^{-3}$ . (a) Initial setup; (b)  $\eta_1 = 10^{21}$  Pa s at time  $t = 10$  Myrs; (c)  $\eta_1 = 10^{22}$  Pa s at time  $t = 20$  Myrs; (d)  $\eta_1 = 10^{23}$  Pa s at time  $t = 20$  Myrs; (e)  $\eta_1 = 10^{25}$  Pa s at time  $t = 20$  Myrs; (f)  $\eta_1 = 10^{27}$  Pa s at time  $t = 20$  Myrs. Right: Velocity measurements as a function of the viscosity contrast between surrounding medium and block for all experiments. Taken from [1258]

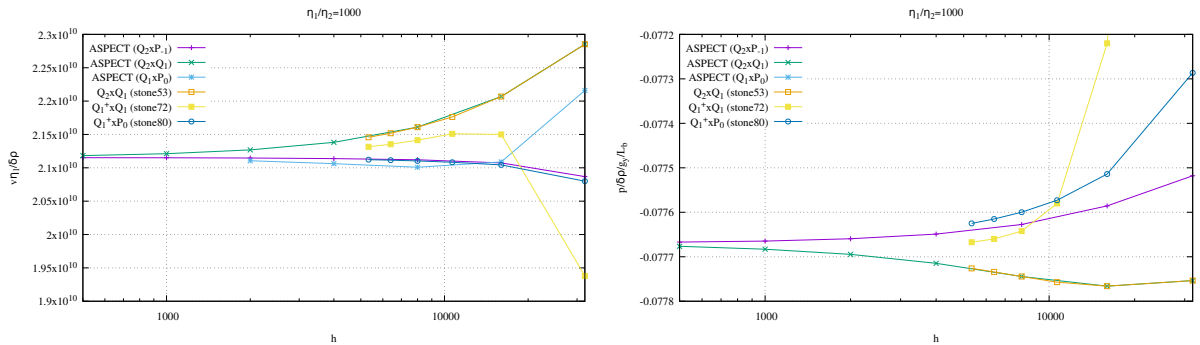
**In three dimensions** Let us look at the sinker experiment from Furuichi *et al.* [431]: The domain is the unit box the origin at the center of the box. A cube with a viscosity  $\eta_1 = \Delta\eta$  and density  $\rho_1 = 1$  was placed at the middle of the domain defined by  $-0.15 \leq x, y, z \leq 0.15$ . The material surrounding the cube has the properties  $\eta_0 = 1$  and  $\rho_0 = 0$ . The body force of the momentum equation was taken as  $(0, 0, -\rho g)$  with  $g = 1$ . Along all walls on the domain, free-slip boundary conditions were employed.



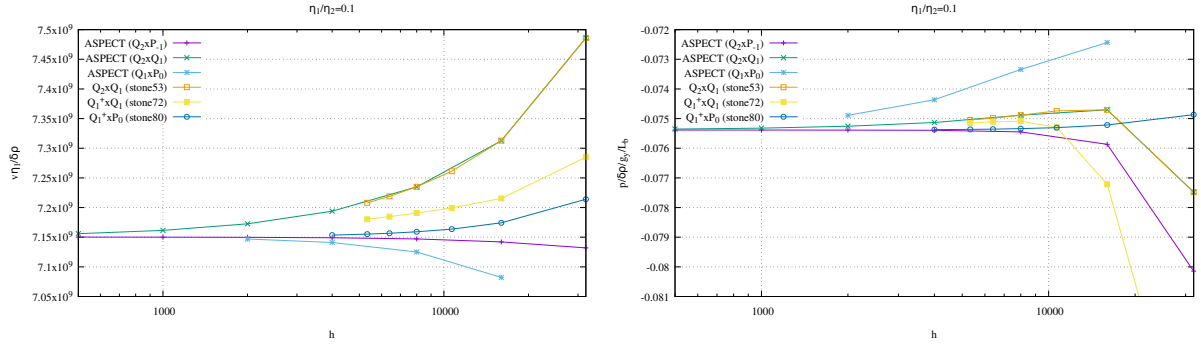
Left: Simulation setup for the 3D falling block (SINKER) problem. The vectors represent computed flow. Taken from Furuichi, May, and Tackley [431] (2011). Right: same experiment in Sanan, May, Mills, et al. [1107] (2022).

**Sinking block results for multiple elements** The setup is slightly altered: the domain is  $512 \times 512 \text{ km}$ . The block has size  $L_b \times L_b = 128 \times 128 \text{ km}$ , and is centered on  $(L_x/2, 3L_y/4)$ . Free slip on all sides. Pressure is volume normalised.  $|g_y| = 10$ . This benchmark is part of ASPECT, and can therefore be run with  $Q_2 \times Q_1$ ,  $Q_2 \times P_{-1}$  and  $Q_1 \times P_0$  elements (although the solver does not converge for the latter at high resolutions). Velocity and pressure are measured in the middle of the block (in the case of the  $Q_1^+ \times P_0$  element the projected pressure  $q$  on the  $Q_1$  is used).

As above, the quantity  $|v|\eta_1/\delta\rho$  velocity is considered, but this time plotted as a function of the resolution for a fixed  $\eta_2$  and for various element types. The quantity  $p/\delta\rho/L_b$  is also plotted:

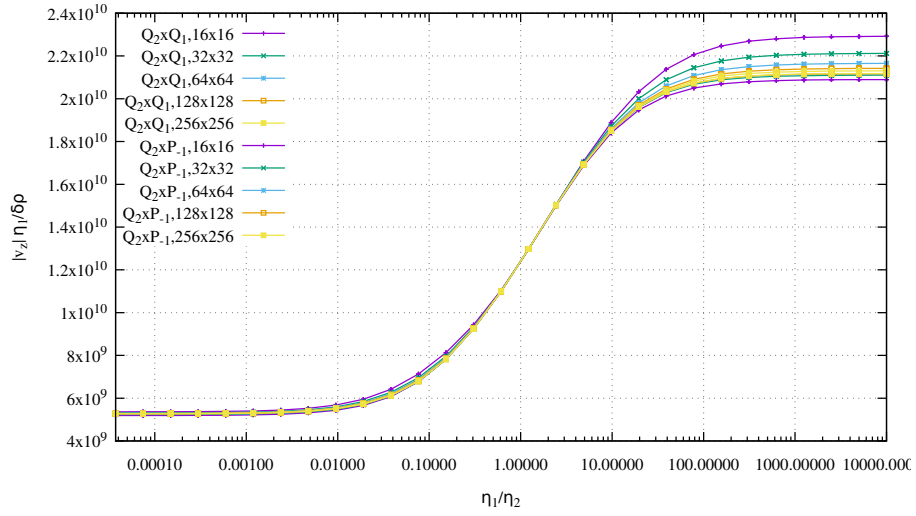


Results for  $\rho_1 = 0$ ,  $\rho_2 = \delta\rho = 8$ ,  $\eta_1 = 10^{21}$  and  $\eta_2 = 10^{18}$



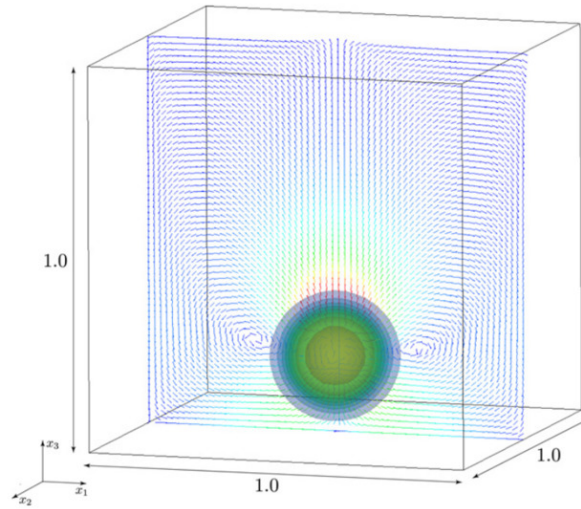
Results for  $\rho_1 = 0$ ,  $\rho_2 = \delta\rho = 8$ ,  $\eta_1 = 10^{21}$  and  $\eta_2 = 10^{22}$

Results obtained with ASPECT with  $\rho_1 = 3200$  and  $\rho_2 = \rho_1 + \delta\rho$  are shown here:



### 12.2.7 The hot blob problem

This is a very similar setup as the 3D sinker from the same authors with higher but more diffusive variation of viscosity. The body force is given by  $(0, 0, \beta T)$  and where the temperature field  $T$  is defined by  $T = \exp(-\gamma(x^2 + y^2 + (z - 0.3)^2))$  with the constant parameters  $\beta = 10^6$  and  $\gamma = 200$ . The temperature-dependent viscosity  $\eta = \exp(-\alpha T)$  is employed with the parameter for viscosity contrast  $\alpha$ .



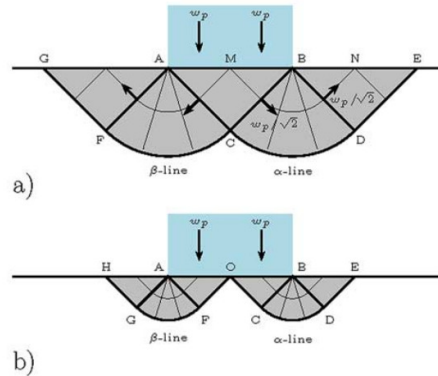
Simulation setting of BLOB problem. Isosurface and vectors represent temperature field and computed flow respectively. Taken from [431]



## 12.2.8 The punch/indenter problem in 2D

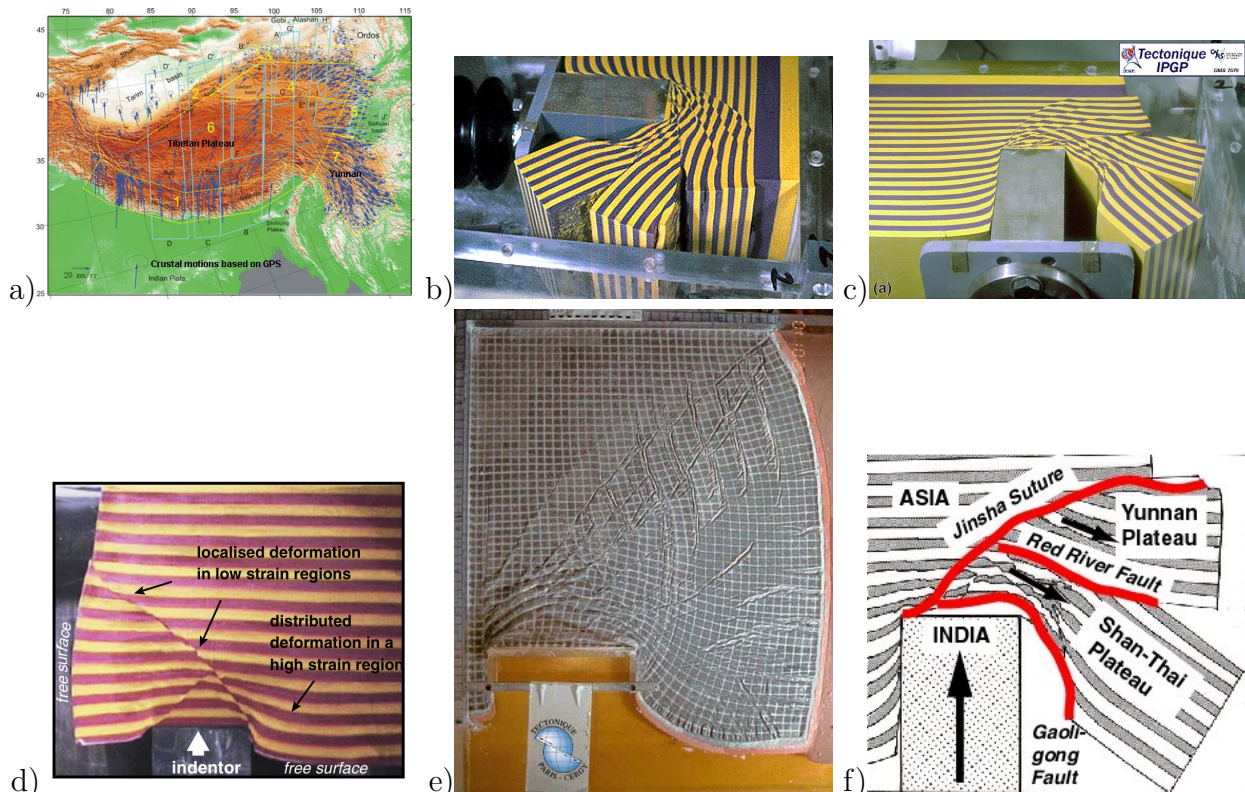
The punch benchmark is one of the few boundary value problems involving plastic solids for which there exists an exact solution. Such solutions are usually either for highly simplified geometries (spherical or axial symmetry, for instance) or simplified material models (such as rigid plastic solids) [661].

In this experiment, a rigid punch indents a rigid plastic half space; the slip line field theory gives exact solutions as shown hereunder:



Two-dimensional rigid punch indenting a rigid plastic half space. (a) Prandtl's rigid plastic solution; (b) Hill's solution. Taken from [1261]

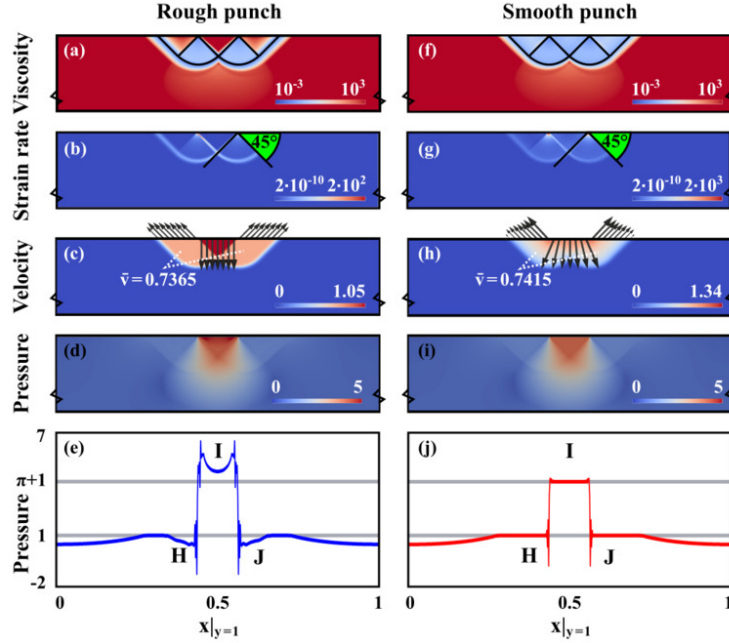
The plane strain formulation of the equations and the detailed solution to the problem were derived in the Appendix of [1261] and are also presented in [449] and in [125, Chapt.6]. The two dimensional punch problem has been extensively studied numerically for the past 40 years [1427, 1015, 1426, 1429, 253, 252, 611, 1393, 160, 1031, 467] and has been used to draw a parallel with the tectonics of eastern China in the context of the India-Eurasia collision [1238, 890, 374] or the European Alps [1054]. It is also worth noting that it has been carried out in one form or another in series of analogue modelling articles concerning the same region, with a rigid indenter colliding with a rheologically stratified lithosphere [1237, 990, 320, 656].



b,c) Model by Tapponnier *et al.* (1982) [1237] or Peltzer & Tapponnier (1988) [990]. Not sure about source for other figures.

Numerically, the one-time step punch experiment is performed on a two-dimensional domain of purely plastic von Mises material. Given that the von Mises rheology yield criterion does not depend on pressure, the density of the material and/or the gravity vector is set to zero. Sides are set to free slip boundary conditions, the bottom to no slip, while a vertical velocity  $(0, -v_p)$  is prescribed at the top boundary for nodes whose  $x$  coordinate is within  $[L_x/2 - \delta/2, L_x/2 + \delta/2]$ .

The analytical solution predicts that the angle of the shear bands stemming from the sides of the punch is  $\pi/4$ , that the pressure right under the punch is  $1 + \pi$ , and that the velocity of the rigid blocks on each side of the punch is  $v_p/\sqrt{2}$  (this is simply explained by invoking conservation of mass).



The punch benchmark results after 500 nonlinear iterations for a rough punch (left column) and a smooth punch (right column). (a,f) Viscosity field with analytical slip lines. (b,g) Strain rate norm  $\dot{\epsilon}_e$  with measured shear band angles. (c,h) Velocity magnitude with velocity vectors along the surface of the domain. (d,i) Pressure field. (e,j) Pressure along the surface of the domain (colored line) and analytical solution values  $\pi + 1$  and 1 (grey lines). Taken from [467]

**Remark.** This benchmark is often mentioned or used in the context of bearing capacity, footings, limit state design/analysis [870, 1418, 477, 478, 758, 1161].

### 12.2.9 Driven cavity with analytical solution

This comes from Elman *et al.* [371](section 3.1.4)<sup>6</sup>. The velocity is prescribed to be

$$\vec{v} = (2y(1 - x^2); -2x(1 - y^2))$$

on the domain  $\Omega = [-1 : 1] \times [-1 : 1]$ . The strainrate tensor is given by:

$$\dot{\epsilon} = \begin{pmatrix} -4xy & -x^2 + y^2 \\ -x^2 + y^2 & 4xy \end{pmatrix}$$

The Stokes equation is then:

$$-\frac{\partial p}{\partial x} + 2\eta(-4y + 2y) = \rho g_x \quad (12.190)$$

$$-\frac{\partial p}{\partial y} + 2\eta(-2x + 4x) = \rho g_y \quad (12.191)$$

<sup>6</sup>actually, not?

where we assume the viscosity  $\eta = 1$  to be constant in space. Assuming  $g_x = 0$ , the first equation is

$$\frac{\partial p}{\partial x} = -4y$$

i.e.

$$p(x, y) = -4yx + f(y)$$

Inserting this in the second equation:

$$4x - f'(y) + 4x = \rho g_y$$

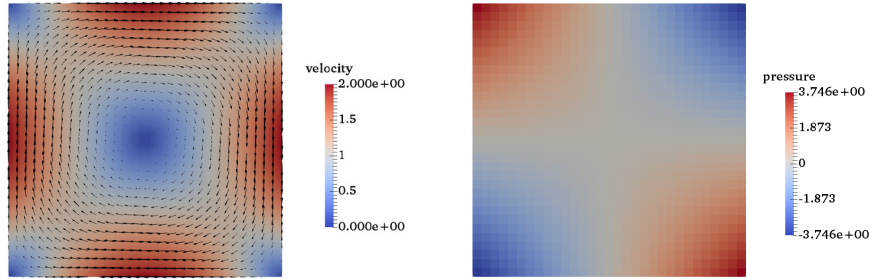
or,

$$-f'(y) + 8x = \rho g_y$$

Assuming  $g_y = -1$ , we get  $\rho = -8x$  and then  $f'(y) = 0$  so  $f(y) = C$  where  $C$  is a constant. Finally the pressure is given by:

$$p(x, y) = -4yx + C$$

We add the following requirement:  $\int_{\Omega} p(x, y) d\Omega = 0$  so that  $C = 0$ .



$$\begin{aligned}
 v_{rms}^2 &= \frac{1}{\Omega} \int_{\Omega} (u^2 + v^2) d\Omega \\
 &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} (u^2 + v^2) dx dy \\
 &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} [4y^2(1-x^2)^2 + 4x^2(1-y^2)^2] dx dy \\
 &= \int_{-1}^{+1} \int_{-1}^{+1} [y^2(1-x^2)^2 + x^2(1-y^2)^2] dx dy \\
 &= \int_{-1}^{+1} \int_{-1}^{+1} [y^2(1-2x+x^2) + x^2(1-2y+y^2)] dx dy \\
 &= \int_{-1}^{+1} \int_{-1}^{+1} y^2 dx dy + \int_{-1}^{+1} \int_{-1}^{+1} 2x^2 y^2 dx dy + \int_{-1}^{+1} \int_{-1}^{+1} x^2 dx dy \\
 &= 2\frac{2}{3} + 2\frac{2}{3}\frac{2}{3} + 2\frac{2}{3} \\
 &= \frac{32}{9} \\
 &= 3.5555
 \end{aligned} \tag{12.192}$$

We can reformulate the benchmark in the unit square  $\Omega = [0 : 1] \times [0 : 1]$ .

$$\begin{aligned}u(x, y) &= (2y - 1)x(1 - x) \\v(x, y) &= -(2x - 1)y(1 - y)\end{aligned}$$

Then

$$\begin{aligned}\dot{\varepsilon}_{xx}(x, y) &= (2y - 1)(1 - 2x) \\ \dot{\varepsilon}_{xy}(x, y) &= x(1 - x) - y(1 - y) \\ \dot{\varepsilon}_{yy}(x, y) &= -(2x - 1)(1 - 2y)\end{aligned}$$

We of course recover

$$\dot{\varepsilon}_{xx}(x, y) + \dot{\varepsilon}_{yy}(x, y) = 0$$

The Stokes equation is then:

$$\begin{aligned}-\frac{\partial p}{\partial x} + 2\eta(-2(2y - 1) - (1 - 2y)) &= \rho g_x \\ -\frac{\partial p}{\partial y} + 2\eta((1 - 2x) + 2(2x - 1)) &= \rho g_y\end{aligned}$$

or

$$\begin{aligned}-\frac{\partial p}{\partial x} + 2\eta(1 - 2y) &= \rho g_x \\ -\frac{\partial p}{\partial y} - 2\eta(1 - 2x) &= \rho g_y\end{aligned}$$

where we assume the viscosity  $\eta = 1$  to be constant in space. Assuming  $g_x = 0$ , the first equation is

$$\frac{\partial p}{\partial x} = 2(1 - 2y)$$

i.e.

$$p(x, y) = 2x(1 - 2y) + f(y)$$

Inserting this in the second equation:

$$4x - f'(y) - 2(1 - 2x) = \rho g_y$$

Assuming  $g_y = -1$ , then

$$8x - 2 - f'(y) = \rho$$

we set  $\rho(x, y) = 8x - 2$  and then  $f'(y) = 0$  so  $f(y) = C$  where  $C$  is a constant. Finally the pressure is given by:

$$p(x, y) = 2x(1 - 2y) + C$$

We add the following requirement:  $\int_{\Omega} p(x, y) d\Omega = 0$

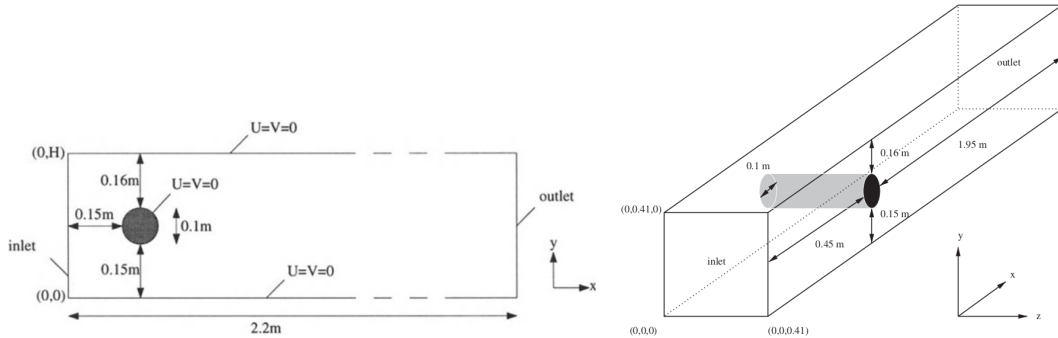
$$\int_0^1 \int_0^1 p(x, y) = 0 \Rightarrow \int_0^1 \int_0^1 [2x(1 - 2y) + C] = 0$$

so that  $C = 0$ .


We also find  $\mathbf{v}_{rms} = \frac{1}{\sqrt{45}} \simeq 0.1490711985$

### 12.2.10 Viscous flow around a cylinder in 2D and 3D

There are many variants of this problem: 2D in Turek [1292], 3D in John [651]. Many studies focus on Navier-Stokes flow since the cylinder generates vortices at high Reynolds numbers. Steady state solutions at low Re are shown here<sup>7</sup>. Note the interesting benchmark for 2D visco-elastic flow in Beuchert & Podlachikov [86].



Left: taken from Turek [1292]; Right: taken from John [651]

 Relevant Literature: Tachibana & Iemoto (1987) [1226] Schäfer & Turek (1996) [1115]

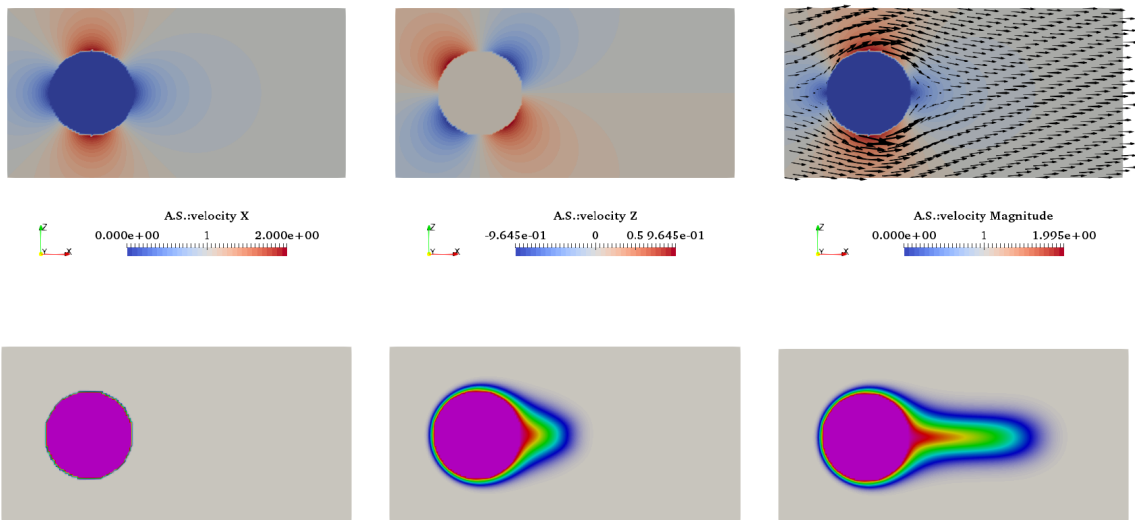
### 12.2.11 Heat flow around a cylinder

The domain is a 2D Cartesian box of size 8x4. The Stokes equations are not solved and the following velocity is prescribed:

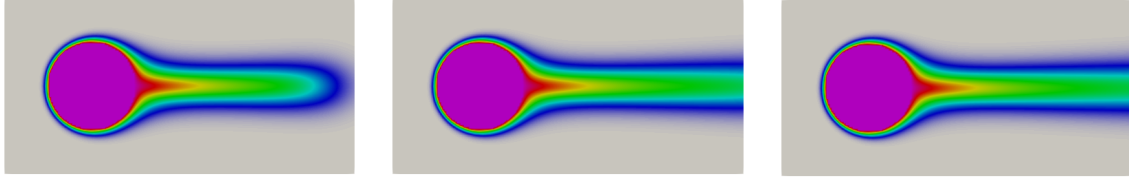
$$u(x, y) = U_{\infty} \left( 1 - \frac{x^2 - y^2}{(x^2 + y^2)^2} \right) \quad (12.193)$$

$$v(x, y) = -2U_{\infty} \frac{xy}{(x^2 + y^2)^2} \quad (12.194)$$

Boundary conditions are as follows:  $T = 0$  is imposed at the top and bottom of the domain.  $T = 1$  is imposed inside a disc centered at (2,2) with radius 1. Further:  $k = 0.01$ ,  $C_p = 1$ ,  $\rho = 1$ , CFL number is 0.1.



<sup>7</sup>upofthetestanddatameasurement



Time evolution of the temperature field. Results obtained with ELEFANT (unpublished)

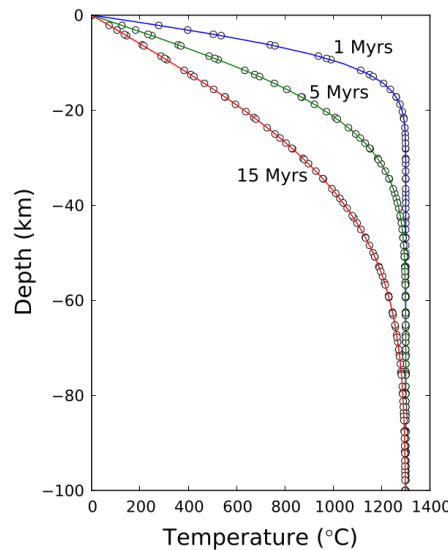
This is carried out in [STONE](#) 65.

### 12.2.12 Thermal diffusion of half-cooling space

This is a simple 1D experiment which solution is (for instance) available in Turcotte & Schubert [1288] and is also presented in Choi *et al.* [238].

The domain is 100km deep.  $T_0 = 0^\circ\text{C}$  is prescribed at the surface and  $T_m = 1300^\circ\text{C}$  is prescribed at the bottom. The initial temperature is  $T(y) = 1300^\circ\text{C}$ . The material is characterised by  $\rho = 1000\text{kg/m}^3$ ,  $C_p = 1000\text{J/kg/K}$ ,  $k = 1\text{J/m/K}$ . The time-dependent solution is given by:

$$T(y, t) = T_0 + (T_0 - T_m) \operatorname{erf} \left( \frac{y}{2\sqrt{kt/\rho C_p}} \right) \quad (12.195)$$

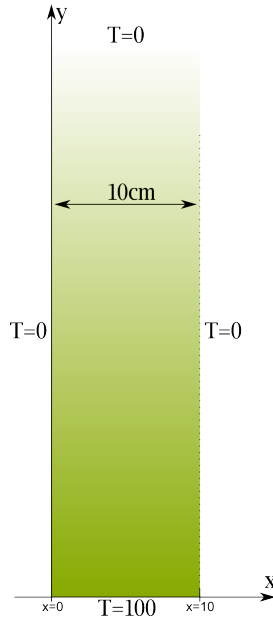


Thermal diffusion of half space cooling plate. The temperature profiles in the analytical solution at 1, 5, and 15 Myrs are plotted in solid lines. The results from DynEarthSol2D are plotted in circles. Taken from [238]

### 12.2.13 Laplace equation on a semi infinite plate

benchmark\_laplace\_plate.tex

This experiment is based on a 2nd year mathematics lecture I give at Utrecht University. One wishes to solve the Laplace equation for temperature on the following plate subject to the indicated boundary conditions:



The temperature satisfies the 2D Laplace equation inside the plate:

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0 \quad (12.196)$$

We could try to solve the equation by using a tentative solution of the form:

$$T(x, y) = \theta(x)\Phi(y) \quad (12.197)$$



We do not *know* the solution is of this form.

We substitute (2) into (1) and obtain:

$$\Phi \frac{\partial^2 \theta}{\partial x^2} + \theta \frac{\partial^2 \Phi}{\partial y^2} = 0$$

Dividing by  $\theta\Phi$  gives:

$$\frac{1}{\theta} \frac{\partial^2 \theta}{\partial x^2} + \frac{1}{\Phi} \frac{\partial^2 \Phi}{\partial y^2} = 0$$

Separation of variables: we say that each term is a constant because the first term is a function of  $x$  only and the second a function of  $y$  only. We then write

$$\frac{1}{\theta} \frac{\partial^2 \theta}{\partial x^2} = -\frac{1}{\Phi} \frac{\partial^2 \Phi}{\partial y^2} = -k^2$$

where  $k$  is called the separation constant. This leads to

$$\frac{\partial^2 \theta}{\partial x^2} + k^2 \theta = 0$$

$$\frac{\partial^2 \Phi}{\partial y^2} - k^2 \Phi = 0$$

- The solution to the first one is  $\theta(x) = \sin kx$  or  $\theta(x) = \cos kx$
- The solution to the second one is  $\Phi(y) = e^{ky}$  or  $\Phi(y) = e^{-ky}$



The general solution writes:

$$T(x, y) = \theta(x)\Phi(y) = \begin{Bmatrix} \sin kx \\ \cos kx \end{Bmatrix} \begin{Bmatrix} e^{ky} \\ e^{-ky} \end{Bmatrix}$$

We can now use the b.c. to find the solution to the Laplace equation.

- Since  $T \rightarrow 0$  when  $y \rightarrow \infty$  then  $e^{ky}$  unacceptable.
- Since  $T = 0$  when  $x = 0$  then  $\cos kx$  unacceptable.

so

$$T(x, y) = \sin(kx) e^{-ky}$$

We finally use  $T = 0$  at  $x = 10$  which leads to  $10k = n\pi$ , i.e.:

$$T(x, y) = \sin\left(\frac{n\pi x}{10}\right) e^{-n\pi y/10}$$

⚠ Problem: the solution does not satisfy  $T(x, 0) = 100$ . However, a linear combination of solutions is still a solution ! Let's find such a combination which satisfies the b.c. at  $y = 0$  :

$$T(x, y) = \sum_{n=1}^{\infty} b_n \sin\left(\frac{n\pi x}{10}\right) e^{-n\pi y/10}$$

We impose then  $T(x, 0) = 100$ :

$$100 = \sum_{n=1}^{\infty} b_n \sin\left(\frac{n\pi x}{10}\right)$$

This is the Fourier sine series of  $f(x) = 100$  with  $l = 10$  (chapter 7.9 of Boas).

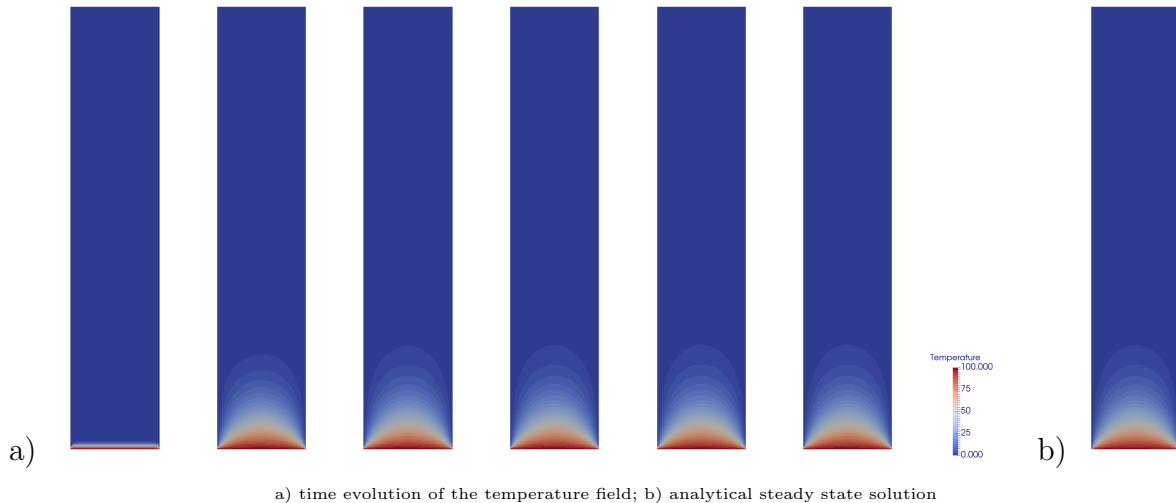
The coefficient  $b_n$  is then given by

$$b_n = \frac{2}{l} \int_0^l f(x) \sin \frac{n\pi x}{l} dx = \frac{2}{10} \int_0^{10} 100 \sin \frac{n\pi x}{10} dx = \begin{cases} 400/n\pi & \text{odd } n \\ 0 & \text{even } n \end{cases}$$

Finally (!):

$$T(x, y) = \frac{400}{\pi} \left( e^{-\pi y/10} \sin\left(\frac{\pi x}{10}\right) + \frac{1}{3} \sin\left(\frac{3\pi x}{10}\right) e^{-3\pi y/10} + \dots \right)$$

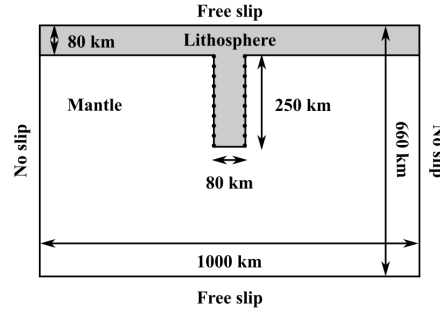
The simulation has been run with a 10x50 domain. All coefficients of the temperature equation are set to 1, and the Stokes equation is not solved. The timestep is fixed to  $dt = 0.1$ . Resolution is 32x160.





## 12.2.14 Slab detachment benchmark

 Relevant Literature: Schmeling (2011) [1116], ASPECT manual [44], Glerum *et al.* [467], [STONE](#) 26.



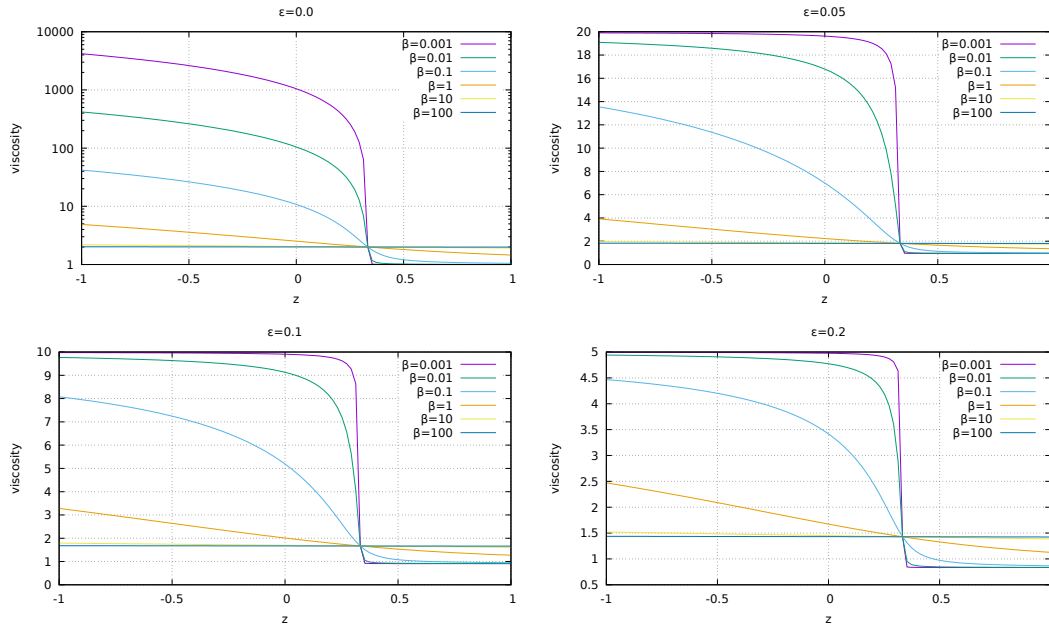
The detachment benchmark model setup of Schmalholz [1116]: a symmetric system of nonlinear viscous lithosphere with a vertical slab extending into a linear viscous mantle. The top and bottom boundaries are free slip, while the vertical boundaries are no slip. Taken from [467].

## 12.2.15 Layered flow with viscosity contrast

The idea behind this benchmark is to construct an analytical solution to the incompressible Stokes equation in the case where the viscosity field showcases a viscosity contrast at location  $y = y_0$  whose amplitude and width can be controlled. The viscosity is defined as

$$\eta(y) = \frac{1}{\frac{1}{\pi} \tan^{-1}\left(\frac{y-y_0}{\beta}\right) + 1/2 + \epsilon}$$

where  $\beta$  and  $\epsilon$  are parameters.



Viscosity profiles for different values of  $\beta$  and  $\epsilon$  for  $y_0 = 1/3$ . When  $\beta$  is very large, the viscosity essentially converges to  $\sim (1/2 + \epsilon)^{-1}$ .  $\beta$  controls the width of the transition while  $\epsilon$  controls the amplitude of the viscosity variation.

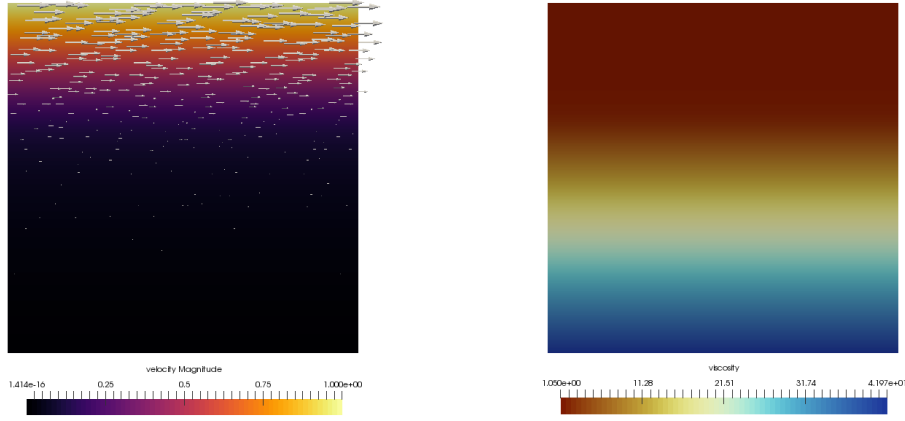
The flow is assumed to take place in an infinitely long pipe (in the horizontal direction) and bound by  $y = -1$  and  $y = 1$ . At the bottom we impose  $v_x(y = -1) = 0$  while we impose  $v_x(y = 1) = 1$  at the top. The density is set to 1 while the gravity is set to zero. Under these assumptions, the flow

velocity and pressure fields are given by:

$$\begin{aligned} v_x(x, y) &= \frac{1}{2\pi} \left( -\beta C_1 \log[\beta^2 + (z - y_0)^2] + 2(z - y_0)C_1 \tan^{-1} \frac{z - y_0}{\beta} + \pi(1 + 2\epsilon)zC_1 + C_2 \right) \\ v_y(x, y) &= 0 \\ p(x, y) &= 0 \end{aligned} \quad (12.198)$$

where  $C_1$  and  $C_2$  are integration constants:

$$\begin{aligned} C_1 &= 2\pi \left[ \beta \log[\beta^2 + (1 + y_0)^2] - 2(1 + y_0) \tan^{-1} \frac{1 + y_0}{\beta} - \beta \log[\beta^2 + (1 - y_0)^2] + 2(1 - y_0) \tan^{-1} \frac{1 - y_0}{\beta} + 2 \right] \\ C_2 &= \left[ \beta \log[\beta^2 + (1 + y_0)^2] - 2(1 + y_0) \tan^{-1} \frac{1 + y_0}{\beta} + \pi(1 + 2\epsilon) \right] C_1, \end{aligned}$$



Velocity and viscosity fields

**Analytical derivations** The flow takes place in the horizontal direction and is infinite in the this direction too so that:

$$\vec{v} = (u(y), 0)$$

The strain rate tensor is then given by:

$$\dot{\epsilon} = \frac{1}{2} \begin{pmatrix} 0 & du/dy \\ du/dy & 0 \end{pmatrix}$$

The momentum equation then becomes:

$$\vec{\nabla} \cdot (2\eta\dot{\epsilon}) - \vec{\nabla}p = \vec{\nabla} \cdot \left[ \eta(y) \begin{pmatrix} 0 & du/dy \\ du/dy & 0 \end{pmatrix} \right] - \vec{\nabla}p = \rho\vec{g}$$

On the vertical axis, when the gravity is zero, the equation is automatically verified when the pressure is zero. On the horizontal axis:

$$\frac{d}{dy} \left( \eta(y) \frac{du}{dy} \right) = 0$$

Then

$$\eta(y) \frac{du}{dy} = C_1$$

or,

$$\frac{du}{dy} = \frac{C_1}{\eta(y)} = C_1 \left( \frac{1}{\pi} \tan^{-1} \frac{y - y_0}{\beta} + 1/2 + \epsilon \right)$$

so that the velocity is given by:

$$u(y) = \frac{1}{\pi} (y \tan^{-1}((y - y_0)/\beta) - y_0 \tan^{-1}((y - y_0)/\beta) - 0.5 * \beta \log(\beta^2 + y^2 - 2yy_0 + y_0^2) + \pi y(\epsilon + 0.5))$$

$$u(z) = \frac{1}{2\pi} \left( -\beta C_1 \log[\beta^2 + (z - y_0)^2] + 2(z - y_0)C_1 \tan^{-1} \frac{z - y_0}{\beta} + \pi(1 + 2\epsilon)zC_1 + C_2 \right)$$

where  $C_1$  and  $C_2$  are integration constants. I wish to impose  $u(z = -1) = 0$  and  $u(z = +1) = 1$ :

$$\frac{1}{2\pi} \left( -\beta C_1 \log[\beta^2 + (-1 - y_0)^2] + 2(-1 - y_0)C_1 \tan^{-1} \frac{-1 - y_0}{\beta} - \pi(1 + 2\epsilon)C_1 + C_2 \right) = 0$$

$$\frac{1}{2\pi} \left( -\beta C_1 \log[\beta^2 + (1 - y_0)^2] + 2(1 - y_0)C_1 \tan^{-1} \frac{1 - y_0}{\beta} + \pi(1 + 2\epsilon)C_1 + C_2 \right) = 1$$

or,

$$-\beta C_1 \log[\beta^2 + (-1 - y_0)^2] + 2(-1 - y_0)C_1 \tan^{-1} \frac{-1 - y_0}{\beta} - \pi(1 + 2\epsilon)C_1 + C_2 = 0$$

$$-\beta C_1 \log[\beta^2 + (1 - y_0)^2] + 2(1 - y_0)C_1 \tan^{-1} \frac{1 - y_0}{\beta} + \pi(1 + 2\epsilon)C_1 + C_2 = 2\pi$$

or,

$$-\beta C_1 \log[\beta^2 + (-1 - y_0)^2] + 2(1 + y_0)C_1 \tan^{-1} \frac{1 + y_0}{\beta} - \pi(1 + 2\epsilon)C_1 + C_2 = 0$$

$$-\beta C_1 \log[\beta^2 + (1 - y_0)^2] + 2(1 - y_0)C_1 \tan^{-1} \frac{1 - y_0}{\beta} + \pi(1 + 2\epsilon)C_1 + C_2 = 2\pi$$

or,

$$-\beta C_1 \log(\beta^2 + (1 + y_0)^2) + 2(1 + y_0)C_1 \tan^{-1}((1 + y_0)/\beta) - \pi(1 + 2\epsilon)C_1 + C_2 = 0$$

$$-\beta C_1 \log(\beta^2 + (1 - y_0)^2) + 2(1 - y_0)C_1 \tan^{-1}((1 - y_0)/\beta) + \pi(1 + 2\epsilon)C_1 + C_2 = 2\pi$$

I can now substract the first line from the second line:

$$\beta C_1 \log(\beta^2 + (1 + y_0)^2) - 2(1 + y_0)C_1 \tan^{-1}((1 + y_0)/\beta) - \beta C_1 \log(\beta^2 + (1 - y_0)^2) + 2(1 - y_0)C_1 \tan^{-1}((1 - y_0)/\beta) + 2\pi(1 + 2\epsilon)C_1 = 2\pi$$

i.e.,

$$C_1 = 2\pi \left[ \beta \log[\beta^2 + (1 + y_0)^2] - 2(1 + y_0) \tan^{-1} \left[ \frac{1 + y_0}{\beta} \right] - \beta \log[\beta^2 + (1 - y_0)^2] + 2(1 - y_0) \tan^{-1} \left[ \frac{1 - y_0}{\beta} \right] + 2\pi(1 + 2\epsilon) \right]$$

and then

$$C_2 = \beta C_1 \log(\beta^2 + (1 + y_0)^2) - 2(1 + y_0)C_1 \tan^{-1}((1 + y_0)/\beta) + \pi(1 + 2\epsilon)C_1$$

## 12.2.16 The annulus convection benchmark # 1

benchmark\_annulus\_converction\_benchmark1.tex

We wish to solve the Stokes equation in an annulus of inner radius  $R_1$  and outer radius  $R_2$  with the following boundary conditions:

- Inner boundary:  $\mathbf{v}_r(R_1, \theta) = 0$
- Outer boundary:  $\mathbf{v}_r(R_2, \theta) = 0$

We then postulate

$$\mathbf{v}_\theta(r, \theta) = f(r) \cos(k\theta)$$

Note that in the case  $k = 0$ , we recover a constant velocity on the inner and outer boundaries.

The divergence of an incompressible vector field in polar coordinates is

$$\frac{1}{r} \frac{\partial(r\mathbf{v}_r)}{\partial r} + \frac{1}{r} \frac{\partial\mathbf{v}_\theta}{\partial \theta} = 0$$

or,

$$\frac{\partial(r\mathbf{v}_r)}{\partial r} + \frac{\partial\mathbf{v}_\theta}{\partial \theta} = 0$$

i.e.,

$$\frac{\partial(r\mathbf{v}_r)}{\partial r} = -\frac{\partial\mathbf{v}_\theta}{\partial \theta} = kf(r) \sin(k\theta)$$

so

$$r\mathbf{v}_r(r, \theta) = k \left[ \int f(r) dr \right] \sin(k\theta)$$

and finally

$$\mathbf{v}_r(r, \theta) = kg(r) \sin(k\theta)$$

with

$$g(r) = \frac{1}{r} \int f(r) dr \quad (12.200)$$

$$g'(r) = -\frac{1}{r^2} \int f(r) dr + \frac{1}{r} f = -\frac{1}{r} g + \frac{1}{r} f = \frac{1}{r} (f - g) \quad (12.201)$$

The boundary conditions lead to

$$\mathbf{v}_r(r = R_1, \theta) = k \left( \frac{A}{2} R_1 + \frac{B}{R_1} \ln R_1 + \frac{C}{R_1} \right) \sin(k\theta) = 0$$

$$\mathbf{v}_r(r = R_2, \theta) = k \left( \frac{A}{2} R_2 + \frac{B}{R_2} \ln R_2 + \frac{C}{R_2} \right) \sin(k\theta) = 0$$

This has to be valid  $\forall \theta$ , so

$$\frac{A}{2} R_1^2 + B \ln R_1 = -C \quad \text{and} \quad \frac{A}{2} R_2^2 + B \ln R_2 = -C$$

leading to

$$\frac{A}{2} + \frac{B}{R_1^2} \ln R_1 = -\frac{C}{R_1^2} \quad \text{and} \quad \frac{A}{2} + \frac{B}{R_2^2} \ln R_2 = -\frac{C}{R_2^2}$$

and finally

$$B = -C \frac{R_2^2 - R_1^2}{R_2^2 \ln R_1 - R_1^2 \ln R_2}$$

Likewise

$$\frac{A}{2}R_1^2 + B \ln R_1 = -C \quad \text{and} \quad \frac{A}{2}R_2^2 + B \ln R_2 = -C$$

yields

$$\frac{A}{2 \ln R_1} R_1^2 + B = -\frac{C}{\ln R_1} \quad \text{and} \quad \frac{A}{2 \ln R_2} R_2^2 + B = -\frac{C}{\ln R_2}$$

or,

$$\begin{aligned} \frac{A}{2 \ln R_1} R_1^2 - \frac{A}{2 \ln R_2} R_2^2 &= -C \left( \frac{1}{\ln R_1} - \frac{1}{\ln R_2} \right) \\ A \left( \frac{R_1^2}{2 \ln R_1} - \frac{R_2^2}{2 \ln R_2} \right) &= -C \left( \frac{1}{\ln R_1} - \frac{1}{\ln R_2} \right) \\ A \left( \frac{R_1^2 \ln R_2}{2} - \frac{R_2^2 \ln R_1}{2} \right) &= -C (\ln R_2 - \ln R_1) \end{aligned}$$

finally

$$A = -C \frac{2(\ln R_2 - \ln R_1)}{R_1^2 \ln R_2 - R_2^2 \ln R_1} = -C \frac{2(\ln R_1 - \ln R_2)}{R_2^2 \ln R_1 - R_1^2 \ln R_2}$$

We set  $\vec{g} = -g_r \vec{e}_r$ . Stokes equation in Polar coordinates (p284 of Schubert, Turcotte and Olson book):

- $r$ -component:

$$\eta \left[ \nabla^2 \mathbf{v}_r - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right] + \frac{\eta}{3} \frac{\partial}{\partial r} \left[ \frac{1}{r} \frac{\partial(r \mathbf{v}_r)}{\partial r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right] - \frac{\partial p}{\partial r} - \rho g_r = 0$$

The second term between brackets in the divergence of the velocity field so it is equal to zero in our case. We end up with

$$\eta \left[ \nabla^2 \mathbf{v}_r - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right] - \frac{\partial p}{\partial r} - \rho g_r = 0$$

- $\theta$ -component:

$$\eta \left[ \nabla^2 \mathbf{v}_\theta + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} \right] + \frac{\eta}{3} \frac{1}{r} \frac{\partial}{\partial \theta} \left[ \frac{1}{r} \frac{\partial(r \mathbf{v}_r)}{\partial r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right] - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0$$

The second term between brackets in the divergence of the velocity field so it is equal to zero in our case. We end up with

$$\eta \left[ \nabla^2 \mathbf{v}_\theta + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} \right] - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0$$

In both equations,  $\nabla^2$  represents the Laplacian of a scalar quantity:

$$\nabla^2 = \frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2}$$

We can then write the two momentum equations for an incompressible Stokes flow in polar coordinates:

$$\begin{aligned} \eta \left[ \frac{\partial^2 \mathbf{v}_r}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_r}{\partial \theta^2} - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} \right] - \frac{\partial p}{\partial r} - \rho g_r &= 0 \\ \eta \left[ \frac{\partial^2 \mathbf{v}_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} \right] - \frac{1}{r} \frac{\partial p}{\partial \theta} &= 0 \end{aligned}$$

We can further choose  $\eta = 1$ , so that

$$\frac{\partial^2 \mathbf{v}_r}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_r}{\partial \theta^2} - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} - \frac{\partial p}{\partial r} - \rho g_r = 0 \quad (12.202)$$

$$\frac{\partial^2 \mathbf{v}_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0 \quad (12.203)$$

Let us define  $f(r) = Ar + B/r$ . We then have

$$\frac{\partial^2 f}{\partial r^2} + \frac{1}{r} \frac{\partial f}{\partial r} - \frac{f}{r^2} = A \left( \frac{\partial^2 r}{\partial r^2} + \frac{1}{r} \frac{\partial r}{\partial r} - \frac{1}{r} \right) + B \left( \frac{\partial^2 r^{-1}}{\partial r^2} + \frac{1}{r} \frac{\partial r^{-1}}{\partial r} - \frac{1}{r^3} \right) = 0 \quad (12.204)$$

Eq. (12.203) simplifies to:

$$\frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0$$

We have

$$\frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} = -k^2 \frac{f(r)}{r^2} \cos(k\theta) \quad \text{and} \quad \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} = \frac{2k^2}{r^2} g(r) \cos(k\theta)$$

so

$$\frac{1}{r} \frac{\partial p}{\partial \theta} = -\frac{k^2}{r^2} f(r) \cos(k\theta) + \frac{2k^2}{r^2} g(r) \cos(k\theta) = k^2 \left( \frac{2g(r) - f(r)}{r^2} \right) \cos(k\theta)$$

and then

$$\frac{\partial p}{\partial \theta} = k^2 \left( \frac{2g(r) - f(r)}{r} \right) \cos(k\theta)$$

This can be integrated with respect to  $\theta$ :

$$p(r, \theta) = k \left( \frac{2g(r) - f(r)}{r} \right) \sin(k\theta) + l(r) = kh(r) \sin(k\theta) + l(r)$$

where  $h(r) = \frac{1}{r}(2g(r) - f(r))$ . We can turn to Eq.(12.202).

$$\begin{aligned} \rho g_r &= \frac{\partial^2 \mathbf{v}_r}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_r}{\partial \theta^2} - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} - \frac{\partial p}{\partial r} \\ &= +kg''(r) \sin(k\theta) + k \frac{g'(r)}{r} \sin(k\theta) - k^3 \frac{g(r)}{r^2} \sin(k\theta) - k \frac{g(r)}{r^2} \sin(k\theta) + k \frac{2f(r)}{r^2} \sin(k\theta) - kh'(r) \sin(k\theta) \\ &= k \sin(k\theta) \left[ g'' + \frac{g'}{r} - \frac{k^2 g}{r^2} - \frac{g}{r^2} + \frac{2f}{r^2} - h' \right] - l'(r) \\ &= k \sin(k\theta) \left[ g'' + \frac{g'}{r} - \frac{k^2 g}{r^2} - \frac{g}{r^2} + \frac{2f}{r^2} + \frac{1}{r^2} (2g - f) - \frac{1}{r} (2g' - f') \right] - l'(r) \\ &= k \sin(k\theta) \left[ g'' + \frac{g'}{r} (1 - 2) - \frac{g}{r^2} (k^2 + 1 - 2) + \frac{f}{r^2} (2 - 1) + \frac{f'}{r} \right] - l'(r) \\ &= k \sin(k\theta) \left[ g'' - \frac{g'}{r} - \frac{g}{r^2} (k^2 - 1) + \frac{f}{r^2} + \frac{f'}{r} \right] - l'(r) \end{aligned} \quad (1)$$

We can further choose  $g_r = 1$ . Note that when  $k = 0$ , we have  $\rho = -l'(r)$ . We then choose  $l'(r) = -\rho_0$  so that the  $k$ -dependent term can be seen as a density perturbation:

$$\rho = k \sin(k\theta) \aleph(r) + \rho_0$$

with

$$\aleph(r) = g'' + \frac{g'}{r} \left( 1 - \frac{2}{r} \right) - \frac{g}{r^2} \left( k^2 + 1 - \frac{4}{r} \right) + \frac{2f}{r^2} \left( 1 - \frac{1}{r} \right) + \frac{f'}{r^2}$$

and

$$f(r) = Ar + \frac{B}{r} \quad (12.206)$$

$$f'(r) = A - \frac{B}{r^2} \quad (12.207)$$

$$g(r) = \frac{A}{2}r + \frac{B}{r} \ln r - \frac{1}{r} \quad (12.208)$$

$$g'(r) = \frac{A}{2} + \frac{B}{r^2}(1 - \ln r) + \frac{1}{r^2} \quad (12.209)$$

$$g''(r) = -\frac{2B}{r^3}(1 - \ln r) - B\frac{1}{r^3} - \frac{2}{r^3} = -\frac{B}{r^3}(3 - 2 \ln r) \quad (12.210)$$

Finally, the pressure is then given by

$$p(r, \theta) = k \left( \frac{2g - f}{r^2} \right) \sin(k\theta) + l(r) = kh(r) \sin(k\theta) + \rho_0 g_r r + Constant$$

We enforce  $p(r = R_2, \theta) = 0$  so that

$$p(r, \theta) = k \left( \frac{2g - f}{r^2} \right) \sin(k\theta) + l(r) = kh(r) \sin(k\theta) + \rho_0 g_r (r - R_2)$$

**Summary of the previous pages :**

$$\mathbf{v}_\theta(r, \theta) = f(r) \cos(k\theta) \quad (12.211)$$

$$\mathbf{v}_r(r, \theta) = g(r)k \sin(k\theta) \quad (12.212)$$

$$p(r, \theta) = kh(r) \sin(k\theta) + \rho_0 g_r (r - R_2) \quad (12.213)$$

$$\rho(r, \theta) = k \sin(k\theta) \aleph(r) + \rho_0 \quad (12.214)$$

$$A = \frac{2(\ln R_1 - \ln R_2)}{R_2^2 \ln R_1 - R_1^2 \ln R_2} \quad (12.215)$$

$$B = \frac{R_2^2 - R_1^2}{R_2^2 \ln R_1 - R_1^2 \ln R_2} \quad (12.216)$$

$$f(r) = Ar + \frac{B}{r} \quad (12.217)$$

$$f'(r) = A - \frac{B}{r^2} \quad (12.218)$$

$$g(r) = \frac{A}{2}r + \frac{B}{r} \ln r - \frac{1}{r} \quad (12.219)$$

$$g'(r) = \frac{A}{2} + \frac{B}{r^2}(1 - \ln r) + \frac{1}{r^2} \quad (12.220)$$

$$g''(r) = -\frac{B}{r^3}(3 - 2 \ln r) \quad (12.221)$$

$$h(r) = \frac{1}{r^2}(2g - f) \quad (12.222)$$

$$\aleph(r) = -g'' - \frac{g'}{r}(1 - \frac{2}{r}) + \frac{g}{r^2}(k^2 + 1 - \frac{4}{r}) - \frac{2f}{r^2}(1 - \frac{1}{r}) - \frac{f'}{r^2} \quad (12.223)$$

**Averagings of fields :**

- Average  $\mathbf{v}_r$  velocity

$$\langle \mathbf{v}_r(r) \rangle = \frac{1}{2\pi} \int_0^{2\pi} \mathbf{v}_r(r, \theta) d\theta = \frac{1}{2\pi} \int_0^{2\pi} g(r)k \sin(k\theta) d\theta = \frac{1}{2\pi} g(r)k \int_0^{2\pi} \sin(k\theta) d\theta = 0$$

since  $k = 0, 2, 4, \dots$

- Average  $\mathbf{v}_\theta$  velocity

$$\langle \mathbf{v}_\theta(r) \rangle = \frac{1}{2\pi} \int_0^{2\pi} \mathbf{v}_\theta(r, \theta) d\theta = \frac{1}{2\pi} \int_0^{2\pi} f(r) \cos(k\theta) d\theta = \frac{1}{2\pi} f(r) \int_0^{2\pi} \cos(k\theta) d\theta = 0$$

since  $k = 0, 2, 4, \dots$

- Root mean square verage  $\mathbf{v}_r$  velocity

$$\begin{aligned} \langle \mathbf{v}_r \rangle_{rms}(r) &= \sqrt{\frac{1}{2\pi} \int_0^{2\pi} \mathbf{v}_r^2 d\theta} \\ &= \sqrt{\frac{1}{2\pi} g(r)^2 k^2 \int_0^{2\pi} \sin^2(k\theta) d\theta} \\ &= \sqrt{\frac{1}{2\pi} g(r)^2 k^2 \int_0^{2\pi} \frac{1}{2} (1 - \cos(2k\theta)) d\theta} \\ &= \sqrt{\frac{1}{2\pi} g(r)^2 k^2 \left( \pi - \frac{1}{2} \int_0^{2\pi} \cos(2k\theta) d\theta \right)} \\ &= \sqrt{\frac{1}{2\pi} g(r)^2 k^2 \left( \pi - \frac{1}{4k} \underbrace{\int_0^{2k\pi} \cos \alpha d\alpha}_{=0} \right)} \\ &= \frac{k|g(r)|}{\sqrt{2}} \end{aligned} \tag{12.224}$$

- Root mean square verage  $\mathbf{v}_\theta$  velocity

$$\begin{aligned} \langle \mathbf{v}_\theta \rangle_{rms}(r) &= \sqrt{\frac{1}{2\pi} \int_0^{2\pi} \mathbf{v}_\theta^2 d\theta} \\ &= \sqrt{\frac{1}{2\pi} f(r)^2 \int_0^{2\pi} \cos^2(k\theta) d\theta} \\ &= \sqrt{\frac{1}{2\pi} f(r)^2 \int_0^{2\pi} \frac{1}{2} (1 + \cos(2k\theta)) d\theta} \\ &= \sqrt{\frac{1}{2\pi} f(r)^2 \left( \pi + \frac{1}{2} \int_0^{2\pi} \cos(2k\theta) d\theta \right)} \\ &= \sqrt{\frac{1}{2\pi} f(r)^2 \left( \pi + \frac{1}{4k} \underbrace{\int_0^{4k\pi} \cos \alpha d\alpha}_0 \right)} \\ &= \frac{|f(r)|}{\sqrt{2}} \end{aligned} \tag{12.225}$$



- Root mean square velocity  $\mathbf{v}_{rms}$

$$\mathbf{v}_{rms} = \sqrt{\frac{1}{V} \int_V (\mathbf{v}_r^2 + \mathbf{v}_\theta^2) dV} \quad (12.226)$$

The volume of the domain is given by

$$V = \pi(R_2^2 - R_1^2)$$

and the sum

$$\begin{aligned} (\mathbf{v}_r^2 + \mathbf{v}_\theta^2) dV &= [(g(r)k \sin(k\theta))^2 + (f(r) \cos(k\theta))^2] r dr d\theta \\ &= [g(r)^2 r dr] [k^2 \sin^2(k\theta) d\theta] + [f(r)^2 r dr] [\cos^2(k\theta) d\theta] \end{aligned}$$

If  $k = 0$ , we have  $\mathbf{v}_\theta = f(r)$  and  $\mathbf{v}_r = 0$  so that

$$\begin{aligned} \mathbf{v}_{rms} &= \sqrt{\frac{1}{V} \int_0^{2\pi} d\theta \int_{R_1}^{R_2} f(r)^2 r dr} \\ &= \sqrt{\frac{2\pi}{\pi(R_2^2 - R_1^2)} \int_{R_1}^{R_2} f(r)^2 r dr} \\ &= \sqrt{\frac{2}{(R_2^2 - R_1^2)} \int_{R_1}^{R_2} \left( Ar + \frac{B}{r} \right)^2 r dr} \\ &= \sqrt{\frac{2}{(R_2^2 - R_1^2)} \int_{R_1}^{R_2} \left( A^2 r^3 + 2ABr + \frac{B^2}{r} \right) dr} \\ &= \sqrt{\frac{2}{(R_2^2 - R_1^2)} \left[ A^2 \frac{r^4}{4} + AB r^2 + B^2 \ln(r) \right]_{R_1}^{R_2}} \\ &= \sqrt{\frac{2}{(R_2^2 - R_1^2)} \left[ \frac{A^2}{4} (R_2^4 - R_1^4) + AB(R_2^2 - R_1^2) + B^2(\ln R_2 - \ln R_1) \right]} \quad (12.227) \end{aligned}$$

If  $k \neq 0$  we have

$$\begin{aligned} \int_0^{2\pi} k^2 \sin^2(k\theta) d\theta &= \frac{1}{2} k^2 \int_0^{2\pi} (1 - \cos(k\theta)) d\theta = \pi k^2 \\ \int_0^{2\pi} \cos^2(k\theta) d\theta &= \frac{1}{2} \int_0^{2\pi} (1 + \cos(k\theta)) d\theta = \pi \end{aligned}$$

so that

$$\mathcal{V} = \int_0^{2\pi} \int_{R_1}^{R_2} (\mathbf{v}_r^2 + \mathbf{v}_\theta^2) r dr d\theta = \pi k^2 \int_{R_1}^{R_2} g(r)^2 r dr + \pi \int_{R_1}^{R_2} f(r)^2 r dr$$

$$\begin{aligned}
\int_{R_1}^{R_2} f(r)^2 r dr &= \int_{R_1}^{R_2} \left( Ar + \frac{B}{r} \right)^2 r dr \\
&= \int_{R_1}^{R_2} \left( A^2 r^3 + 2ABr + \frac{B^2}{r} \right) dr \\
&= \left[ A^2 \frac{r^4}{4} + AB r^2 + B^2 \ln(r) \right]_{R_1}^{R_2} \\
&= \frac{A^2}{4} (R_2^4 - R_1^4) + AB (R_2^2 - R_1^2) + B^2 (\ln R_2 - \ln R_1) \\
\\
\int_{R_1}^{R_2} g(r)^2 r dr &= \int_{R_1}^{R_2} \left( \frac{A}{2} r + \frac{B}{r} \ln r + \frac{C}{r} \right)^2 r dr \\
&= \int_{R_1}^{R_2} \left( \frac{A^2}{4} r^2 + AB \ln r + AC + \frac{2BC}{r^2} \ln r + \frac{B^2}{r^2} (\ln r)^2 + \frac{C^2}{r^2} \right) r dr \\
&= \int_{R_1}^{R_2} \left( \frac{A^2}{4} r^3 + AB r \ln r + AC r + \frac{2BC}{r} \ln r + \frac{B^2}{r} (\ln r)^2 + \frac{C^2}{r} \right) dr \\
&= \int_{R_1}^{R_2} \left( \frac{A^2}{4} r^3 + AC r + \frac{C^2}{r} \right) dr + E + F + G \\
&= \left[ \frac{A^2}{16} r^4 + \frac{1}{2} AC r^2 + C^2 \ln r \right]_{R_1}^{R_2} + E + F + G \\
&= \frac{A^2}{16} (R_2^4 - R_1^4) + \frac{AC}{2} (R_2^2 - R_1^2) + C^2 (\ln R_2 - \ln R_1) + E + F + G
\end{aligned}$$

(12.228)

$$\begin{aligned}
E &= 2BC \int_{R_1}^{R_2} \frac{1}{r} \ln r \, dr \\
&= BC \int_{\ln R_1}^{\ln R_2} 2X dX \quad X = \ln r, \, dX = dr/r \\
&= BC[X^2]_{\ln R_1}^{\ln R_2} \\
&= BC[(\ln R_2)^2 - (\ln R_1)^2]
\end{aligned}$$

$$\begin{aligned}
F &= B^2 \int_{R_1}^{R_2} \frac{1}{r} (\ln r)^2 dr \\
&= B^2 \int_{\ln R_1}^{\ln R_2} X^2 dX \quad X = \ln r, \, dX = dr/r \\
&= \frac{B^2}{3} [X^3]_{\ln R_1}^{\ln R_2} \\
&= \frac{B^2}{3} [(\ln R_2)^3 - (\ln R_1)^3]
\end{aligned}$$

$$\begin{aligned}
G &= AB \int_{R_1}^{R_2} r \ln r \, dr \\
&= AB \left[ \frac{1}{2} r^2 \ln r \right]_{R_1}^{R_2} - AB \int_{R_1}^{R_2} \frac{1}{2} r^2 \frac{1}{r} dr \\
&= \frac{AB}{2} [R_2^2 \ln R_2 - R_1^2 \ln R_1] - \frac{AB}{2} \int_{R_1}^{R_2} r \, dr \\
&= \frac{AB}{2} [R_2^2 \ln R_2 - R_1^2 \ln R_1] - \frac{AB}{4} [r^2]_{R_1}^{R_2} \\
&= \frac{AB}{2} [R_2^2 \ln R_2 - R_1^2 \ln R_1] - \frac{AB}{4} (R_2^2 - R_1^2)
\end{aligned}$$

$$\begin{aligned}
\mathcal{V}/\pi &= \frac{A^2}{4} (R_2^4 - R_1^4) + AB(R_2^2 - R_1^2) + B^2(\ln R_2 - \ln R_1) \\
&+ \frac{A^2 k^2}{16} (R_2^4 - R_1^4) + \frac{ACk^2}{2} (R_2^2 - R_1^2) + C^2 k^2 (\ln R_2 - \ln R_1) \\
&+ BCk^2 [(\ln R_2)^2 - (\ln R_1)^2] \\
&+ \frac{B^2 k^2}{3} [(\ln R_2)^3 - (\ln R_1)^3] \\
&+ \frac{ABk^2}{2} [R_2^2 \ln R_2 - R_1^2 \ln R_1] - \frac{ABk^2}{4} (R_2^2 - R_1^2) \\
&= \frac{A^2}{16} (4 + k^2) (R_2^4 - R_1^4) + \left( \frac{AB}{4} (4 - k^2) + \frac{ACk^2}{2} \right) (R_2^2 - R_1^2) \\
&+ (B^2 + C^2 k^2) (\ln R_2 - \ln R_1) \\
&+ BCk^2 [(\ln R_2)^2 - (\ln R_1)^2] \\
&+ \frac{B^2 k^2}{3} [(\ln R_2)^3 - (\ln R_1)^3] \\
&+ \frac{ABk^2}{2} [R_2^2 \ln R_2 - R_1^2 \ln R_1]
\end{aligned} \tag{12.229}$$

$$v_{rms} = \sqrt{\frac{1}{V} \mathcal{V}} = \sqrt{\frac{1}{2(R_2^2 - R_1^2) \pi} \mathcal{V}}$$

### Computing the strain rate and stress tensors :

Since we know the viscosity, the velocity field and the pressure field, we can also compute the full stress tensor  $\boldsymbol{\sigma} = -p\mathbf{1} + 2\eta\dot{\boldsymbol{\epsilon}}$ . In this benchmark we have set  $\eta = 1$  so:  $\boldsymbol{\sigma} = -p\mathbf{1} + 2\dot{\boldsymbol{\epsilon}}$ . We start with the velocity gradient:

$$\begin{aligned} \vec{\nabla} \vec{v} &= \begin{pmatrix} \frac{\partial v_r}{\partial r} & \frac{1}{r} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r} \\ \frac{\partial v_\theta}{\partial r} & \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} + \frac{v_r}{r} \end{pmatrix} \\ &= \begin{pmatrix} g'k \sin(k\theta) & \frac{1}{r} g k^2 \cos(k\theta) - \frac{1}{r} f \cos(k\theta) \\ f' \cos(k\theta) & -\frac{1}{r} f k \sin(k\theta) + \frac{1}{r} g k \sin(k\theta) \end{pmatrix} \\ &= \begin{pmatrix} g'k \sin(k\theta) & \frac{1}{r} (g k^2 - f) \cos(k\theta) \\ f' \cos(k\theta) & \frac{1}{r} (g - f) k \sin(k\theta) \end{pmatrix} \end{aligned} \quad (12.230)$$

The strain rate is then given by:

$$\dot{\boldsymbol{\epsilon}} = \frac{1}{2}(\nabla \mathbf{v} + \nabla \mathbf{v}^T) = \begin{pmatrix} g'k \sin(k\theta) & \frac{1}{2r} (r f' + g k^2 - f) \cos(k\theta) \\ \frac{1}{2r} (r f' + g k^2 - f) \cos(k\theta) & \frac{1}{r} (g - f) k \sin(k\theta) \end{pmatrix} \quad (12.231)$$

Let us verify once again that the flow is incompressible:

$$\vec{\nabla} \cdot \vec{v} = g'(r)k \sin(k\theta) + \frac{1}{r}(g(r) - f(r))k \sin(k\theta) = \frac{1}{r}(r g'(r) + g(r) - f(r))k \sin(k\theta) = 0$$

since  $g'(r) = \frac{1}{r}(f(r) - g(r))$ .

I can now write the full stress tensor:

$$\boldsymbol{\sigma} = -p\mathbf{1} + 2\dot{\boldsymbol{\epsilon}} = \begin{pmatrix} -p + 2g'k \sin(k\theta) & \frac{1}{r}(r f' + g k^2 - f) \cos(k\theta) \\ \frac{1}{r}(r f' + g k^2 - f) \cos(k\theta) & -p + \frac{2}{r}(g - f)k \sin(k\theta) \end{pmatrix} \quad (12.232)$$

On the boundaries, i.e.  $r = R_1$  or  $r = R_2$ , the function  $g$  is exactly zero, so that the stress  $\boldsymbol{\sigma}_b$  on the boundaries is given by

$$\boldsymbol{\sigma}_b = \begin{pmatrix} -p + 2g'k \sin(k\theta) & \frac{1}{r}(r f' - f) \cos(k\theta) \\ \frac{1}{r}(r f' - f) \cos(k\theta) & -p - \frac{2}{r} f k \sin(k\theta) \end{pmatrix} \quad (12.233)$$

Furthermore I can use the identity  $g' = (f - g)/r$  which simplifies to  $g' = f/r$  on the boundaries:

$$\boldsymbol{\sigma}_b = \begin{pmatrix} -p + 2\frac{f}{r}k \sin(k\theta) & \frac{1}{r}(r f' - f) \cos(k\theta) \\ \frac{1}{r}(r f' - f) \cos(k\theta) & -p - \frac{2}{r} f k \sin(k\theta) \end{pmatrix} \quad (12.234)$$

Also,  $h(r) = \frac{1}{r}(2g - f)$  simplifies to  $h(r) = -f/r$  so

$$p = k h(r) \sin(k\theta) + \rho_0 g_r (r - R_2) = -k \frac{f}{r} \sin(k\theta) + \rho_0 g_r (r - R_2)$$

Finally

$$\begin{aligned}\boldsymbol{\sigma}_b &= \begin{pmatrix} k\frac{f}{r}\sin(k\theta) - \rho_0 g_r(r - R_2) + 2\frac{f}{r}k\sin(k\theta) & \frac{1}{r}(rf' - f)\cos(k\theta) \\ \frac{1}{r}(rf' - f)\cos(k\theta) & k\frac{f}{r}\sin(k\theta) - \rho_0 g_r(r - R_2) - \frac{2}{r}fk\sin(k\theta) \end{pmatrix} \\ &= \begin{pmatrix} k\frac{3f}{r}\sin(k\theta) - \rho_0 g_r(r - R_2) & \frac{1}{r}(rf' - f)\cos(k\theta) \\ \frac{1}{r}(rf' - f)\cos(k\theta) & -k\frac{f}{r}\sin(k\theta) - \rho_0 g_r(r - R_2) \end{pmatrix}\end{aligned}$$

The traction along the normal is given by

$$\boldsymbol{\sigma}_n = \boldsymbol{\sigma}_b \cdot \vec{n} = -\boldsymbol{\sigma}_b \cdot \vec{e}_r = - \begin{pmatrix} k\frac{3f}{r}\sin(k\theta) - \rho_0 g_r(r - R_2) \\ \frac{1}{r}(rf' - f)\cos(k\theta) \end{pmatrix}$$

**Strain rate tensor in Cartesian coordinates** The analytical expressions for the strain rate components in polar coordinates are:

$$\dot{\epsilon}_{rr} = g'(r)k\sin k\theta \quad (12.235)$$

$$\dot{\epsilon}_{r\theta} = \dot{\epsilon}_{\theta r} = \frac{1}{2} \left( \frac{1}{r}g(r)k^2\cos k\theta + f'(r)\cos k\theta - \frac{f(r)}{r}\cos k\theta \right) \quad (12.236)$$

$$= \frac{1}{2} \left( \frac{1}{r}g(r)k^2 + f'(r) - \frac{f(r)}{r} \right) \cos k\theta \quad (12.237)$$

$$\dot{\epsilon}_{\theta\theta} = -\frac{1}{r}kf(r)\sin k\theta + \frac{g(r)}{r}k\sin k\theta \quad (12.238)$$

$$= \frac{g(r) - f(r)}{r}k\sin k\theta \quad (12.239)$$

Their counterparts in Cartesian coordinates are obtained with Eqs. (2.154).

**Could we find a steady state temperature field that goes along?** We can start from the pure advection equation:

$$\rho C_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = 0$$

At steady state we are left with

$$\vec{v} \cdot \vec{\nabla} T = 0$$

I postulate then

$$T(r, \theta) = l(r)(\alpha \cos k\theta + \beta \sin k\theta)$$

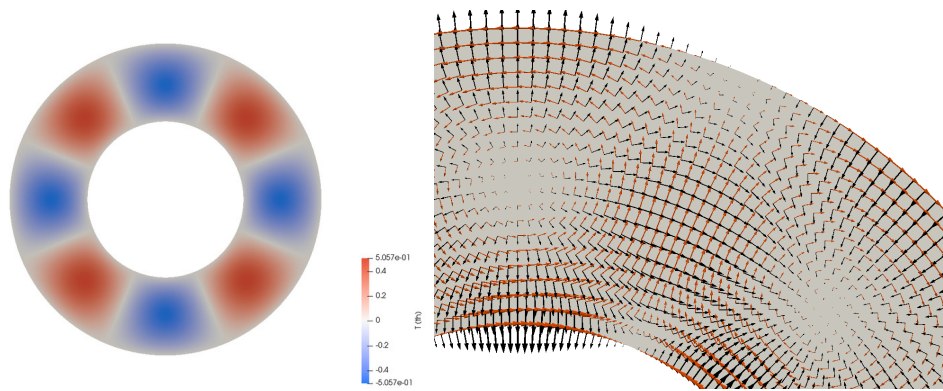
Inserting this into the equation above I arrive at the conclusion that  $\beta \neq 0$  is a dead end. It then simply follows that

$$T(r, \theta) = l(r) \cos k\theta$$

which yields  $T(r, \theta) = rg(r) \cos(k\theta)$ . and

$$\vec{\nabla} T = (f(r) \cos k\theta; -g(r)k \sin k\theta)$$

I can then plot the temperature gradient (in black) next to the velocity field in red (right figure) and show that these are indeed always perpendicular:



However, what is deeply puzzling is the fact that the temperature is exactly zero on both boundaries (because  $g(r)$  is by construction) and yet the temperature field is not zero in the middle, in the absence of heat source ... ?

### 12.2.17 The annulus convection benchmark #2 - no slip bc

We seek an exact solution to the incompressible Stokes equations for an isoviscous, isothermal fluid in an annulus. Given the geometry of the problem, we work in polar coordinates. We denote the orthonormal basis vectors by  $\vec{e}_r$  and  $\vec{e}_\theta$ , the inner radius of the annulus by  $R_1$  and the outer radius by  $R_2$ . Further, we assume that the viscosity  $\eta$  is constant, which we set to  $\eta = 1$ , we set the gravity vector to  $\vec{g} = -g_r \vec{e}_r$  with  $g_r = 1$ .

Given these assumptions, the incompressible Stokes equations in the annulus are (see Schubert, Turcotte & Olson [1140]):

$$A_r = \frac{\partial^2 \mathbf{v}_r}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_r}{\partial \theta^2} - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} - \frac{\partial p}{\partial r} = \rho g_r \quad (12.240)$$

$$A_\theta = \frac{\partial^2 \mathbf{v}_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0 \quad (12.241)$$

$$\frac{1}{r} \frac{\partial(r \mathbf{v}_r)}{\partial r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} = 0 \quad (12.242)$$

Equations (12.240) and (12.241) are the momentum equations in polar coordinates while Equation (12.242) is the mass conservation equation (also called continuity equation). The components of the velocity are obtained from the stream function  $\Psi$  as follows:

$$\mathbf{v}_r = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \quad \mathbf{v}_\theta = -\frac{\partial \Psi}{\partial r}$$

where  $\mathbf{v}_r$  is the radial component and  $\mathbf{v}_\theta$  is the tangential component of the velocity vector.

The stream function is defined for incompressible (divergence-free) flows in 2D (as well as in 3D with axisymmetry). The stream function can be used to plot streamlines, which represent the trajectories of particles in a steady flow. From calculus it is known that the gradient vector  $\nabla \Psi$  is normal to the curve  $\Psi = C$ . It can be shown that everywhere  $\vec{\mathbf{v}} \cdot \nabla \Psi = 0$  using the formula for  $\vec{u}$  in terms of  $\Psi$  which proves that level curves of  $\Psi$  are streamlines:

$$\vec{\mathbf{v}} \cdot \nabla \Psi = \mathbf{v}_r \frac{\partial \Psi}{\partial r} + \mathbf{v}_\theta \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \frac{\partial \Psi}{\partial r} - \frac{\partial \Psi}{\partial r} \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = 0$$

In polar coordinates the curl of a vector  $\vec{A}$  is <sup>8</sup>:

$$\vec{\nabla} \times \vec{A} = \frac{1}{r} \left( \frac{\partial(r A_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} \right) \vec{e}_z$$

Taking the curl of vector  $\vec{A}$  (see Eqs. (12.240) and (12.241)) yields:

$$\frac{1}{r} \left( \frac{\partial(r A_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} \right) = \frac{1}{r} \left( -\frac{\partial(\rho g_r)}{\partial \theta} \right)$$

Multiplying on each side by  $r$

$$\frac{\partial(r A_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} = -\frac{\partial(\rho g_r)}{\partial \theta}$$

If we now replace  $A_r$  and  $A_\theta$  by their expressions as a function of velocity and pressure, we will see that the pressure terms cancel out (which is one of the advantages of working with stream line formulations).

Let us assume the following separation of variables  $\boxed{\Psi(r, \theta) = \phi(r)\xi(\theta)}$ . Then

$$\mathbf{v}_r = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{\phi \xi'}{r} \quad \mathbf{v}_\theta = -\frac{\partial \Psi}{\partial r} = -\phi' \xi$$

<sup>8</sup>[https://en.wikipedia.org/wiki/Del\\_in\\_cylindrical\\_and\\_spherical\\_coordinates](https://en.wikipedia.org/wiki/Del_in_cylindrical_and_spherical_coordinates)

Let us first express  $A_r$  and  $A_\theta$  as functions of  $\phi$  and  $\xi$ :

$$\begin{aligned}
A_r &= \frac{\partial^2 \mathbf{v}_r}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_r}{\partial \theta^2} - \frac{\mathbf{v}_r}{r^2} - \frac{2}{r^2} \frac{\partial \mathbf{v}_\theta}{\partial \theta} - \frac{\partial p}{\partial r} \\
&= \frac{\partial^2}{\partial r^2} \left( \frac{\phi \xi'}{r} \right) + \frac{1}{r} \frac{\partial}{\partial r} \left( \frac{\phi \xi'}{r} \right) + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} \left( \frac{\phi \xi'}{r} \right) - \frac{1}{r^2} \left( \frac{\phi \xi'}{r} \right) - \frac{2}{r^2} \frac{\partial}{\partial \theta} (-\phi' \xi) - \frac{\partial p}{\partial r} \\
&= \left( \frac{\phi''}{r} - 2 \frac{\phi'}{r^2} + 2 \frac{\phi}{r^3} \right) \xi' + \left( \frac{\phi'}{r^2} - \frac{\phi}{r^3} \right) \xi + \frac{\phi}{r^3} \xi''' - \frac{\phi \xi'}{r^3} + \frac{2}{r^2} \phi' \xi' - \frac{\partial p}{\partial r} \\
&= \frac{\phi'' \xi'}{r} + \frac{\phi' \xi'}{r^2} + \frac{\phi \xi'''}{r^3} - \frac{\partial p}{\partial r} \\
\frac{\partial A_r}{\partial \theta} &= \frac{\phi'' \xi''}{r} + \frac{\phi' \xi''}{r^2} + \frac{\phi \xi''''}{r^3} - \frac{\partial^2 p}{\partial r \partial \theta} \tag{12.243}
\end{aligned}$$

$$\begin{aligned}
A_\theta &= \frac{\partial^2 \mathbf{v}_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} - \frac{1}{r} \frac{\partial p}{\partial \theta} \\
&= \frac{\partial^2}{\partial r^2} (-\phi' \xi) + \frac{1}{r} \frac{\partial}{\partial r} (-\phi' \xi) + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} (-\phi' \xi) + \frac{2}{r^2} \frac{\partial}{\partial \theta} \left( \frac{\phi \xi'}{r} \right) - \frac{1}{r^2} (-\phi' \xi) - \frac{1}{r} \frac{\partial p}{\partial \theta} \\
&= -\phi''' \xi - \frac{\phi'' \xi}{r} + \frac{\phi'(\xi - \xi'')}{r^2} + \frac{2\phi \xi''}{r^3} - \frac{1}{r} \frac{\partial p}{\partial \theta} \\
r A_\theta &= -\phi''' \xi r - \phi'' \xi + \frac{\phi'(\xi - \xi'')}{r} + \frac{2\phi \xi''}{r^2} - \frac{\partial p}{\partial \theta} \\
\frac{\partial(r A_\theta)}{\partial r} &= -(\phi''' \xi + \phi'''' \xi r) - \phi''' \xi + (\xi - \xi'') \left( \frac{\phi''}{r} - \frac{\phi'}{r^2} \right) + \left( \frac{2\phi' \xi''}{r^2} - 2 \frac{2\phi \xi'''}{r^3} \right) - \frac{\partial^2 p}{\partial \theta \partial r} \\
&= -2\phi''' \xi - \phi'''' \xi r + \frac{1}{r} \phi''(\xi - \xi'') + \frac{1}{r^2} (-\phi'(\xi - \xi'') + 2\phi' \xi'') + \frac{1}{r^3} (-4\phi \xi''') - \frac{\partial^2 p}{\partial \theta \partial r} \\
&= -2\phi''' \xi - \phi'''' \xi r + \frac{\phi''}{r}(\xi - \xi'') + \frac{\phi'}{r^2}(-\xi + 3\xi'') + \frac{\phi}{r^3}(-4\xi''') - \frac{\partial^2 p}{\partial \theta \partial r} \tag{12.244}
\end{aligned}$$

### No slip boundary conditions

No-slip boundary conditions inside and outside impose that all components of the velocity must be zero on both boundaries, i.e.

$$\vec{\mathbf{v}}(r = R_1, \theta) = \vec{\mathbf{v}}(r = R_2, \theta) = \vec{0}$$

Due to the separation of variables, and choosing  $\xi(\theta) = \cos(k\theta)$  we have

$$\mathbf{v}_r(r, \theta) = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{\phi \xi'}{r} = -\frac{1}{r} \phi(r) k \sin(k\theta) \tag{12.245}$$

$$\mathbf{v}_\theta(r, \theta) = -\frac{\partial \Psi}{\partial r} = -\phi'(r) \xi(\theta) = -\phi'(r) \cos(k\theta) \tag{12.246}$$

The velocity divergence is given by

$$\vec{\nabla} \cdot \vec{\mathbf{v}} = \frac{1}{r} \frac{\partial(r \mathbf{v}_r)}{\partial r} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial \theta} = \frac{1}{r} (-\phi'(r) k \sin(k\theta) + \phi'(r) k \sin(k\theta)) = 0$$

so the flow is indeed incompressible.

Since  $\xi$  is a function of  $\theta$  it is obvious that the only way to insure no-slip boundary conditions for any  $\theta$  value is to have all the following four conditions satisfied

$$\phi(R_1) = \phi'(R_1) = 0 \tag{12.247}$$

$$\phi(R_2) = \phi'(R_2) = 0 \tag{12.248}$$



We could then choose

$$\phi(r) = (r - R_1)^2(r - R_2)^2\mathcal{F}(r) \quad (12.249)$$

$$\phi'(r) = 2(r - R_1)(r - R_2)^2\mathcal{F}(r) + 2(r - R_1)^2(r - R_2)\mathcal{F}(r) + (r - R_1)^2(r - R_2)^2\mathcal{F}'(r) \quad (12.250)$$

which are indeed identically zero on both boundaries. Here  $\mathcal{F}(r)$  is any (smooth enough) function of  $r$ . A generic form for  $\Psi$  could then be

$$\boxed{\Psi(r, \theta) = (r - R_1)^2(r - R_2)^2\mathcal{F}(r) \cos(k\theta)}$$

In what follows I take  $\mathcal{F}(r) = 1$  for simplicity. Then

$$\begin{aligned} \phi(r) &= (r - R_1)^2(r - R_2)^2 \\ &= (r^2 - 2rR_1 + R_1^2)(r^2 - 2rR_2 + R_2^2) \\ &= \underbrace{1}_a r^4 + \underbrace{(-2R_1 - 2R_2)}_b r^3 + \underbrace{(R_1^2 + R_2^2 + 4R_1R_2)}_c r^2 + \underbrace{(-2R_1R_2^2 - 2R_1^2R_2)}_d r + \underbrace{R_1^2R_2^2}_e \\ &= ar^4 + br^3 + cr^2 + dr + e \end{aligned} \quad (12.251)$$

and then

$$\phi'(r) = 2(r - R_1)(r - R_2)^2 + 2(r - R_1)^2(r - R_2) \quad (12.252)$$

$$= 2(r - R_1)(r - R_2)(r - R_2 + r - R_1) \quad (12.253)$$

$$= 4(r - R_1)(r - R_2) \left( r - \frac{R_1 + R_2}{2} \right) \quad (12.254)$$

$$\phi''(r) = 8 \left( r - \frac{R_1 + R_2}{2} \right)^2 + 4(r - R_1)(r - R_2) \quad (12.255)$$

$$\phi'''(r) = 24 \left( r - \frac{R_1 + R_2}{2} \right) \quad (12.256)$$

$$\phi''''(r) = 24 \quad (12.257)$$

$$\mathbf{v}_r(r, \theta) = -\frac{1}{r}(r - R_1)^2(r - R_2)^2 k \sin(k\theta) \quad (12.258)$$

$$\mathbf{v}_\theta(r, \theta) = -4(r - R_1)(r - R_2) \left( r - \frac{R_1 + R_2}{2} \right) \cos(k\theta) \quad (12.259)$$

In the end the functions  $\phi$  and  $\xi$  are of the form:

$$\xi(\theta) = \cos(k\theta) \quad (12.260)$$

$$\phi(r) = ar^4 + br^3 + cr^2 + dr + e \quad (12.261)$$

with

$$\begin{aligned} \xi'(\theta) &= -k \sin(k\theta) \\ \xi''(\theta) &= -k^2 \cos(k\theta) = -k^2 \xi(\theta) \\ \xi'''(\theta) &= k^3 \sin(k\theta) \\ \xi''''(\theta) &= k^4 \cos(k\theta) = k^4 \xi(\theta) \\ \phi'(r) &= 4ar^3 + 3br^2 + 2cr + d \\ \phi''(r) &= 12ar^2 + 6br + 2c \\ \phi'''(r) &= 24ar + 6b \\ \phi''''(r) &= 24a \end{aligned}$$

## Finding the pressure and density fields

We start from the relationship

$$A_\theta = \frac{\partial^2 \mathbf{v}_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial \mathbf{v}_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \mathbf{v}_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial \mathbf{v}_r}{\partial \theta} - \frac{\mathbf{v}_\theta}{r^2} - \frac{1}{r} \frac{\partial p}{\partial \theta} = -\phi''' \xi - \frac{\phi'' \xi}{r} + \frac{\phi'(\xi - \xi'')}{r^2} + \frac{2\phi \xi''}{r^3} - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0$$

Then, after multiplying all by  $r^3$  we have

$$\begin{aligned} & r^2 \frac{\partial p}{\partial \theta} \\ &= -r^3 \phi''' \xi - r^2 \phi'' \xi + r \phi'(\xi - \xi'') + 2\phi \xi'' \\ &= \cos(k\theta) \left[ -r^3(24ar + 6b) - r^2(12ar^2 + 6br + 2c) + r(4ar^3 + 3br^2 + 2cr + d)(1 + k^2) + 2(ar^4 + br^3 + cr^2 + dr + e)k^3 \sin(k\theta) \right. \\ &= \cos(k\theta) \left[ -24ar^4 - 6br^3 - 12ar^4 - 6br^3 - 2cr^2 + (4ar^4 + 3br^3 + 2cr^2 + dr)(1 + k^2) - 2k^2(ar^4 + br^3 + cr^2 + dr + e) \right. \\ &= \cos(k\theta) \left[ (-24 - 12 + 4 + 4k^2 - 2k^2)ar^4 + (-6 - 6 + 3 + 3k^2 - 2k^2)br^3 + (-2 + 2 + 2k^2 - 2k^2)cr^2 + (1 - k^2)dr - 2k^2e \right] \\ &= \cos(k\theta) \left[ 2(k^2 - 16)ar^4 + (k^2 - 9)br^3 + (1 - k^2)dr - 2k^2e \right] \end{aligned}$$

or,

$$\frac{\partial p}{\partial \theta} = \cos(k\theta) \left[ 2(k^2 - 16)ar^2 + (k^2 - 9)br + (1 - k^2)\frac{d}{r} - 2k^2\frac{e}{r^2} \right]$$

i.e. after integration with respect to  $\theta$ :

$$p(r, \theta) = \frac{1}{k} \sin(k\theta) \left[ 2(k^2 - 16)ar^2 + (k^2 - 9)br + (1 - k^2)\frac{d}{r} - 2k^2\frac{e}{r^2} \right] + f(r)$$

For simplicity we set  $f(r) = 0$  and then

$$\frac{\partial p}{\partial r} = \frac{1}{k} \sin(k\theta) \left[ 4(k^2 - 16)ar + (k^2 - 9)b - (1 - k^2)\frac{d}{r^2} + 4k^2\frac{e}{r^3} \right]$$

and since we will need it later:

$$r^3 \frac{\partial p}{\partial r} = \frac{1}{k} \sin(k\theta) \left[ 4(k^2 - 16)ar^4 + (k^2 - 9)br^3 - (1 - k^2)dr + 4k^2e \right]$$

We now turn to

$$A_r = \frac{\phi'' \xi'}{r} + \frac{\phi' \xi'}{r^2} + \frac{\phi \xi'''}{r^3} - \frac{\partial p}{\partial r} = \rho g_r$$

or, after multiplying both sides by  $r^3$ :

$$\begin{aligned} & r^2 \phi'' \xi' + r \phi' \xi' + \phi \xi''' - r^3 \frac{\partial p}{\partial r} \\ &= r^2(12ar^2 + 6br + 2c)(-k \sin(k\theta)) + r(4ar^3 + 3br^2 + 2cr + d)(-k \sin(k\theta)) + (ar^4 + br^3 + cr^2 + dr + e)k^3 \sin(k\theta) \\ &\quad - \frac{1}{k} \sin(k\theta) \left[ 4(k^2 - 16)ar^4 + (k^2 - 9)br^3 - (1 - k^2)dr + 4k^2e \right] \\ &= \sin(k\theta) \left[ -k^2(12ar^4 + 6br^3 + 2cr^2) - k^2(4ar^4 + 3br^3 + 2cr^2 + dr) + k^4(ar^4 + br^3 + cr^2 + dr + e) \right. \\ &\quad \left. - \sin(k\theta) \left[ 4(k^2 - 16)ar^4 + (k^2 - 9)br^3 - (1 - k^2)dr + 4k^2e \right] \right] \\ &= \sin(k\theta) \left[ (-12k^2 - 4k^2 + k^4 - 4(k^2 - 16))ar^4 + (-6k^2 - 3k^2 + k^4 - (k^2 - 9))br^3 \right. \\ &\quad \left. + (-2k^2 - 2k^2 + k^4)cr^2 + (-k^2 + k^4 + (1 - k^2))dr + (k^4 - 4k^2e) \right] \\ &= \sin(k\theta) \left[ (k^4 - 20k^2 + 64)ar^4 + (k^4 - 10k^2 + 9)br^3 + (k^4 - 4k^2)cr^2 + (k^4 - 2k^2 + 1)dr + k^2(k^2 - 4)e \right] \end{aligned}$$

So, assuming  $g_r = 1$ ,

$$\rho(r, \theta) = \frac{\sin(k\theta)}{k} \frac{Aar^4 + Bbr^3 + Ccr^2 + Ddr + Ee}{r^3}$$

with

$$\begin{aligned}
A &= (k^2 - 4)(k^2 - 16) \\
B &= k^4 - 10k^2 + 9 \\
C &= k^2(k^2 - 4) \\
D &= (k^2 - 1)(k^2 - 1) \\
E &= k^2(k^2 - 4)
\end{aligned} \tag{12.266}$$

In the end:

$$\mathbf{v}_r(r, \theta) = -\frac{1}{r}(\textcolor{brown}{a}r^4 + \textcolor{blue}{b}r^3 + \textcolor{teal}{c}r^2 + \textcolor{red}{d}r + \textcolor{violet}{e})k \sin(k\theta) \tag{12.267}$$

$$\mathbf{v}_\theta(r, \theta) = -(4ar^3 + 3br^2 + 2cr + d) \cos(k\theta) \tag{12.268}$$

$$p(r, \theta) = \frac{1}{k} \sin(k\theta) \left[ 2(k^2 - 16)ar^2 + (k^2 - 9)br + (1 - k^2)\frac{d}{r} - 2k^2\frac{e}{r^2} \right] \tag{12.269}$$

$$\rho(r, \theta) = \frac{\sin(k\theta)}{k} \frac{A\textcolor{brown}{a}r^4 + B\textcolor{blue}{b}r^3 + C\textcolor{teal}{c}r^2 + D\textcolor{red}{d}r + E\textcolor{violet}{e}}{r^3} \tag{12.270}$$

### Finding the density field - alternate & easier take

We start this time from

$$\frac{\partial(rA_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} = -\frac{\partial \rho g_r}{\partial \theta}$$

or, with  $g_r = 1$ ,

$$\begin{aligned}
-2\phi''' \xi - \phi'''' \xi r + \frac{\phi''}{r}(\xi - \xi'') + \frac{\phi'}{r^2}(-\xi + 3\xi'') + \frac{\phi}{r^3}(-4\xi'') - \frac{\partial^2 p}{\partial \theta \partial r} - \frac{\phi'' \xi''}{r} - \frac{\phi' \xi'''}{r^2} - \frac{\phi \xi''''}{r^3} + \frac{\partial^2 p}{\partial r \partial \theta} &= -\frac{\partial \rho}{\partial \theta} \\
-2\phi''' \xi - \phi'''' \xi r + \frac{\phi''}{r}(\xi - \xi'') + \frac{\phi'}{r^2}(-\xi + 3\xi'') + \frac{\phi}{r^3}(-4\xi'') - \frac{\phi'' \xi''}{r} - \frac{\phi' \xi'''}{r^2} - \frac{\phi \xi''''}{r^3} &= -\frac{\partial \rho}{\partial \theta}
\end{aligned}$$

Then we note that  $\xi'' = -k^2 \xi$  and  $\xi'''' = k^4 \xi$  so that

$$\begin{aligned}
-2\phi''' \xi - \phi'''' \xi r + \frac{\phi''}{r}(\xi + k^2 \xi) + \frac{\phi'}{r^2}(-\xi - 3k^2 \xi) + \frac{\phi}{r^3}(4k^2 \xi) - \frac{-k^2 \phi'' \xi}{r} - \frac{-k^2 \phi' \xi}{r^2} - \frac{k^4 \phi \xi}{r^3} &= -\frac{\partial \rho}{\partial \theta} \\
\xi \left[ -2\phi''' - \phi'''' r + \frac{\phi''}{r}(1 + k^2) + \frac{\phi'}{r^2}(-1 - 3k^2) + \frac{\phi}{r^3}(4k^2) - \frac{-k^2 \phi''}{r} - \frac{-k^2 \phi'}{r^2} - \frac{k^4 \phi}{r^3} \right] &= -\frac{\partial \rho}{\partial \theta} \\
\xi(\theta) \left[ -2\phi''' - \phi'''' r + \frac{\phi''}{r}(1 + 2k^2) + \frac{\phi'}{r^2}(-1 - 2k^2) + \frac{\phi}{r^3}(4k^2 - k^4) \right] &= -\frac{\partial \rho}{\partial \theta}
\end{aligned}$$

i.e.

$$\rho(r, \theta) = -\frac{1}{k} \sin(k\theta) \frac{\mathcal{G}(r)}{r^3}$$

with

$$\begin{aligned}
\mathcal{G}(r) &= -\phi'''' r^4 - 2\phi''' r^3 + \phi'' r^2(1 + 2k^2) + \phi' r(-1 - 2k^2) + \phi(4k^2 - k^4) \\
&= -24ar^4 - 2(24ar + 6b)r^3 + (12ar^2 + 6br + 2c)r^2(1 + 2k^2) \\
&\quad + (4ar^3 + 3br^2 + 2cr + d)r(-1 - 2k^2) + (ar^4 + br^3 + cr^2 + dr + e)(4k^2 - k^4) \\
&= -24\textcolor{brown}{a}r^4 - 2(24\textcolor{brown}{a}r^4 + 6\textcolor{blue}{b}r^3) + (12\textcolor{brown}{a}r^4 + 6\textcolor{blue}{b}r^3 + 2\textcolor{teal}{c}r^2)(1 + 2k^2) \\
&\quad + (4\textcolor{brown}{a}r^4 + 3\textcolor{blue}{b}r^3 + 2\textcolor{teal}{c}r^2 + \textcolor{red}{d}r)(-1 - 2k^2) + (\textcolor{brown}{a}r^4 + \textcolor{blue}{b}r^3 + \textcolor{teal}{c}r^2 + \textcolor{red}{d}r + \textcolor{violet}{e})(4k^2 - k^4) \\
&= -A\textcolor{brown}{a}r^4 - B\textcolor{blue}{b}r^3 - C\textcolor{teal}{c}r^2 - D\textcolor{red}{d}r - E\textcolor{violet}{e}
\end{aligned}$$

with

$$\begin{aligned}
-A &= -24 - 48 + 12(1 + 2k^2) + 4(-1 - 2k^2) + (4k^2 - k^4) \\
&= -24 - 48 + 12 + 24k^2 - 4 - 8k^2 + 4k^2 - k^4 \\
&= -64 + 20k^2 - k^4 \\
&= -(k^2 - 4)(k^2 - 16) \\
-B &= -12 + 6(1 + 2k^2) + 3(-1 - 2k^2) + (4k^2 - k^4) \\
&= -12 + 6 + 12k^2 - 3 - 6k^2 + 4k^2 - k^4 \\
&= -(k^4 - 10k^2 + 9) \\
-C &= 2(1 + 2k^2) + 2(-1 - 2k^2) + (4k^2 - k^4) \\
&= 2 + 4k^2 - 2 - 4k^2 + 4k^2 - k^4 \\
&= 4k^2 - k^4 \\
&= -k^2(k^2 - 4) \\
-D &= (-1 - 2k^2) + (4k^2 - k^4) \\
&= -1 - 2k^2 + 4k^2 - k^4 \\
&= -1 + 2k^2 - k^4 \\
&= -(k^2 - 1)(k^2 - 1) \\
-E &= 4k^2 - k^4 \\
&= -k^2(k^2 - 4)
\end{aligned}$$

so in the end we recover

$$\begin{aligned}
\rho(r, \theta) &= \frac{1}{k} \sin(k\theta) \frac{A \textcolor{brown}{a} r^4 + B \textcolor{blue}{b} r^3 + C \textcolor{teal}{c} r^2 + D \textcolor{red}{d} r + E \textcolor{red}{e}}{r^3} \\
A &= (k^2 - 4)(k^2 - 16) \\
B &= k^4 - 10k^2 + 9 \\
C &= k^2(k^2 - 4) \\
D &= (k^2 - 1)(k^2 - 1) \\
E &= k^2(k^2 - 4)
\end{aligned}$$

The pressure is obtained as presented before.

### Root mean square velocity

$$\begin{aligned}
\mathbf{v}_{rms}^2 &= \iint (\mathbf{v}_r^2 + \mathbf{v}_\theta^2) r dr d\theta \\
&= \iint \left[ \left( \frac{\phi \xi'}{r} \right)^2 + (-\phi' \xi)^2 \right] r dr d\theta \\
&= \iint \left( \frac{\phi \xi'}{r} \right)^2 r dr d\theta + \iint (-\phi' \xi)^2 r dr d\theta \\
&= \underbrace{\int_{R_1}^{R_2} \left( \frac{\phi}{r} \right)^2 r dr}_{I_1} \underbrace{\int_0^{2\pi} (\xi')^2 d\theta}_{I_2} + \underbrace{\int_{R_1}^{R_2} (\phi')^2 r dr}_{I_3} \underbrace{\int_0^{2\pi} \xi^2 d\theta}_{I_4}
\end{aligned}$$

$$\begin{aligned}
I_1 &= \\
I_2 &= \\
I_3 &= \\
I_4 &=
\end{aligned}$$

Unfinished! Also compute strain rate tensor, stress, total mass, etc ...

## Free slip boundary conditions

what follows needs to be checked!!!

Before postulating the form of  $\phi(r)$ , let us now turn to the boundary conditions that the flow must fulfill, i.e. free-slip on both boundaries. Two conditions must be met:

- $\mathbf{v} \cdot \mathbf{n} = 0$  (no flow through the boundaries) which yields  $u(r = R_1) = 0$  and  $u(r = R_2) = 0$ , :

$$\frac{1}{r} \frac{\partial \Psi}{\partial \theta}(r = R_1, R_2) = 0 \quad \forall \theta$$

which gives us the first constraint since  $\Psi(r, \theta) = \phi(r)\xi(\theta)$ :

$$\phi(r = R_1) = \phi(r = R_2) = 0$$

- $(\boldsymbol{\sigma} \cdot \mathbf{n}) \times \mathbf{n} = \mathbf{0}$  (the tangential stress at the boundary is zero) which imposes:  $\sigma_{\theta r} = 0$ , with

$$\sigma_{\theta r} = 2\eta \cdot \frac{1}{2} \left( \frac{\partial v}{\partial r} - \frac{v}{r} + \frac{1}{r} \frac{\partial u}{\partial \theta} \right) = \eta \left( \frac{\partial}{\partial r} \left( -\frac{\partial \Psi}{\partial r} \right) - \frac{1}{r} \left( -\frac{\partial \Psi}{\partial r} \right) + \frac{1}{r} \frac{\partial}{\partial \theta} \left( \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \right) \right)$$

Finally  $\Psi$  must fulfill (on the boundaries!):

$$\begin{aligned}
-\frac{\partial^2 \Psi}{\partial r^2} + \frac{1}{r} \frac{\partial \Psi}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \Psi}{\partial \theta^2} &= 0 \\
-\phi'' \xi + \frac{1}{r} \phi' \xi + \frac{1}{r^2} \phi \xi'' &= 0
\end{aligned}$$

or,

$$-\phi'' + \frac{1}{r} \phi' - k^2 \frac{1}{r^2} \phi = 0$$

Note that this equation is a so-called Euler Differential Equation<sup>9</sup>. Since we are looking for a solution  $\phi$  such that  $\phi(R_1) = \phi(R_2) = 0$  then the 3rd term of the equation above is by definition zero on the boundaries. We have to ensure the following equality on the boundary:

$$-\phi'' + \frac{1}{r} \phi' = 0 \quad \text{for } r = R_1, R_2$$

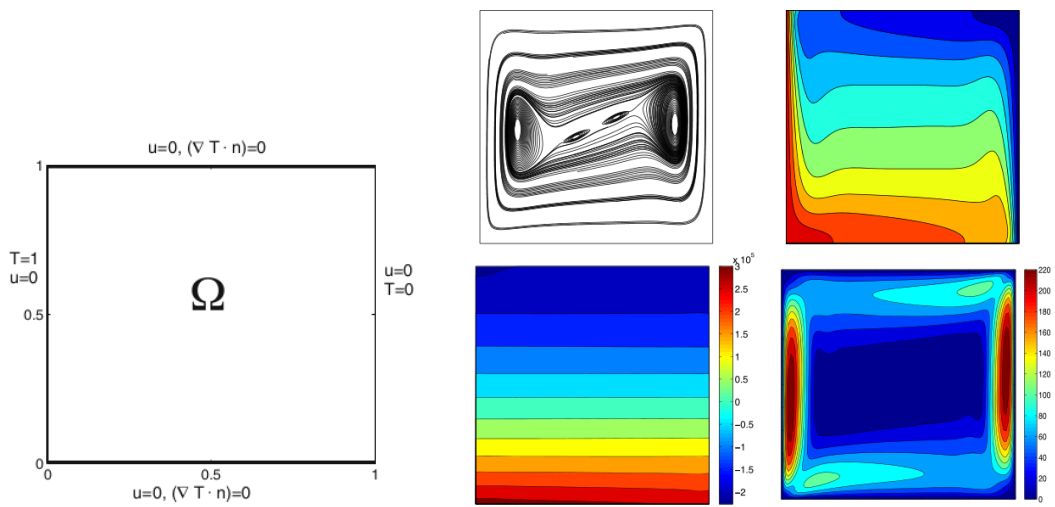
The solution of this ODE is of the form  $\phi(r) = ar^2 + b$  and it becomes evident that it cannot satisfy  $\phi(r = R_1) = \phi(r = R_2) = 0$ .

Separation of variables leads to solutions which cannot fulfill the free slip boundary conditions

<sup>9</sup><http://mathworld.wolfram.com/EulerDifferentialEquation.html>

### 12.2.18 Rayleigh-Bénard convection for silicon oil

This originates in Jenkins *et al.* (2014) [643].



## 12.2.19 Rayleigh-Taylor experiment of van Keken *et al.* (1997)

benchmark\_vaks97.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/benchmark\\_vaks97](https://github.com/cedrict/fieldstone/tree/master/images/benchmark_vaks97)

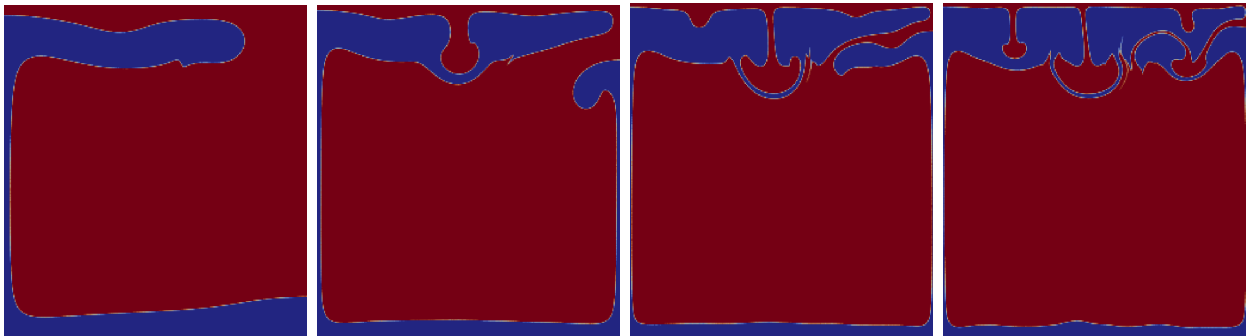
This numerical experiment was first presented in van Keken, King, Schmeling, Christensen, Neumeister, and Doin [1309] (1997). It consists of an isothermal Rayleigh-Taylor instability in a two-dimensional box of size  $L_x = 0.9142$  and  $L_y = 1$ .

Two Newtonian fluids are present in the system: the buoyant layer is placed at the bottom of the box and the interface between both fluids is given by

$$y(x) = 0.2 + 0.02 \cos\left(\frac{\pi x}{L_x}\right) \quad (12.271)$$

The bottom fluid is parametrised by its density  $\rho_1 = 1000$  and its viscosity  $\eta_1$ , while the layer above is parametrised by  $\rho_2 = 1010$  and  $\eta_2 = 100$ . This experiment is to be carried out for various viscosity contrasts between the two layers, i.e.  $\eta_1 = \{1, 10, 100\}$ .

No-slip boundary conditions are applied at the bottom and at the top of the box while free-slip boundary conditions are applied on the sides. Gravity is pointing downwards with  $|\vec{g}| = 10$ .



Example of time evolution obtained with the ELEFANT code for a 256x256 grid with the particle-in-cell method.

in images/benchmark\_vaks97/

One can measure the following quantities:

- the root mean square velocity  $\mathbf{v}_{rms}$  in the domain as a function of time:

$$\mathbf{v}_{rms} = \sqrt{\frac{1}{L_x L_y} \int \int |\vec{v}|^2 dx dy} \quad (12.272)$$

- the maximum (or local maxima) of the  $\mathbf{v}_{rms}$  and its (their) corresponding time(s)
- the growth rate of the instability at  $t = 0$ . From linear stability analysis, the analytical growth rate can be calculated [1038, 1037]:  $\gamma_{th} = 0.01094019$ , which is valid for an infinitesimal perturbation. For each model run, the growth rate  $\gamma$  is measured by fitting the  $\mathbf{v}_{rms}$  and/or (?) the maximum vertical velocity measurements for a short time  $t$ .

- the total mass of the system  $M(t)$  as a function of time. Since there is no chemical diffusion in the system (pure advection equation) the amount of material in the system is to remain constant, and therefore its mass too.

$$M(t) = \int \int \rho(x, y, t) dx dy \quad (12.273)$$

Given the layout described in the previous paragraph, the exact analytical initial mass  $M_0$  of the system is given by

$$M_0 = 0.9142 \times (0.2 \times 1000 + 0.8 \times 1010) = 921.5136$$

The average density is then

$$\langle \rho \rangle_0 = \frac{M_0}{L_x L_y} = 1008$$

We will then measure the relative mass error as a function of time

$$\delta M(t) = \frac{M(t) - M_0}{M_0}$$

which is equal to

$$\langle \delta \rho \rangle(t) = \frac{\langle \rho \rangle(t) - \langle \rho \rangle_0}{\langle \rho \rangle_0}$$

- the length of the interface between the fluids. At startup it is given by

$$\mathcal{L}(0) = \int_0^L \sqrt{1 + (dy/dx)^2} dx$$

with  $y(x) = 0.2 + 0.02 \cos(\pi x/L)$ . Using WolframAlpha, we find

$$\mathcal{L}(0) = \frac{L}{\pi} \int_0^\pi \sqrt{1 + \left(0.02 \frac{\pi}{L} \sin(\pi x/L)\right)^2} dx \simeq 0.9152786349$$

**Instantaneous results** . Results obtained with **STONE** ?? ( $Q_2 \times Q_1$  element) and **STONE** ?? ( $P_2^+ \times P_{-1}$  element), both with mesh fitted so as to follow the interface between both fluids.



|             |                                    | $\eta_1 = 100$ | $\eta_1 = 10$ | $\eta_1 = 1$ |
|-------------|------------------------------------|----------------|---------------|--------------|
| $\min(u)$   | Stone 25                           |                |               | -4563.5      |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\max(u)$   | Stone 25                           |                |               | 1054.03      |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\min(v)$   | Stone 25                           |                |               |              |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\max(v)$   | Stone 25                           |                |               | 2070.9       |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\max( v )$ | Stone 25                           | 422.125        |               | 4563.67      |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $v_{rms}$   | Stone 25                           | 185.2947       |               | 1441.87      |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\min(p)$   | Stone 25                           | -5048.16       |               | -5048.7345   |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |
| $\max(p)$   | Stone 25                           | 5033.6947      |               | 5032.329     |
|             | Stone 93 ( $P_2^+ \times P_{-1}$ ) |                |               |              |
|             | Aspect                             |                |               |              |

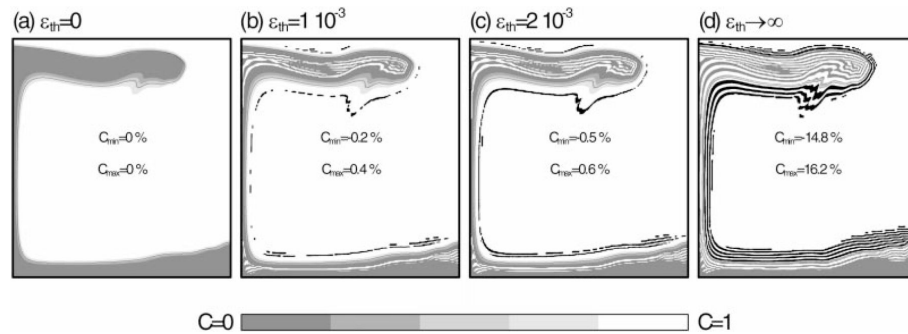
| Code   | Grid                     | Method       | Growth Rate $\gamma$ | t (max vrms) | max (vrms) | publicat |
|--------|--------------------------|--------------|----------------------|--------------|------------|----------|
| HS     | 41x41                    | extrapolated | 0.010996             | 211.1        | 0.0030958  | [1309]   |
|        | 61x64                    | extrapolated | 0.011109             | 209.17       | 0.0031022  |          |
|        | 81x81                    | extrapolated | 0.011177             | 208.99       | 0.0030916  |          |
| CND    | 32x32                    |              | 0.01106              | 208.4        | 0.003092   | [1309]   |
|        | 48x48                    |              | 0.01106              | 208.5        | 0.0030943  |          |
| SK     | 80x80                    |              | 0.01130              | 215.67       | 0.00299279 | [1309]   |
|        | 120x120                  |              | 0.01127              | 206.38       | 0.0028922  |          |
|        | 160x160                  |              | 0.01179              | 207.84       | 0.0028970  |          |
| PvK    | 30x30                    | splines      | 0.01185              | 213.38       | 0.00300    | [1309]   |
|        | 50x50                    | splines      | 0.01198              | 211.81       | 0.003016   |          |
|        | 80x80                    | splines      | 0.01207              | 210.75       | 0.003050   |          |
|        | 100x100                  | splines      | 0.01211              |              |            |          |
|        | 30x30                    | C1 element   | 0.01253              | 210.59       | 0.003100   |          |
|        | 80x80                    | C1 element   | 0.01225              | 207.05       | 0.003091   |          |
| ISMM   | 160x160                  |              | 0.00991              | 230.1        | 0.003093   | [1176]   |
| ISMM   | 120x120                  |              | 0.00998              | 226.1        | 0.003133   |          |
| MSOU   | 160x160                  |              | 0.00993              | 231.4        | 0.003085   |          |
| MSOU   | 120x120                  |              | 0.01020              | 227.6        | 0.003134   |          |
| FSOU   | 160x160                  |              | 0.01111              | 217.3        | 0.003118   |          |
| FSOU   | 120x120                  |              | 0.01159              | 213.1        | 0.003151   |          |
|        | 64x64                    | 5            | 0.01112              | 206.5        | 0.003041   | [1229]   |
|        |                          | 15           | 0.01117              | 208.8        | 0.003098   |          |
|        |                          | 40           | 0.01115              | 209.9        | 0.003110   |          |
|        | 128x128                  | 5            | 0.01113              | 208.1        | 0.003079   |          |
|        |                          | 15           | 0.01110              | 208.9        | 0.003097   |          |
|        |                          | 40           | 0.01109              | 209.2        | 0.003102   |          |
|        | 120x132                  | Level sets   | 0.01252              | 211.2        | 0.00301    | [1218]   |
|        | 67m res.                 |              |                      | 215.3        | 0.003106   | [238]    |
|        | 100m res.                |              |                      | 215.28       | 0.003101   |          |
| LaCoDe | 1808 elts (10754 dofs)   |              | 0.01221              | 215          | 0.003110   | [322]    |
|        | 7093 elts (2592 dofs)    |              | 0.01222              | 212          | 0.003080   |          |
|        | 17960 elts (107468 dofs) |              | 0.01222              | 211          | 0.003075   |          |

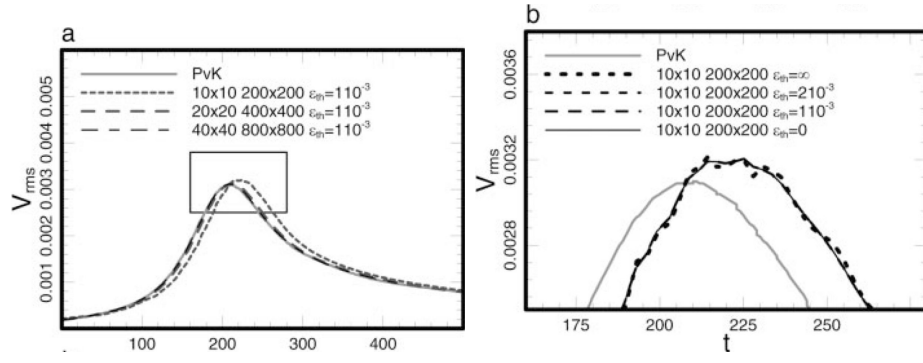
Selected Quantities for the Isoviscous Rayleigh-Taylor problem. HS: FDM, stream function formulation, particles. PvK: FEM, stream function, marker-chain.

SK: FEM, ConMan code, compositional field. CND: spline method, stream function formulation, particles.

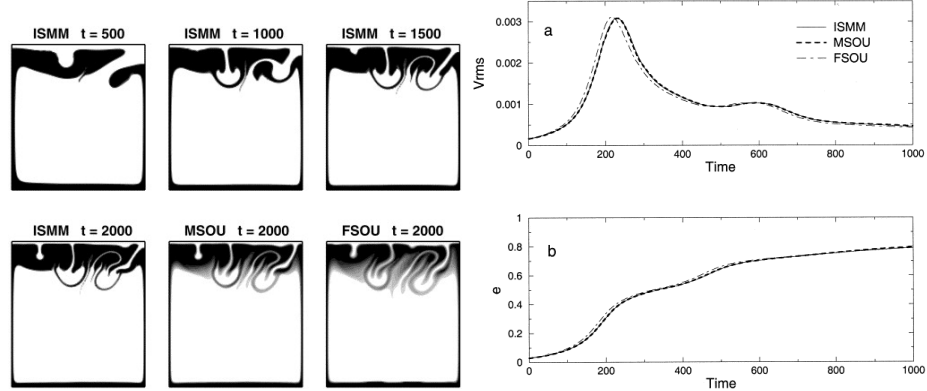
List of literature showcasing results of the van Keken *et al.* (1997) [1309] setup:

- de Smet *et al.* (2000) [323]. No table of results, only figures:

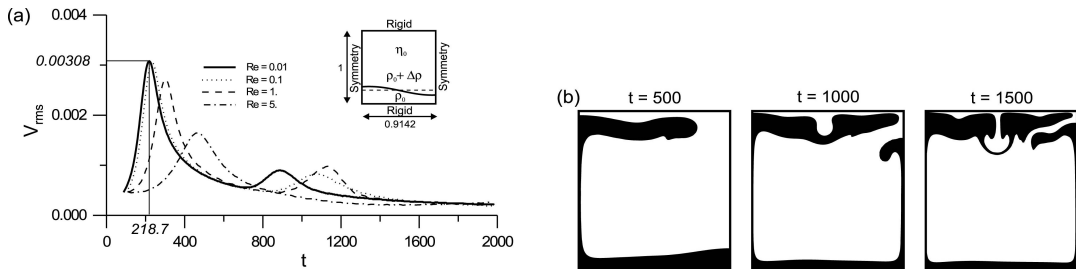




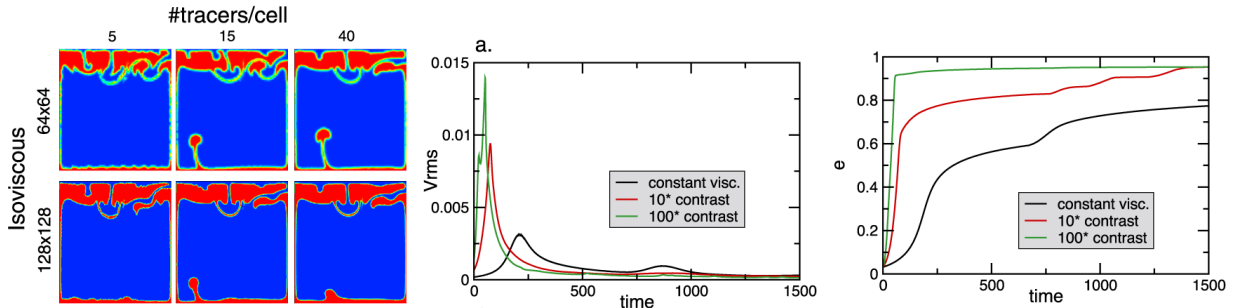
- Soboutia *et al.* (2001) [1176]. Results reported in table above.



- Babeyko *et al.* (2002) [35]. No table of results, only one figure:

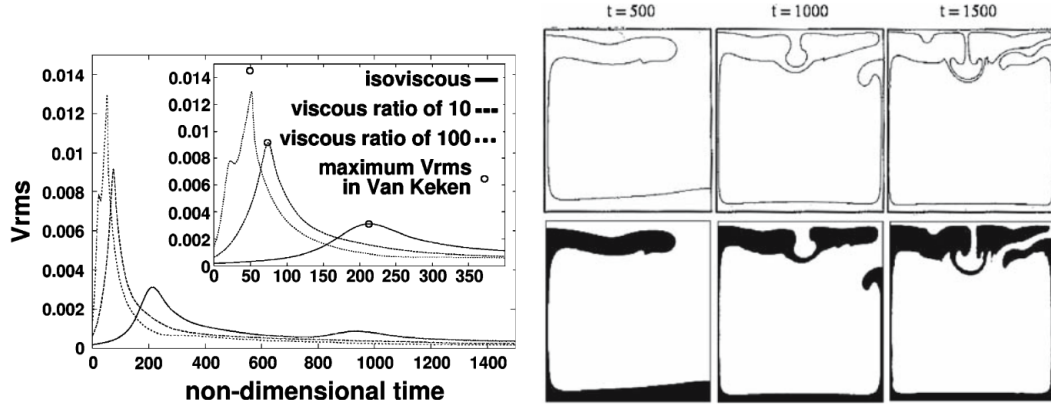


- Tackley & King (2003) [1229]. Performed with grid resolutions of 64x64 or 128x128 and with either 5, 15, or 40 tracers per cell (on average). Results reported in table above.

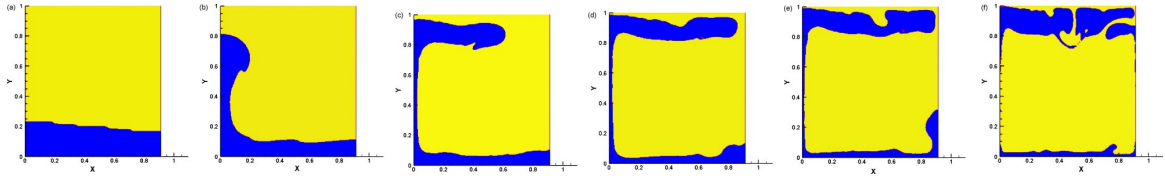


There are also results for non-isoviscous cases.

- Bourgouin *et al.* (2006) [124]. No table of results, only figures. Additional results for non-isoviscous in the paper.

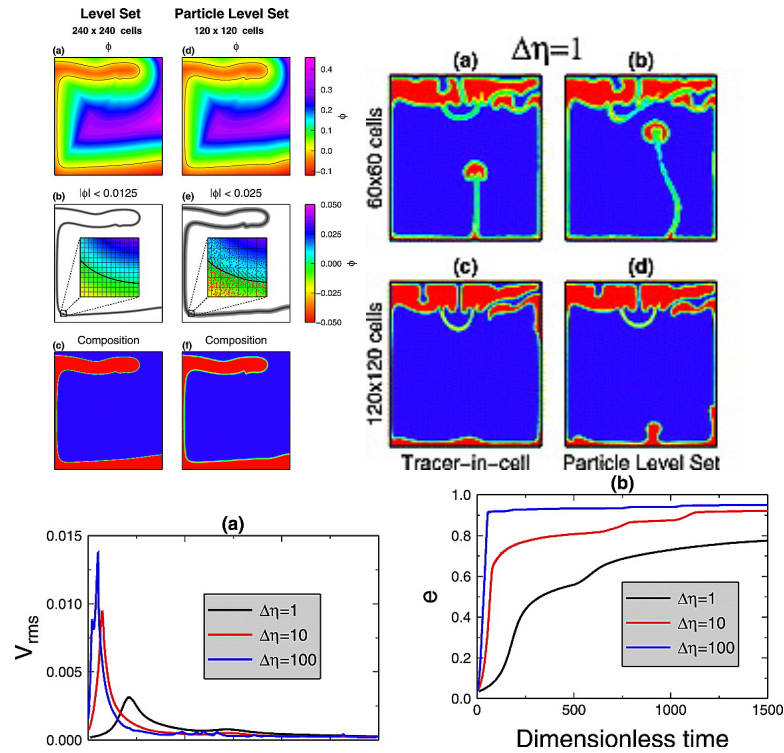


- Quinteros *et al.* (2009) [1029].

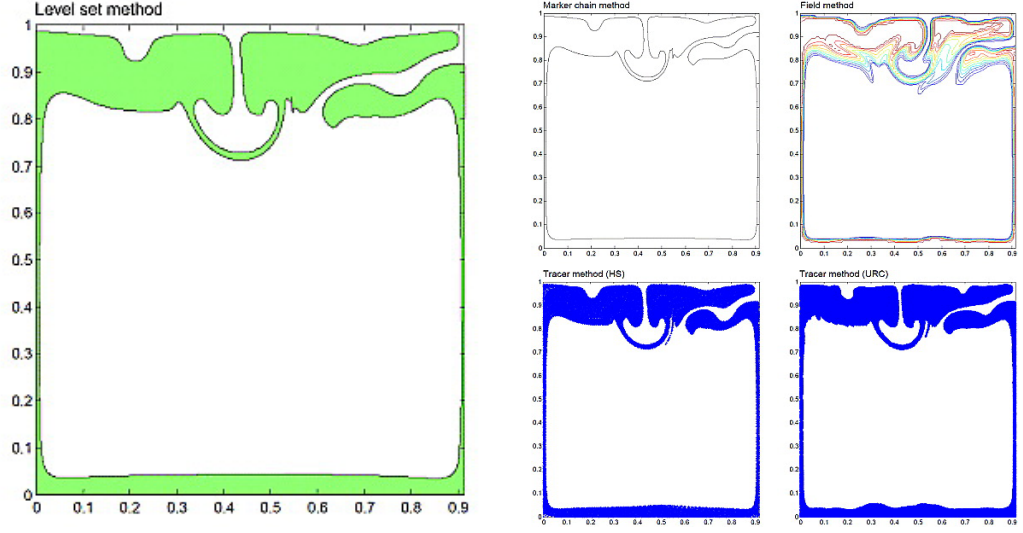


”Different snapshots from the domain evolution that are shown in Fig. 10 were compared with the ones published by van Keken *et al.* (1997). The evolution shown in this chapter and in the van Keken paper are identical for all the compared time steps.”

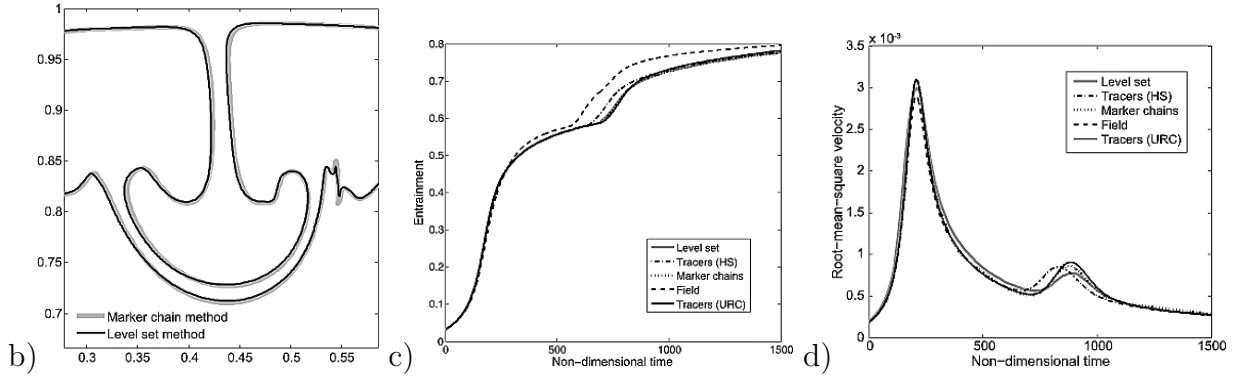
- Samuel & Evonuk (2010) [1103].



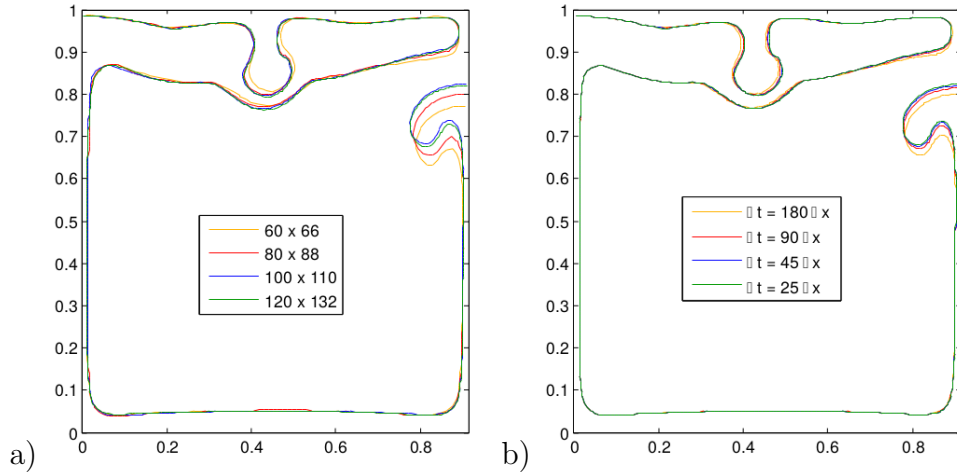
- Suckale *et al.* (2010) [1218].



The Rayleigh-Taylor instability at  $t = 1500$  computed by (left) the level set method on a  $300 \times 330$  grid, compared to the best results of (right) the four codes compared by van Keken *et al.*

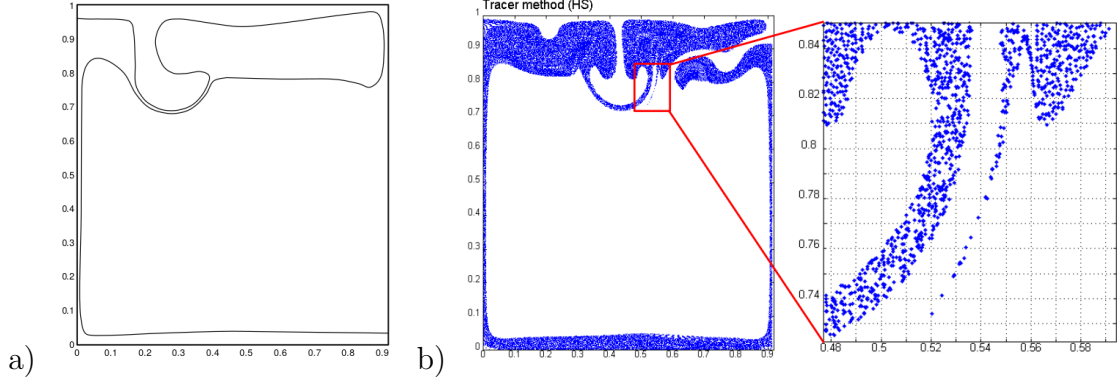


b) Detailed comparison of the level set (thin black line) and the marker chain approach (thick grey line) for the isothermal and isoviscous Rayleigh-Taylor instability at nondimensional time  $t = 1500$ . The plotted interfaces represent a zoom onto the instability descending from the top downwards in the middle of the box. The two methods yield an almost identical interface. c) Evolution of the entrainment of the buoyant fluid over time as computed by the five different codes. The level set computation was done on a  $160 \times 176$  grid. d) Evolution of the root mean square velocity of the interface over time as computed by the five different codes. The level set computation was done on a  $160 \times 176$  grid.



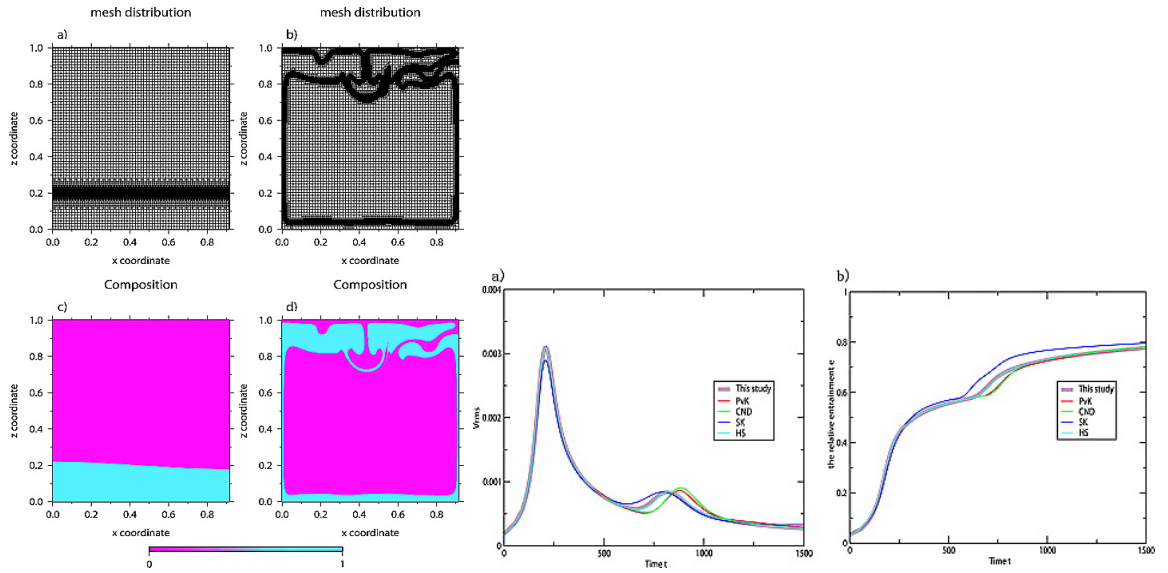
Taken from the supplementary material: a) Convergence test for the isothermal and isoviscous Rayleigh-Taylor instability, benchmark problem 3. A lack of convergence is easiest to identify during the phases of rapid rise of an instability. We illustrate this for the rise of the secondary instability on the right side of the box at time  $t=1000$  and four different grid sizes:  $60 \times 66$ ,  $80 \times 88$ ,  $100 \times 110$ , and  $120 \times 132$ . We observe convergence for grid sizes above  $100 \times 110$ .

b) Convergence test for the isothermal and isoviscous Rayleigh-Taylor instability, benchmark problem 3. We illustrate this convergence test for the rise of the secondary instability at time  $t = 1000$ . The four interfaces were computed based on the time steps:  $\Delta t = 180\Delta x$ ,  $\Delta t = 90\Delta x$ ,  $\Delta t = 45\Delta x$ , and  $\Delta t = 25\Delta x$ . We observe convergence for time steps  $\Delta t \leq 25\Delta x$ .



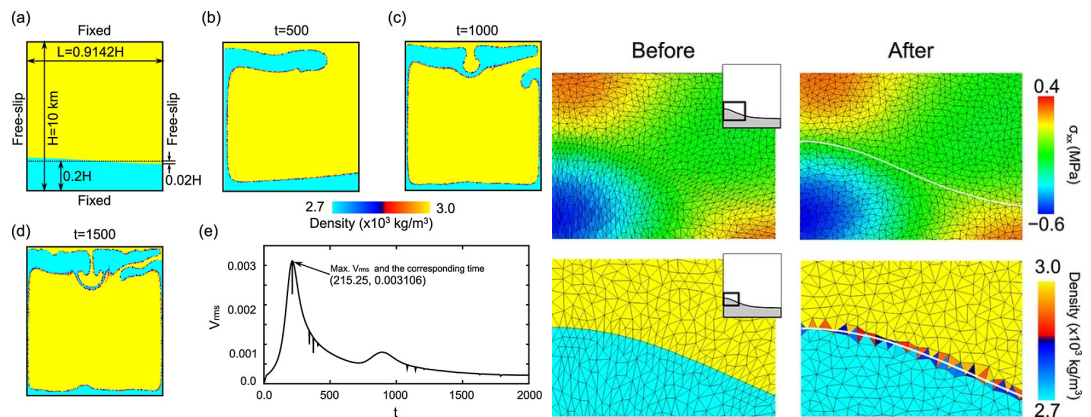
a) The isothermal Rayleigh-Taylor instability with viscosity contrast 10 at non-dimensional time  $t=500$ . The computation was done with a grid resolution of  $250 \times 275$ . b) The Rayleigh-Taylor instability as computed by the HS-tracer method at time  $t=1500$ . The equations of motion for this simulation were solved on an  $81 \times 81$  grid. The right panel is a zoom onto the peak located left of the descending instability. Each blue dot represents one particle and the grid represents a rough estimate of the scale at which the flow field is approximated correctly.

- Leng & Zhong (2011) [769].



Left: The mesh distribution and chemical composition for case RT1 at two different times: (a and c)  $t = 0$  and (b and d)  $t = 1500$ . Right: a) The root mean square velocity  $v_{rms}$  and (b) the relative entrainment of the buoyant material, with time for case RT1. The corresponding benchmark results from van Keken *et al.* [1997] are also plotted.

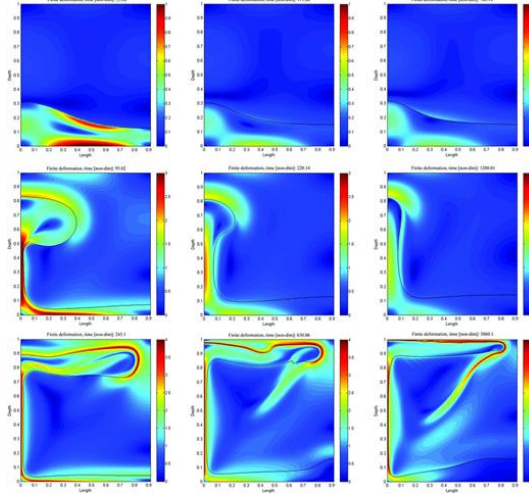
- Vynnytska *et al.* (2013) [1333]. Results only for 100 viscosity contrast
- Choi *et al.* (2013) [238].





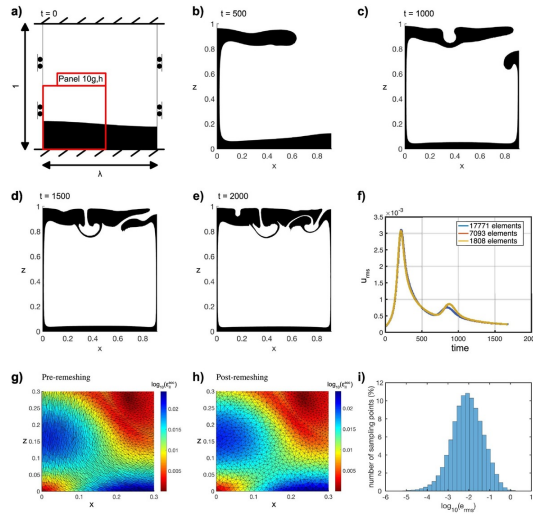
Left: Rayleigh-Taylor instability. (a) Model setup. Snapshots of the density at dimensionless time of (b) 500, (c) 1000, and (d) 1500. (e) Plot of  $v_{rms}$  versus dimensionless time,  $t$ . The resolution is about 0.6km. Right:  $\sigma_{xx}$  and density fields before and after the first remeshing with about 1km resolution. The white lines in the “After” images denote the original phase boundary before remeshing. The thick-lined box in the inset shows the location of the zoomed-in part of the domain.

- Fuchs & Schmeling (2013) [421].



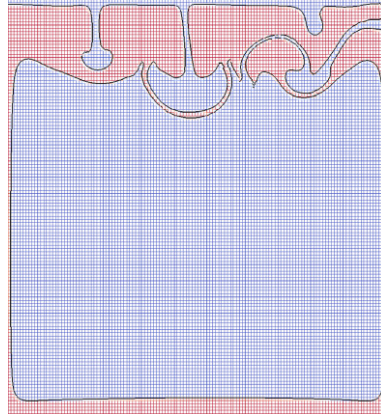
Finite deformation field for different stages of the diapirism with three different viscosity ratios and an initial non-dimensional thickness of the buoyant layer  $h_2 = 0.2$ . Left column: viscosity ratio  $m = 0.1$ . Middle column:  $m = 1$ . Right column:  $m = 10$ . Top row: pillow stage. Middle row: rising stage. Lower row: final stage.

- de Montserrat *et al.* (2019) [322].



a-e) Temporal evolution of the Rayleigh-Taylor instability. f) Evolution of  $v_{rms}$ . Remeshing of the domain is necessary when the mesh becomes highly distorted. Note that the red lines overlap with the blue line. g) Second invariant of the accumulated strain in a mesh with heavily distorted elements, and h) interpolated into a new high-quality mesh. i) Histogram showing the logarithm of the error between the accumulated square root of second invariant of the strain rate, pre- and post-remeshing.

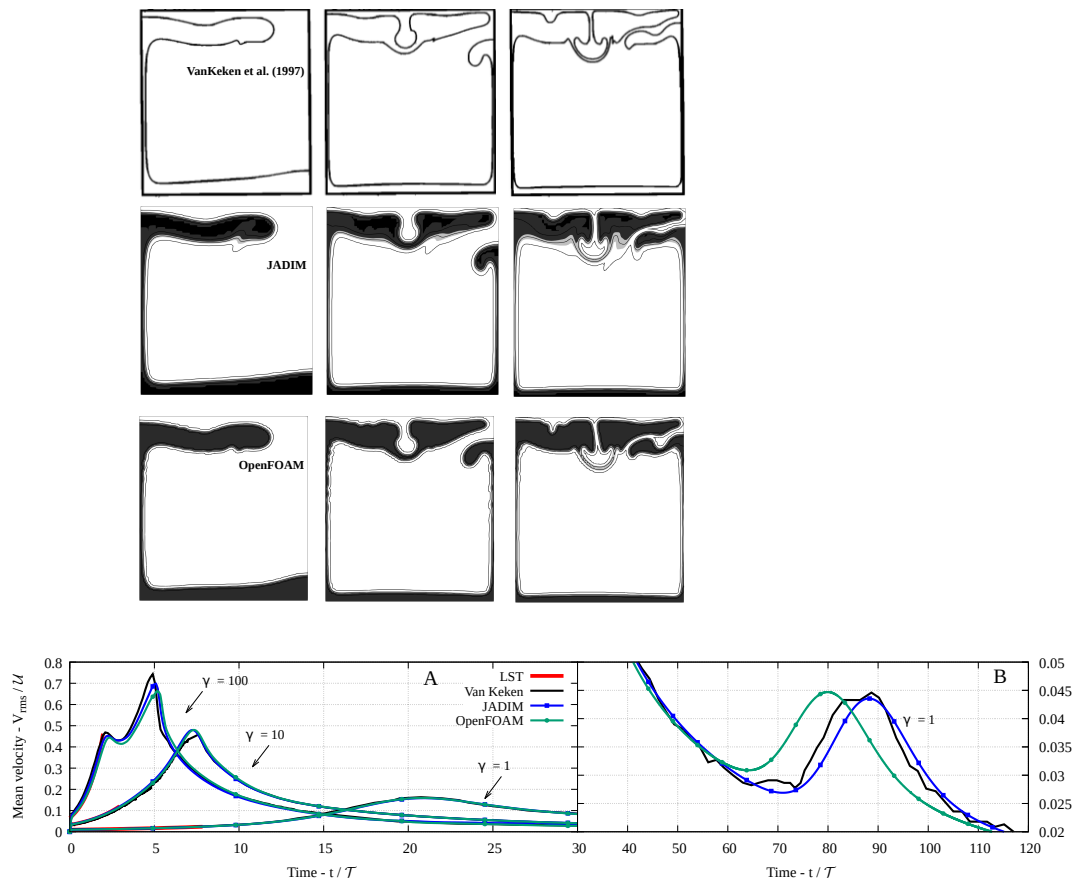
- Robey & Puckett (2019) [1079], Robey (2019) [1078] (PhD thesis).



Computed solution of the van Keken isoviscous Rayleigh-Taylor problem at time  $t = 2000$  on a uniform grid of  $128 \times 128$  cells.

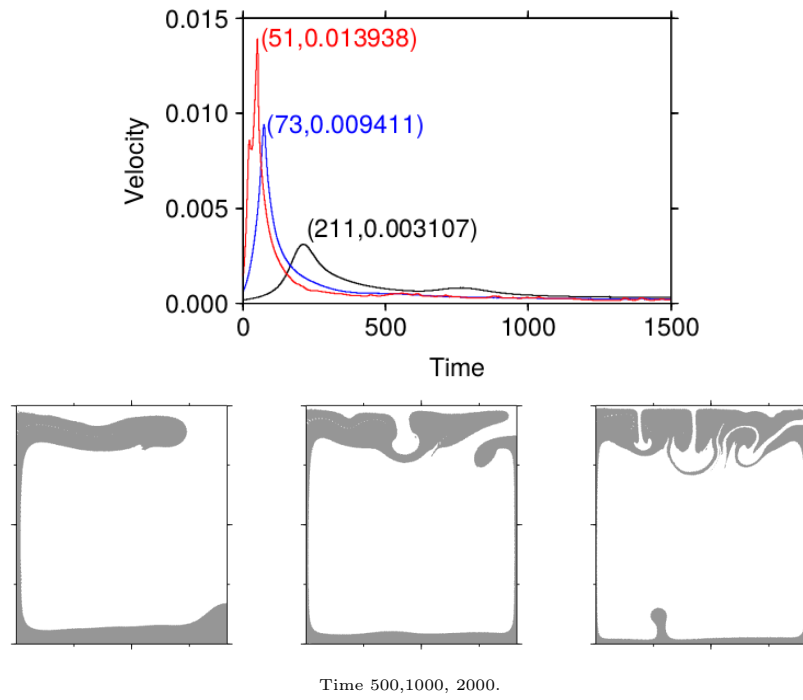
- Louis-Napoleon *et al.* (2020) [811].

Raw data available in `./images/benchmark_vaks97/louis_napoleon_etal.`

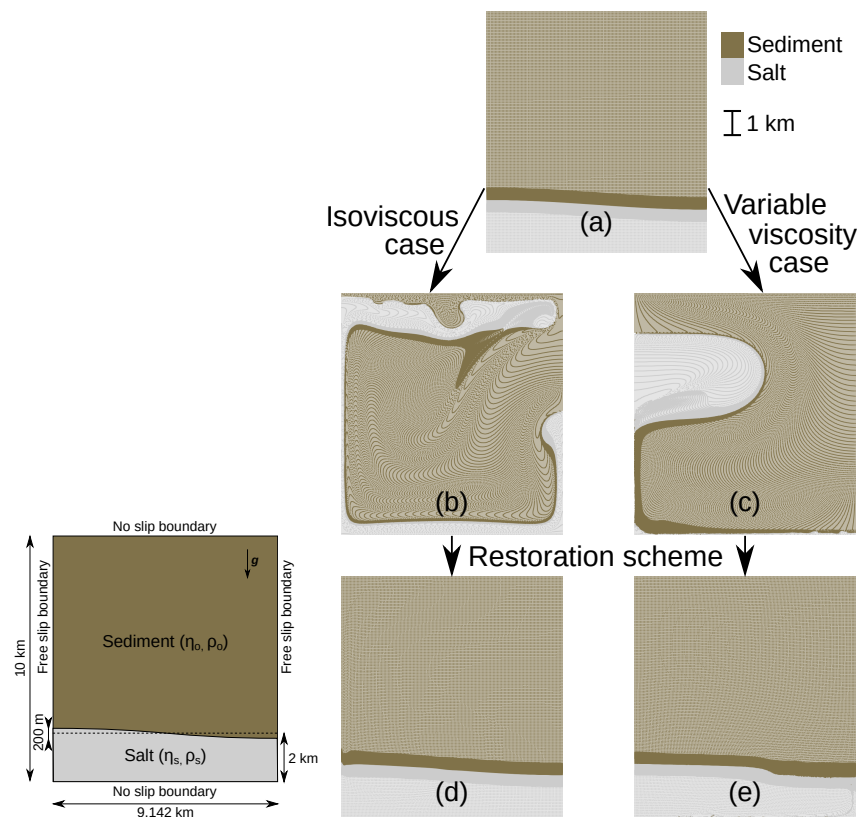


- Maierova (2012) [825] (phd thesis)

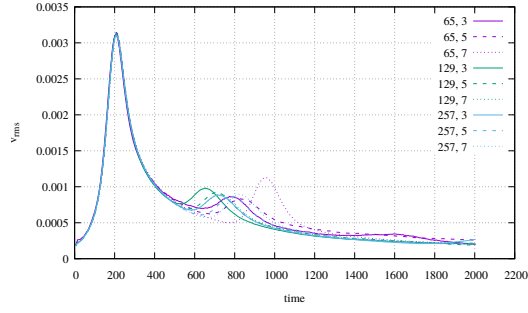




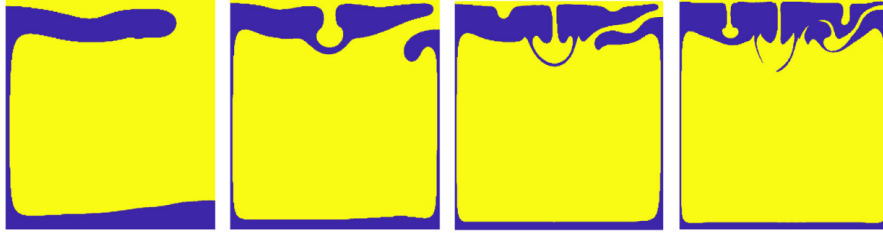
- Schuh-Senlis *et al.* (2020) [1144]. The setup in the paper is inspired by [1309] but results are therefore not consistent with those of [1309].



- MVEP2 code, courtesy of Marcel Thielmann



- Burcet, Oliveira, Afonso, and Zlotnik [175] (2024).

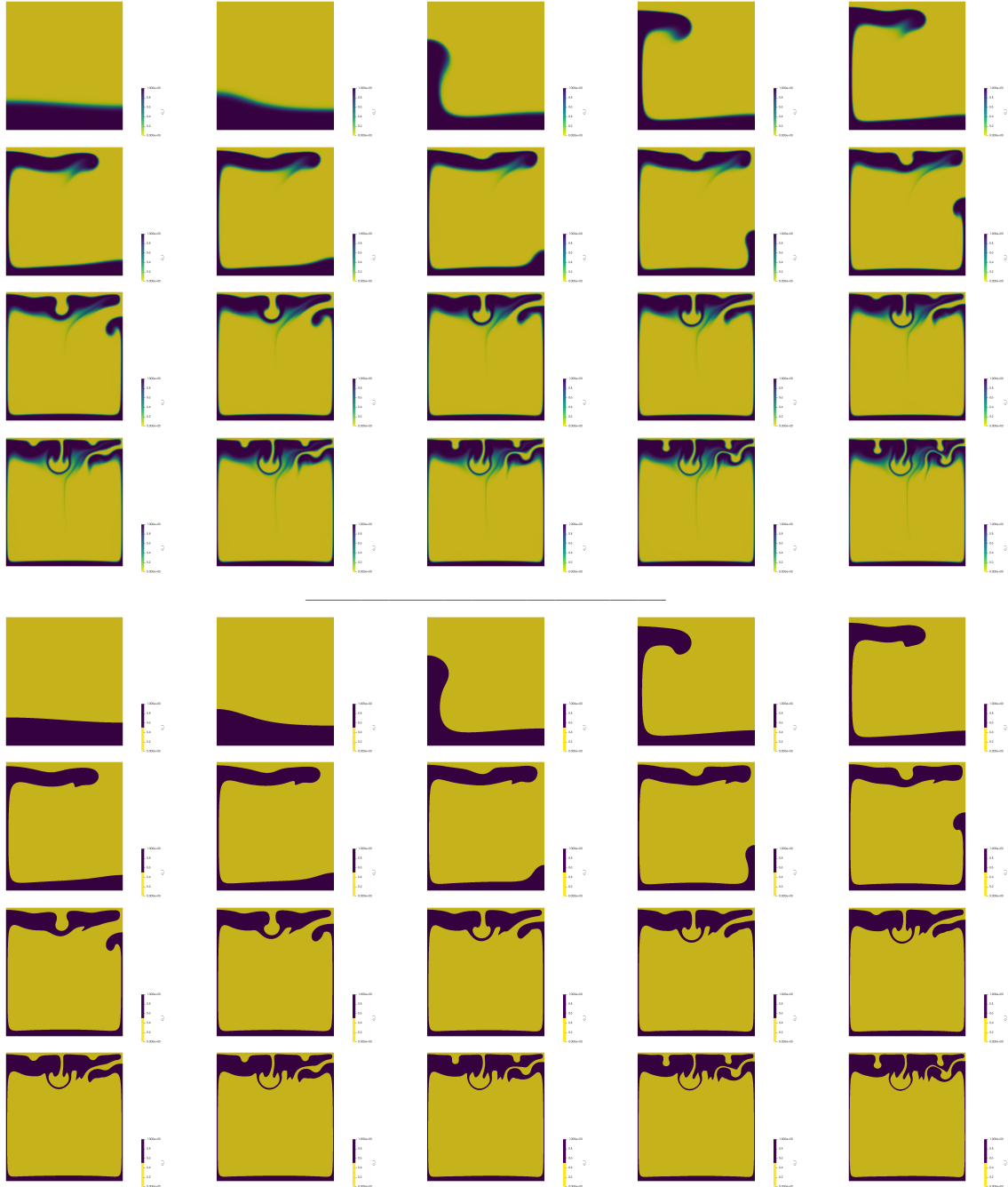


**Table 2**  
Comparison of measurable results in the Rayleigh-Taylor benchmark 4.3.

| Method              | $\varphi$ | $\mu_{\text{rms}}^{\text{max}}$ | $t(\mu_{\text{rms}}^{\text{max}})$ | Ref. |
|---------------------|-----------|---------------------------------|------------------------------------|------|
| FCFV (this study)   | 0.011093  | 0.003016                        | 209.35                             | -    |
| ALE                 | 0.01222   | 0.003075                        | 211                                | [58] |
| Level sets          | 0.01252   | 0.00301                         | 211.2                              | [59] |
| Tracers (HS)        | 0.011177  | 0.003092                        | 208.99                             | [56] |
| Tracers (CND)       | 0.01106   | 0.003094                        | 208.5                              | [56] |
| Field (SK)          | 0.01179   | 0.002897                        | 207.84                             | [56] |
| Marker Chains (PvK) | 0.01225   | 0.003091                        | 207.05                             | [56] |

- Logg *et al.* (2012) [806]: only Entrainment of a Dense Layer by Thermal convection.

- ASPECT manual [44].



Time evolution of the system, smooth compositional field, CFL=0.25, mesh refinement level 8. Top block is with 256 colors, bottom block is the same data but only using two colors. From  $t=0$  until  $t=2000$ , every 100s.

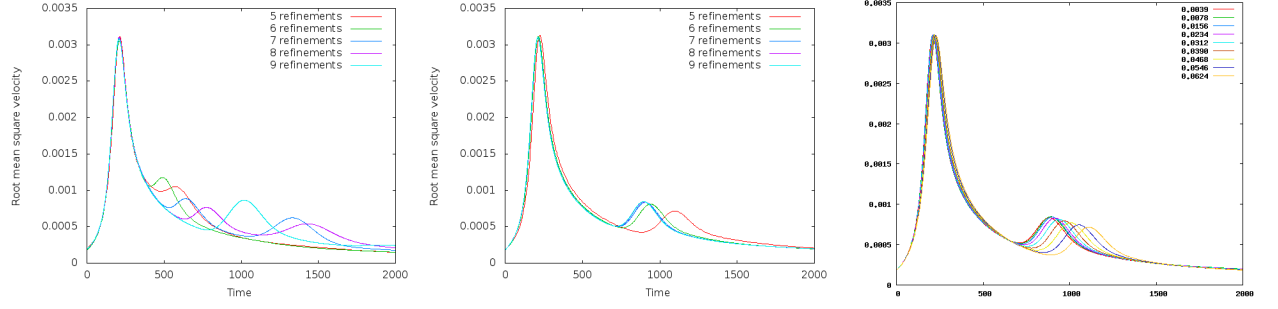
When compositional fields are used the solution has been shown (see ASPECT manual) to be very sensitive to the resolution. In an attempt to remedy this issue, an approach was devised, which consists in replacing the original discontinuous initial condition with a smoothed out version. The function in the input file

```
set Function expression = if((z>0.2+0.02*cos(pi*x/0.9142)) , 0 , 1)
```

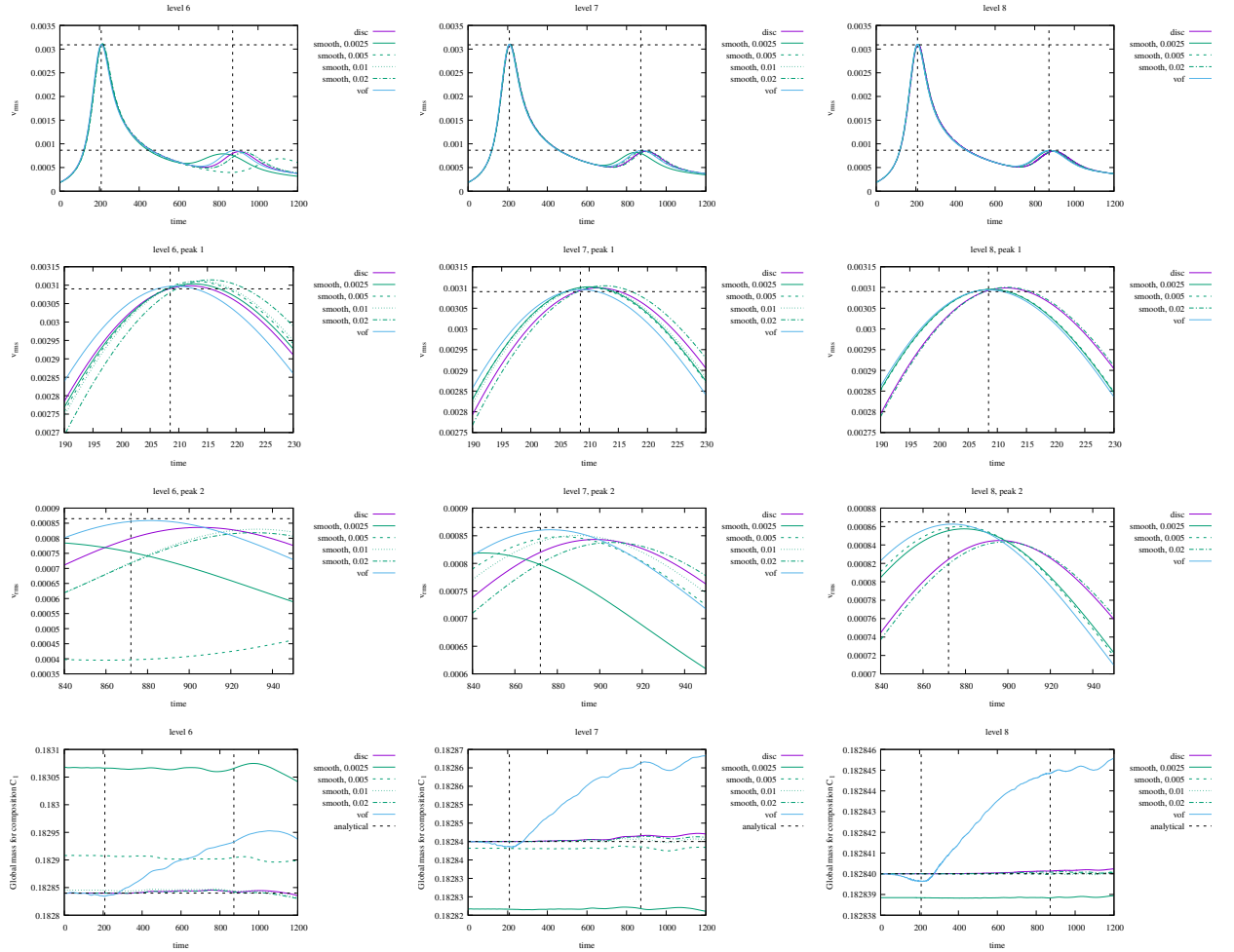
is then replaced by

```
set Function expression = 0.5*(1+tanh((0.2+0.02*cos(pi*x/0.9142)-z)/0.02))
```

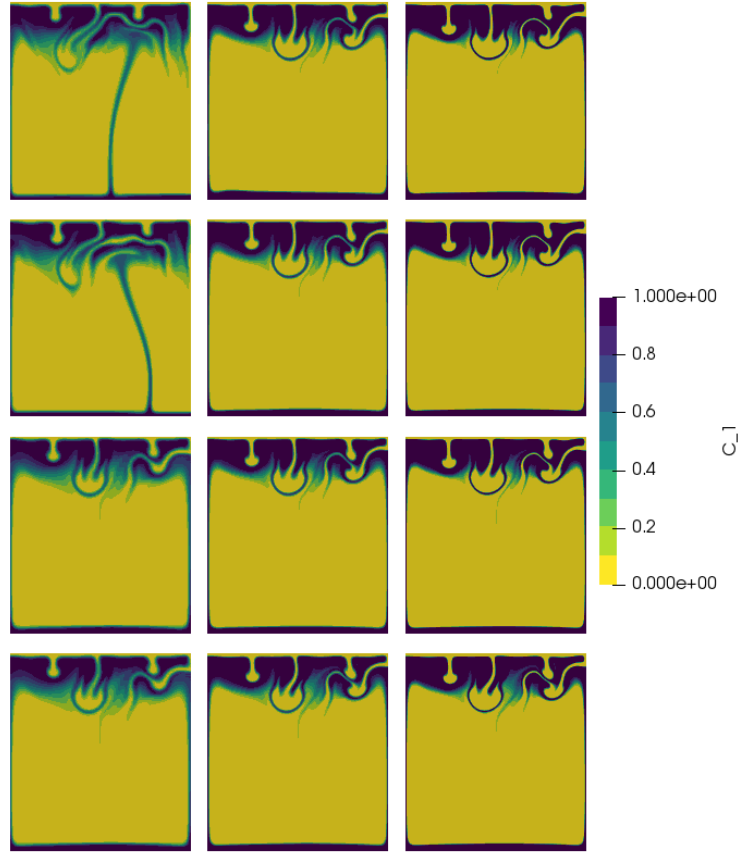
The last number on the line is a parameter which controls the 'thickness' of the interface. When it becomes small we recover a discontinuous implementation (this very much depends on the mesh size).



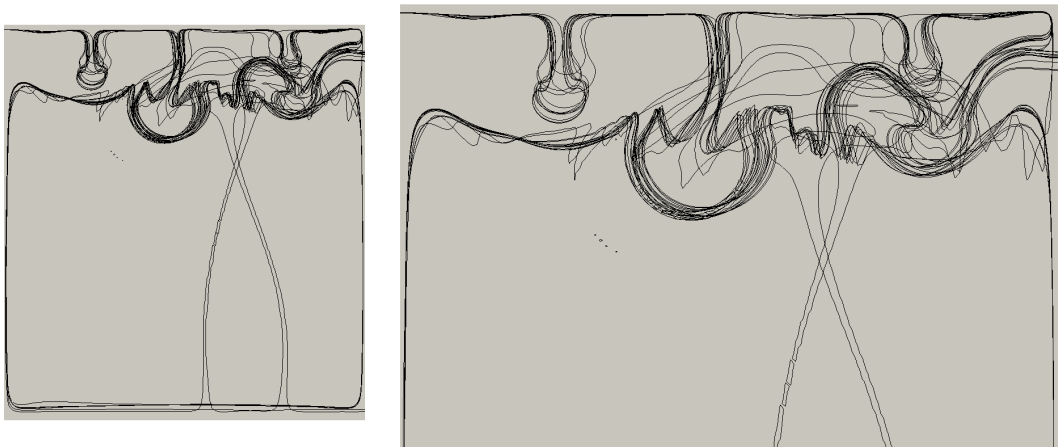
Taken from the ASPECT manual. Root mean square measurements with discontinuous (left) and smoothed, continuous (middle) initial conditions for the compositional field: 5 global refinements correspond to a  $32 \times 32$  mesh, 9 refinements to a  $512 \times 512$  mesh. Right: influence of smoothing parameter (level 7).



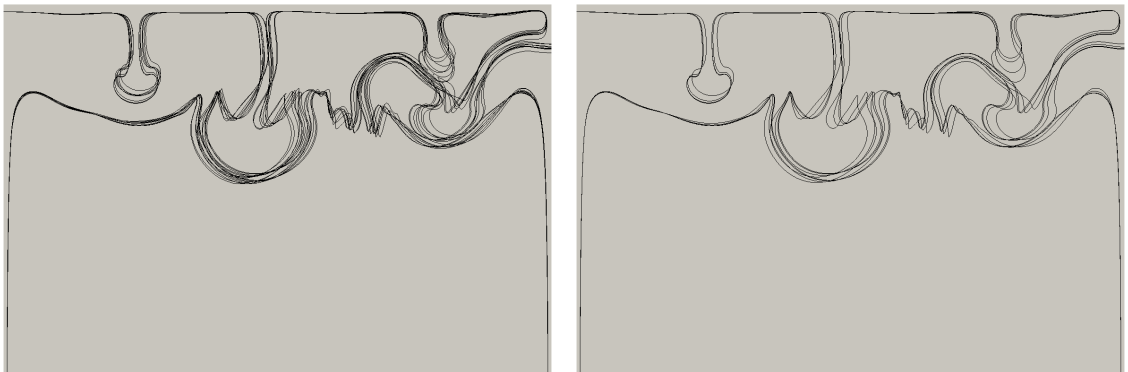
Results obtained for both approaches and vof, on various meshes and for various smoothing parameters. The dashed lines indicate the [STONE 95](#) results.



Results for smoothed approach. Left to right: level 6,7,8. Top to bottom: smoothing parameter 0.0025, 0.005, 0.01, 0.02. We see that level6 results in combination with a small smoothing parameter yield a very different outcome.

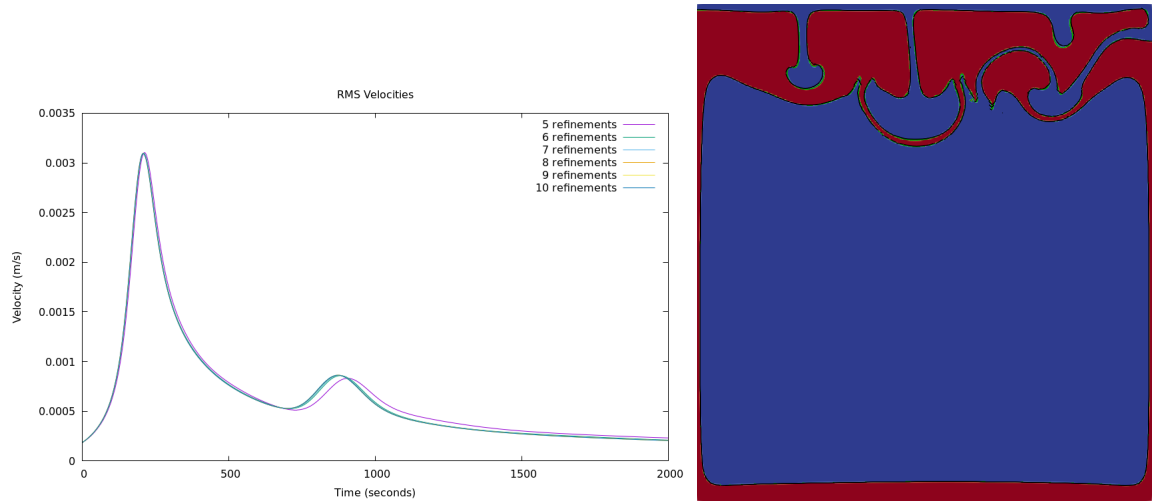


0.5 compositional field  $C_1$  isocontours for all 12 simulations.



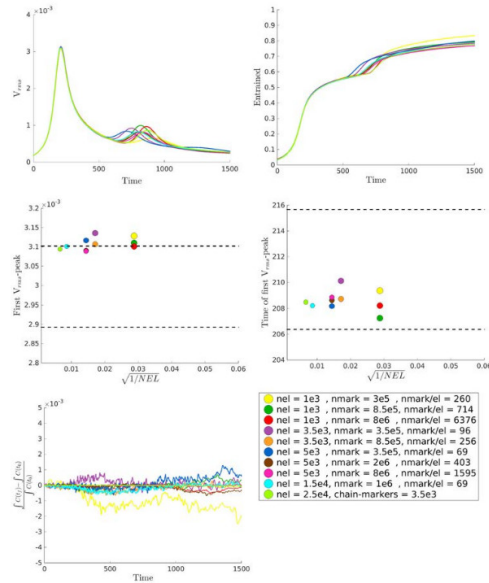
0.5 compositional field  $C_1$  isocontours for levels 7 and 8 (left), and level 8 (right)

The code can also rely on the VOF method to solve the advection equation for the computational fields:

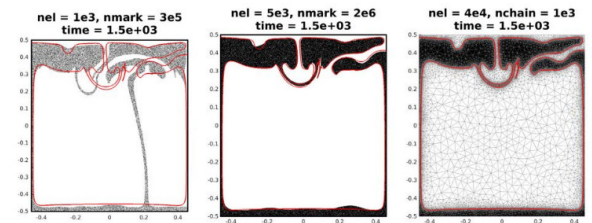


Left: Evolution of the root mean square velocity as a function of time for computations of the van Keken problem made with the VOF interface tracking algorithm with five different global mesh refinements (from a  $32 \times 32$  mesh to a  $1024 \times 1024$  mesh. Right: The results of two computations of the van Keken problem made with the VOF interface tracking algorithm overlaid upon each other at  $t_{\text{end}} = 2000$ . This visualization shows the reconstructed boundary between the two materials at the final time  $t_{\text{end}}$  as computed on a uniform grid with 7 and 8 levels of refinement. The boundaries between the materials are displayed as contours of the fields  $\tilde{\psi}^7(t_{\text{end}})$  (black) and  $\tilde{\psi}^8(t_{\text{end}})$  (bright green), which are generated by the visualization postprocessor. The contours for the reconstructed material boundaries are superimposed on a color gradient visualization of the material composition for the computation with 8 levels of refinement in order to make the regions with each fluid type more evident.

- Mulyokova phd thesis [919]



**Figure 4.19:** Benchmark-quantities for the isoviscous Rayleigh-Taylor instability set-up, obtained with the tracer-ratio method at different resolutions (according to legend) and marker-chain method with a conforming grid. In the legend, ‘nel’ means number of mechanical elements, ‘nmark’ number of markers in the tracer-ratio method, and ‘nchain’ number of markers in the marker-chain method. Black dashed lines in plots showing values and points in time of the first  $v_{\text{rms}}$ -peak (middle row) indicate the range of values for the respective quantities presented in Van Keken et al. (1997).



**Figure 4.18:** Comparison of the results for the Rayleigh-Taylor benchmark obtained with the tracer-ratio methods at two different resolutions (first two figures), and those obtained with the marker-chain method with a conforming grid, at time  $t = 1500$ . For the latter, the mesh is shown with gray lines delineating element-edges. The red line delineates the marker-chain, and is plotted for comparison in all three figures. Resolutions are given in the titles of each figure, where ‘nel’ means number of mechanical elements, ‘nmark’ number of markers in the tracer-ratio method, and ‘nchain’ number of markers in the marker-chain method.



| Publication                        | $\eta^*$ | Stokes method | Transport method | vrms | png | max time | peak 1 time | peak 1 vrms | peak 2 time | peak 2 vrms | growth rate | observation              |
|------------------------------------|----------|---------------|------------------|------|-----|----------|-------------|-------------|-------------|-------------|-------------|--------------------------|
| van Keken <i>et al.</i> [1309]     | 1        |               |                  | ✓    | ✓   | 2000     | 211.1       | 0.0030958   |             |             | 0.010996    | HS 41x41                 |
|                                    | 1        |               |                  |      |     |          | 209.17      | 0.0031022   |             |             | 0.011109    | HS 61x61                 |
|                                    | 1        |               |                  |      |     |          | 208.99      | 0.0030916   |             |             | 0.011177    | HS 81x81                 |
|                                    | 1        |               |                  |      |     |          | 208.4       | 0.003092    |             |             | 0.01106     | CND 32x32                |
|                                    | 1        |               |                  |      |     |          | 208.5       | 0.0030943   |             |             | 0.01106     | CND 48x48                |
|                                    | 1        |               |                  |      |     |          | 215.67      | 0.00299279  |             |             | 0.01130     | SK 80x80                 |
|                                    | 1        |               |                  |      |     |          | 206.38      | 0.0028922   |             |             | 0.01127     | SK 120x120               |
|                                    | 1        |               |                  |      |     |          | 207.84      | 0.0028970   |             |             | 0.01179     | SK 160x160               |
|                                    | 1        |               |                  |      |     |          |             |             |             |             | 0.01220     | SK 160x160               |
|                                    | 1        |               |                  |      |     |          | 213.38      | 0.00300     |             |             | 0.01185     | PvK 30x30 splines        |
|                                    | 1        |               |                  |      |     |          | 211.81      | 0.003016    |             |             | 0.01198     | PvK 50x50 splines        |
|                                    | 1        |               |                  |      |     |          | 210.75      | 0.003050    |             |             | 0.01207     | PvK 80x80 splines        |
|                                    | 1        |               |                  |      |     |          |             |             |             |             | 0.01211     | PvK 100x100 splines      |
|                                    | 1        |               |                  |      |     |          | 210.59      | 0.003100    |             |             | 0.01253     | PvK 30x30 C1             |
|                                    | 1        |               |                  |      |     |          | 207.05      | 0.003091    |             |             | 0.01225     | PvK 80x80 C1             |
| de Smet <i>et al.</i> [323]        |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Sobouti <i>et al.</i> [1176]       |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Babeyko <i>et al.</i> [35]         |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Tackley & King [1229]              |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Bourgouin <i>et al.</i> [124]      |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Quinteros <i>et al.</i> [1029]     |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Samuel & Evonuk [1103]             |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Suckale <i>et al.</i> [1218]       | 1        | FDM           | level set        | ✓    | ✓   | 1500     | 211.2       | 0.00301     | ?           | ?           | 0.01252     | 120x132                  |
| Suckale <i>et al.</i> [1218]       | 10       | FDM           | level set        | ✓    | ✓   | 500      | ?           | ?           | ?           | ?           | 0.04809     | 250x275                  |
| Leng & Zhong [769]                 | 1        | FEM           | tracers          | ✓    | ✓   | 1500     | ?           | ?           | ?           | ?           | ?           | AMR(6+3)                 |
| Maierova PhD thesis [825]          | 1        |               |                  | Y    | Y   | 2000     | 211         | 0.003107    |             |             |             |                          |
|                                    | 0.1      |               |                  | Y    | N   | 2000     | 73          | 0.009411    |             |             |             |                          |
|                                    | 100      |               |                  | Y    | N   | 2000     | 51          | 0.013938    |             |             |             |                          |
| Vynnytska <i>et al.</i> [1333]     |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Choi <i>et al.</i> [238]           |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Fuchs & Schmeling [421]            |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Mulyukova PhD thesis[919]          |          |               |                  |      |     |          |             |             |             |             |             |                          |
| de Montserrat <i>et al.</i> [322]  |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Robey & Puckett [1079]             |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Robey PhD thesis [1078]            |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Trim <i>et al.</i> (2021) [1282]   |          |               |                  |      |     |          |             |             |             |             |             |                          |
| Louis-Napoleon <i>et al.</i> [811] | 1        | FEM           | VOF & CF         |      |     |          |             |             |             |             |             | JADIM, OpenFOAM & ASPECT |
| ASPECT manual [44]                 | 1        | FEM           | CF, PIC          | ✓    | ✓   |          |             |             |             |             |             |                          |

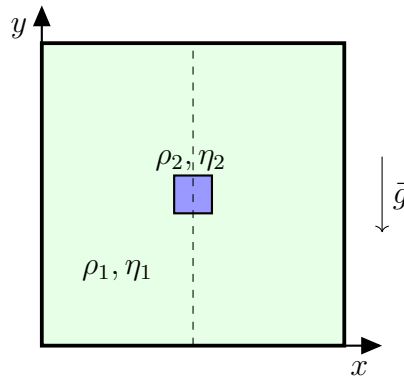


## 12.2.20 (Instantaneous) Sinking block (2D)

benchmark\_sinking\_block.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/sinking\\_block](https://github.com/cedrict/fieldstone/tree/master/images/sinking_block)

The domain is a unit square. Fluids are such that  $\rho_1 = 1$ ,  $\eta_1 = 1$  and  $\rho_2 = 1.01$ ,  $\eta_2 = 1000$ . Boundary conditions are either free slip or no slip on all sides. Pressure is normalised so that the volume average is zero. Gravity points downwards with  $|\vec{g}| = 1$ . Profile measurements are carried out on the dashed line.



When using ASPECT, it is good to remember that a compositional field is used, which 'lives' on the nodes of the FE grid. Part of the input file is shown here:

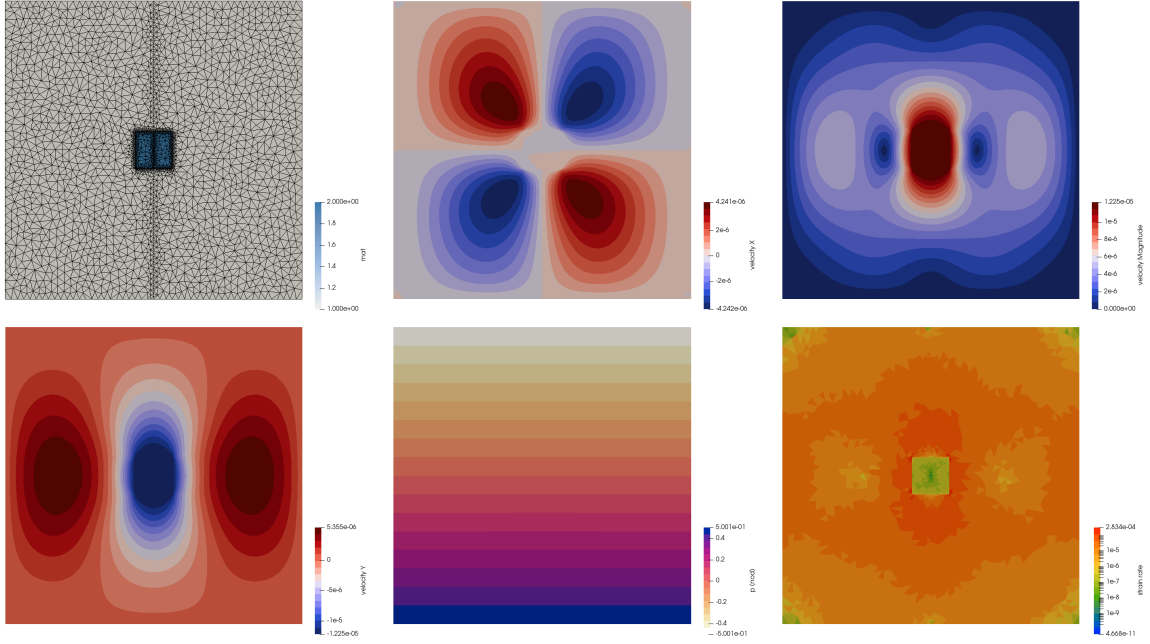
```
subsection Compositional fields
 set Number of fields = 1
end

subsection Initial composition model
 set Model name = function
 subsection Function
 set Variable names = x,y
 set Function constants = p=0.5
 set Function expression = if(abs(x-p)<0.0625 && abs(y-p)<0.0625 , 1, 0)
 end
end

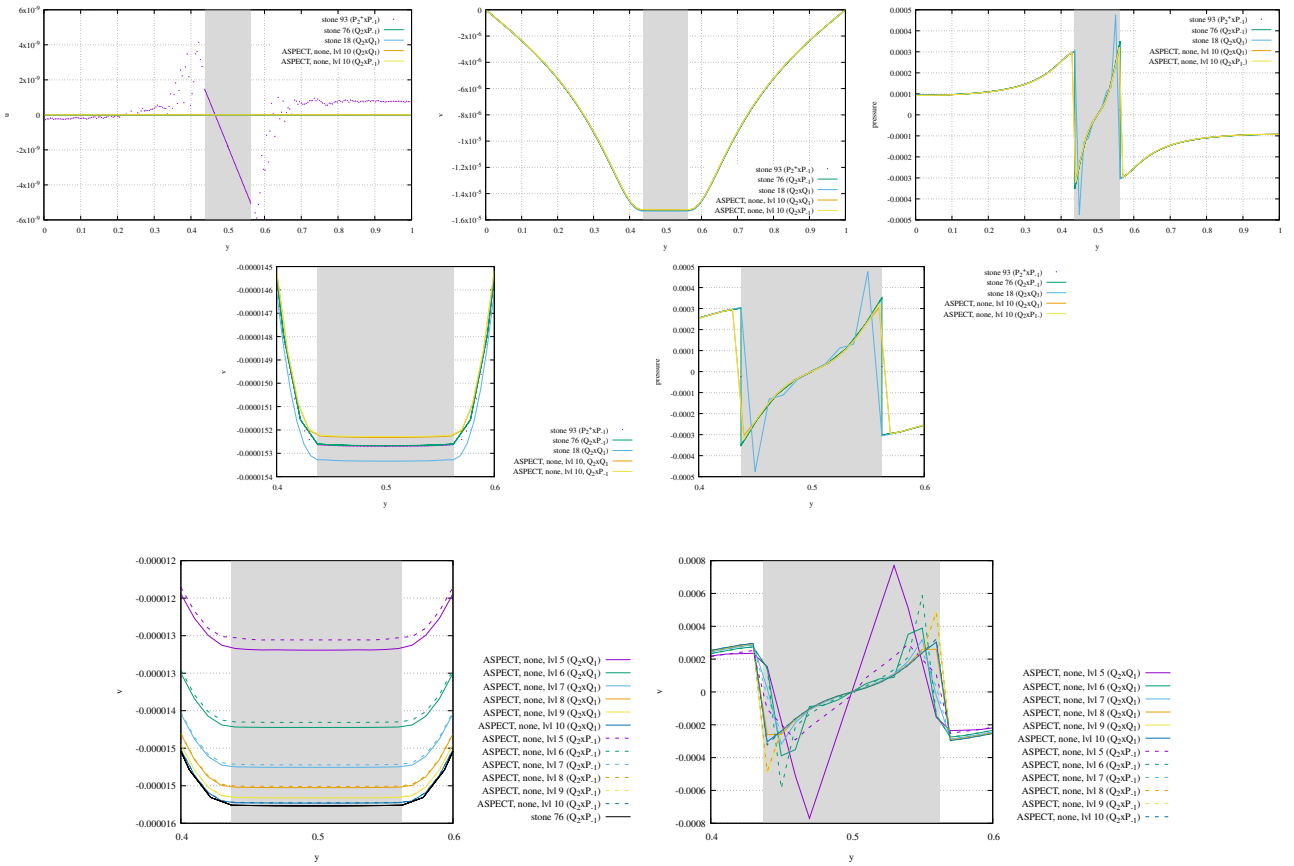
subsection Material model
 subsection Simple model
 set Density differential for compositional field 1 = 0.01
 set Composition viscosity prefactor = 1000
 end
end
```

The value of the composition (and therefore the density and viscosity values) on a quadrature point is obtained via interpolation and averaging, which is different than the Stone codes where the density and viscosity are elemental quantities.

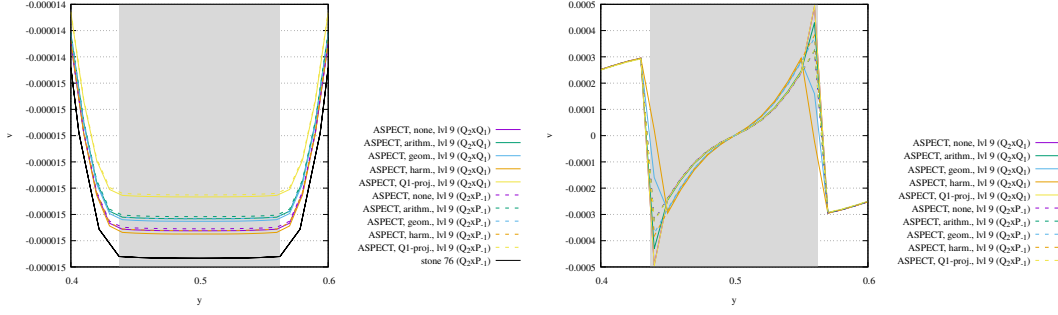
## Free-slip boundary conditions



Results obtained with Stone 93.



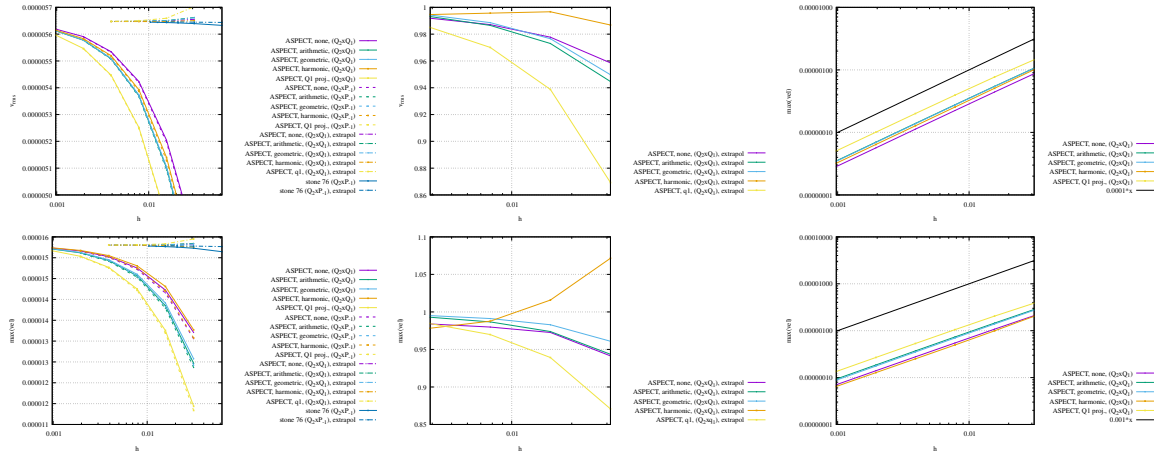
ASPECT results with various global mesh refinement. No averaging



ASPECT results with various averagings. level 9.

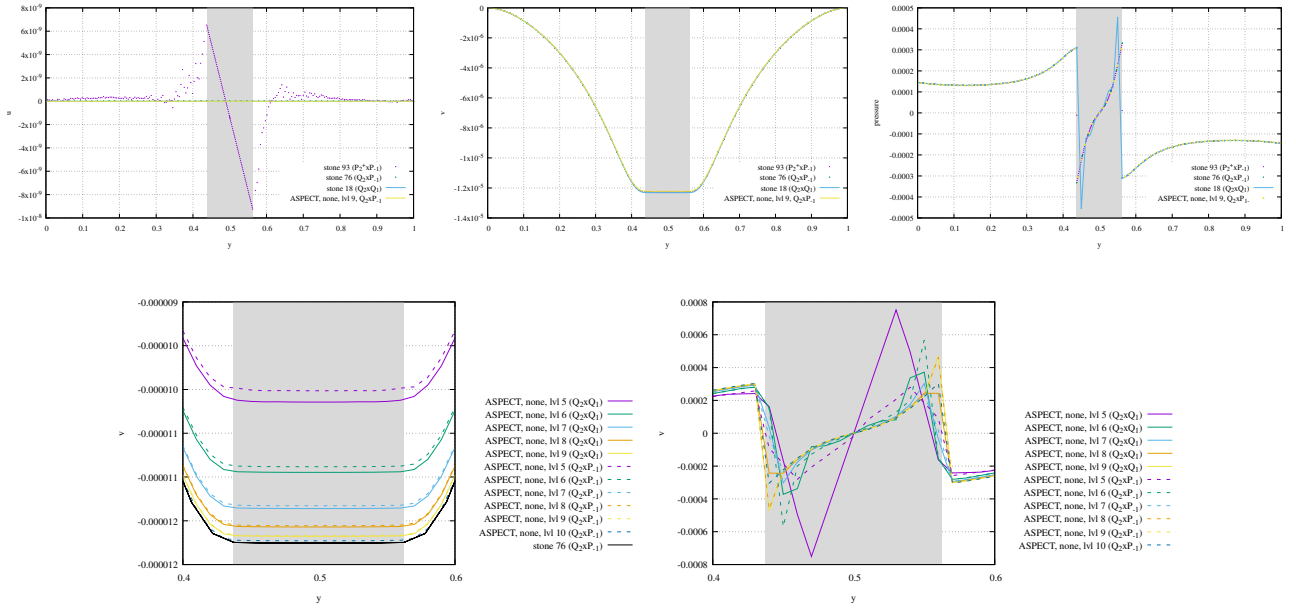
|                 | FS                                                                                                                | NS |
|-----------------|-------------------------------------------------------------------------------------------------------------------|----|
| $\min(u)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\max(u)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\min(v)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\max(v)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\max( v )$     | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $v_{rms}$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\min(p)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $\max(p)$       | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |
| $ v (0.5, 0.5)$ | Stone 18 ( $Q_2 \times Q_1$ )<br>Stone 93 ( $P_2^+ \times P_{-1}$ )<br>Stone 76 ( $Q_2 \times P_{-1}$ )<br>Aspect |    |

Using error extrapolation (see Section 9.12), one can compute an estimate of the resolution independent value of the vrms of maximum velocity for example:

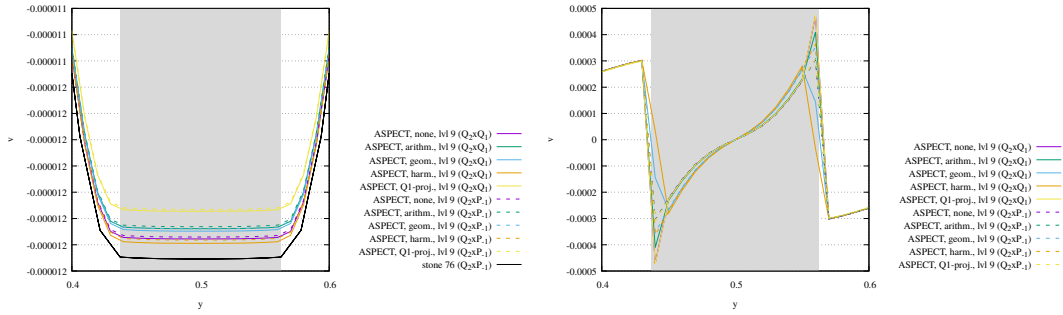


We find that the rates are near unity.  
 TODO: write material model in ASPECT to bypass compositions!

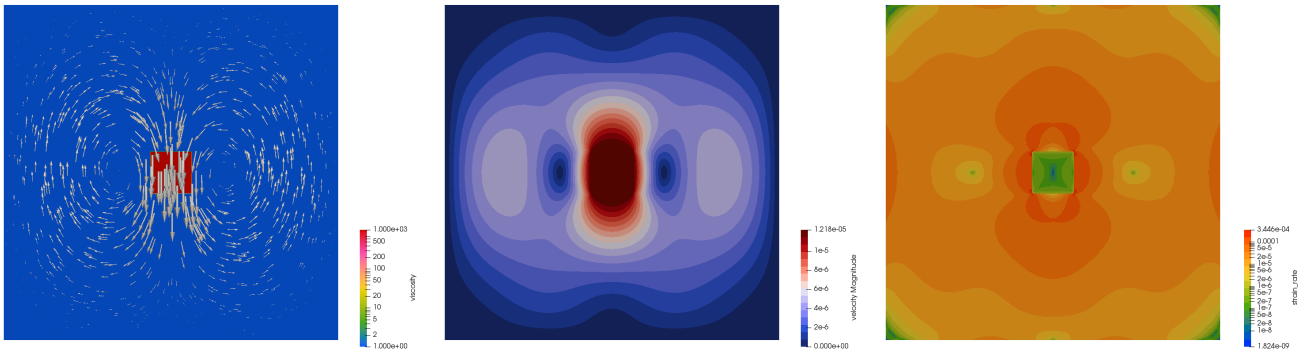
## No-slip boundary conditions



ASPECT results with various global mesh refinement. No averaging



ASPECT results with various averagings. level 9.



Obtained with ASPECT, level 9.

TODO: finish analysis

## 12.2.21 (Instantaneous) Stokes sphere (3D)

benchmark\_stokes\_sphere\_3D.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/stokes\\_sphere3D](https://github.com/cedrict/fieldstone/tree/master/images/stokes_sphere3D)

This is a simple experiment without an analytical solution. The idea here is simple: to design a small number of Stokes sphere-related experiments and provide for them (very) high-resolution results obtained with various codes so as to turn these into benchmarks. The domain is chosen to be a unit cube. Gravity is such that  $\vec{g} = (0, 0, -g)$  with  $g = 1$ . The sphere is in the middle of the domain and has a radius  $R = 0.123456798$ . The fluid has a density  $\rho_f = 1$  and viscosity  $\eta_f = 1$ . The sphere has a density  $\rho_s = \rho_f + \delta\rho$  and a viscosity  $\eta_s = 10^m \eta_f$ . Default values for  $\delta\rho$  and  $m$  are set to 0.01 and 3 respectively.

Concerning boundary conditions, we distinguish three cases:

- FS: free slip boundary conditions are imposed on all 6 sides;
- NS: no slip boundary conditions are imposed on all 6 sides;
- OT: free slip boundary conditions are prescribed on the sides and bottom, but the top surface is open.
- BO: both top and bottom are open (still with  $u = 0$ ) and no-slip is prescribed on the sides.

In the FS and NS case the null space of the pressure will need to be addressed and we require that the average pressure over the domain is zero, i.e.

$$\iiint_{\Omega} p(x, y, z) dx dy dz = 0$$

The following quantities are reported:

- the root mean square velocity  $v_{rms}$  over the whole domain;
- the minimum and maximum velocity and pressure in the domain (i.e.  $u_{min,max}$ ,  $v_{min,max}$ ,  $w_{min,max}$  and  $p_{min,max}$ );
- the velocity in the center of the sphere (maybe).

The factors which are expected to influence these measurements are:

- the resolution, especially if hexahedral elements are used;
- the quadrature rule, especially if the material properties are directly prescribed on these;
- the type of numerical method and their order (think  $Q_1 \times P_0$  vs  $Q_2 \times Q_1$  vs  $Q_2 \times P_{-1}$  vs ... for finite elements)
- whether full or reduced densities are used (except for OT case);
- the parameter  $m$  which controls the rigidity of the sphere with regards to the surrounding fluid;

- the relative density difference between the fluid and the sphere.

Stokes' law was derived by George Gabriel Stokes in 1851. It describes the frictional force a sphere with a density different than the surrounding fluid experiences in a laminar flowing viscous medium around it. By equating the frictional force term  $6\pi\eta_f R\mathbf{v}_s$  with the buoyancy force  $4/3\pi R^3\delta\rho g$ , we arrive at the following settling velocity:

$$\mathbf{v}_s = \frac{2}{9} \frac{\delta\rho R^2 g}{\eta_f}$$

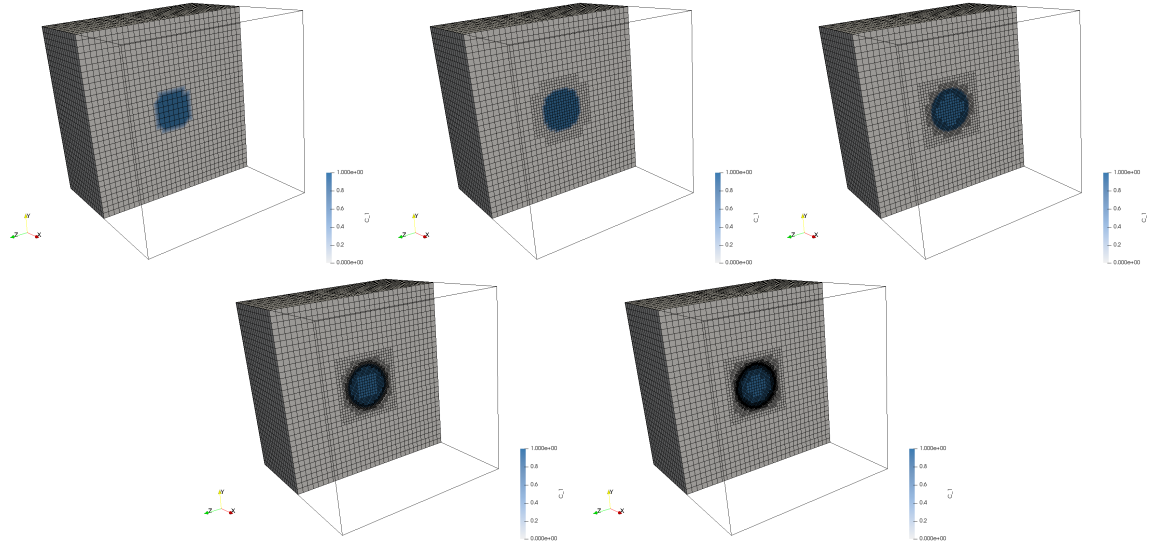
All the measurements above will then be adimensionalised by dividing all velocities by  $\mathbf{v}_s$  and pressures by  $p_{ref} = \rho_f g L_z = 1$ . Note that Stokes law is derived in an infinite fluid so that the recovered sphere velocity measurements are not expected to match this analytical value exactly. In our case we have

$$\mathbf{v}_s = \frac{2}{9} \frac{0.01 \cdot 0.123456789^2 \cdot 1}{1} \simeq 0.00003387017 = 3.387017 \cdot 10^{-5}$$

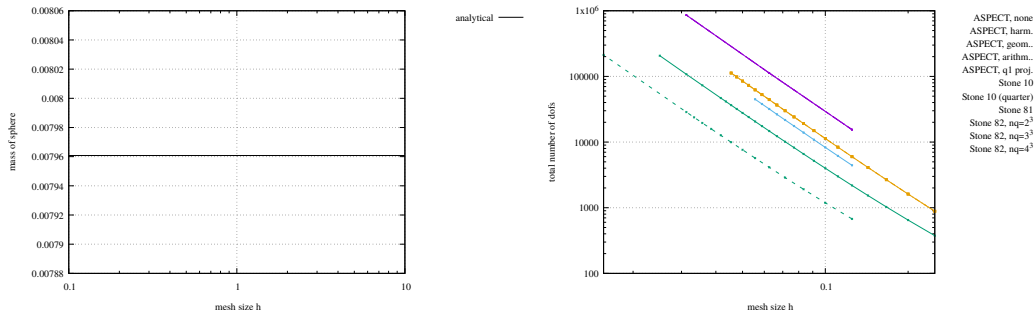
Note that as noted in the Gale manual<sup>10</sup> a correction can be made to this velocity when the sphere is itself viscous, see problem 2 p65-66 of second English edition of 'Fluid Mechanics' by Landau & Lifshitz, volume 6 of Theoretical Physics.

---

<sup>10</sup><https://geodynamics.org/cig/software/gale/>

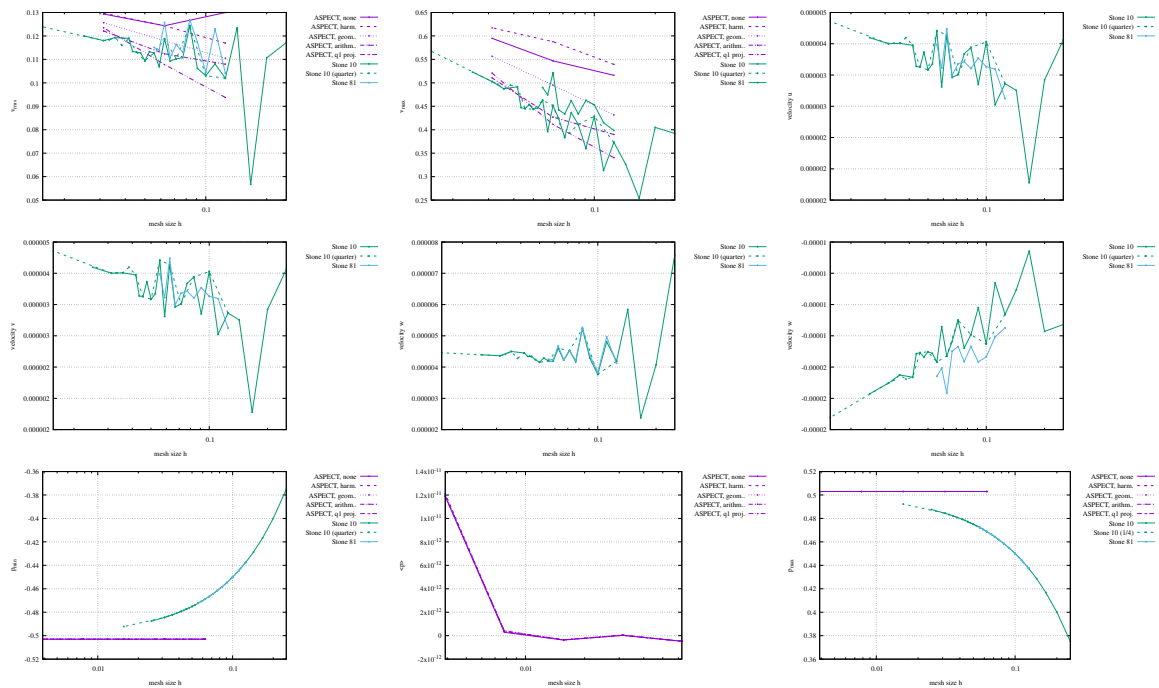


Octree-based ASPECT meshes. The background mesh is  $16^3$  and refinement is allowed to take place 4 times. Note that a special output is automatically generated in the code which subdivides all elements in 8 for visualisation purposes only.

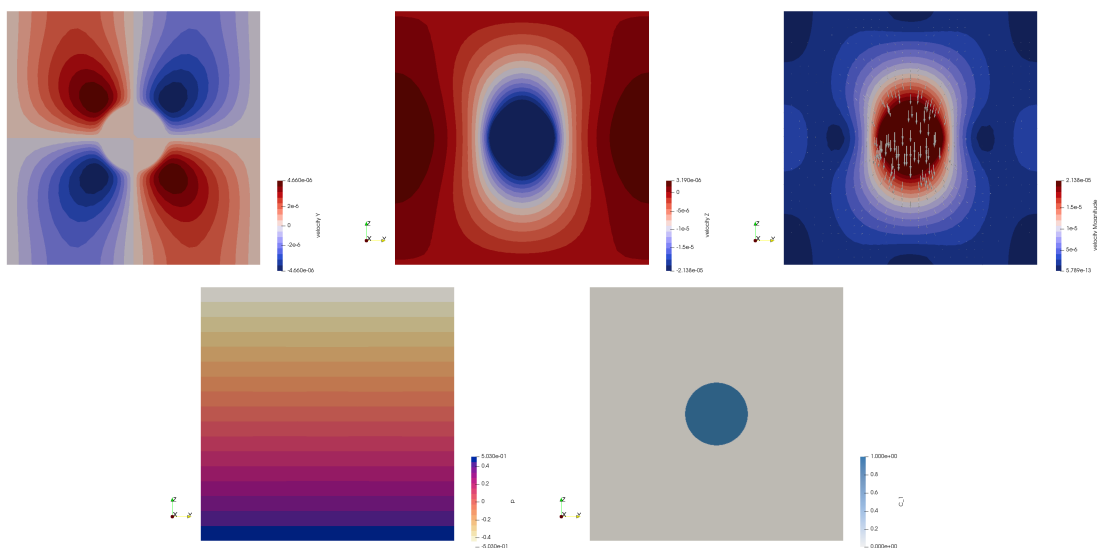




## FS results

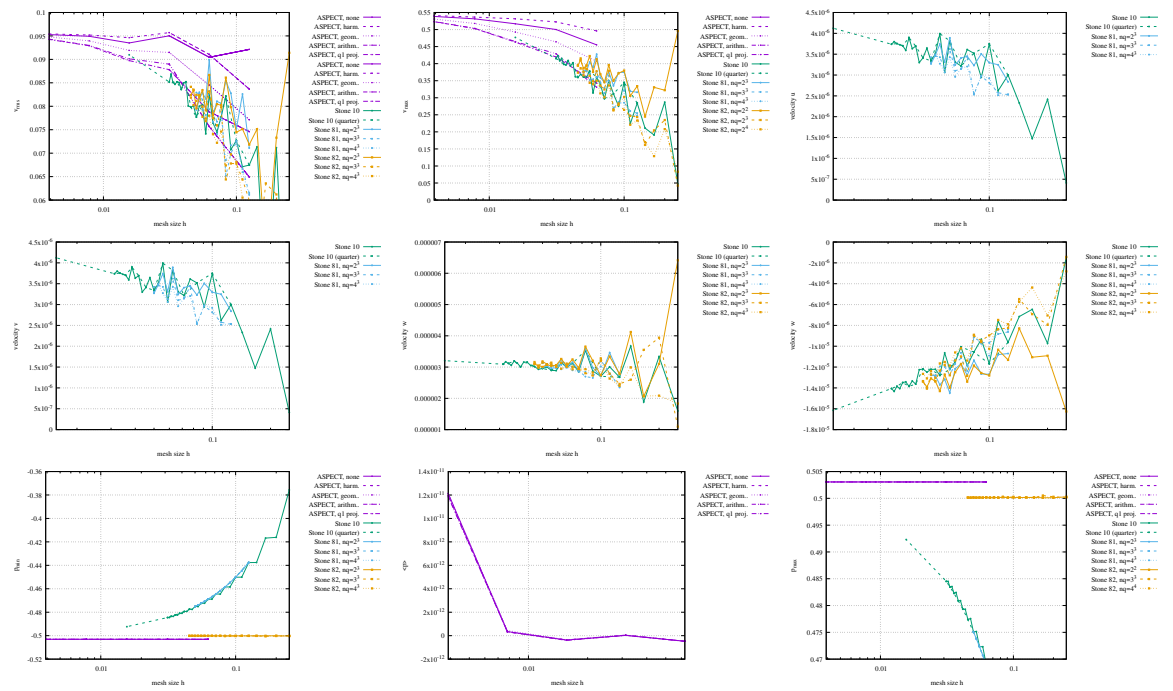


Measurements obtained with ASPECT and STONE 10 for various averaging schemes.

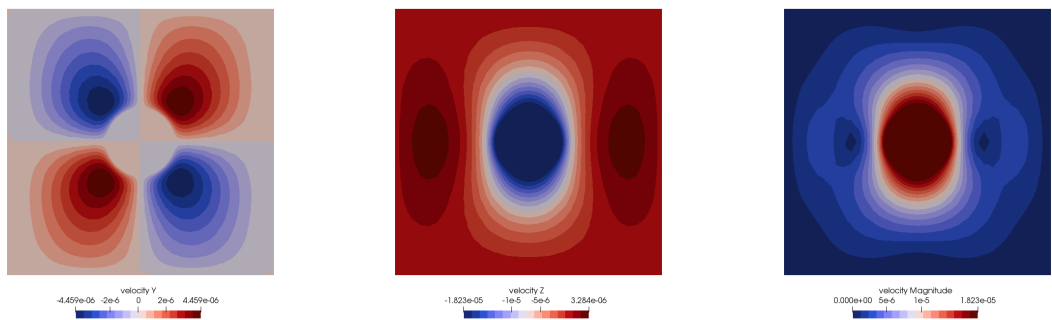


High-resolution solution cross section at  $x = 0.5$

NS results

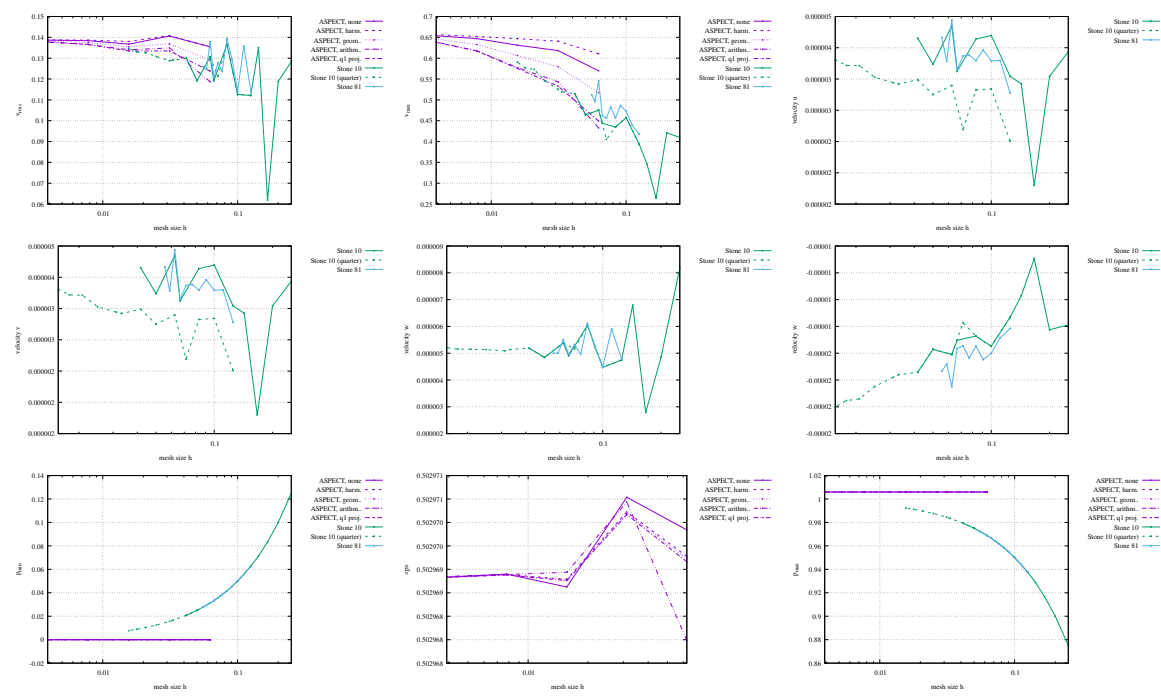


Measurements obtained with ASPECT and STONE 10 for various averaging schemes.

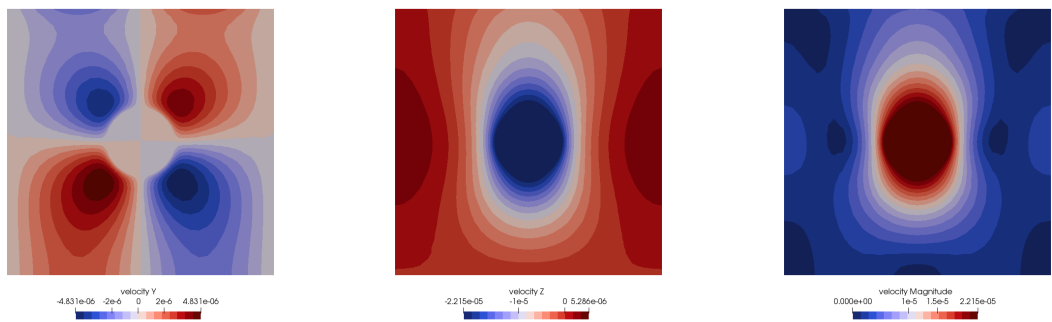


4+4 mesh

OT results



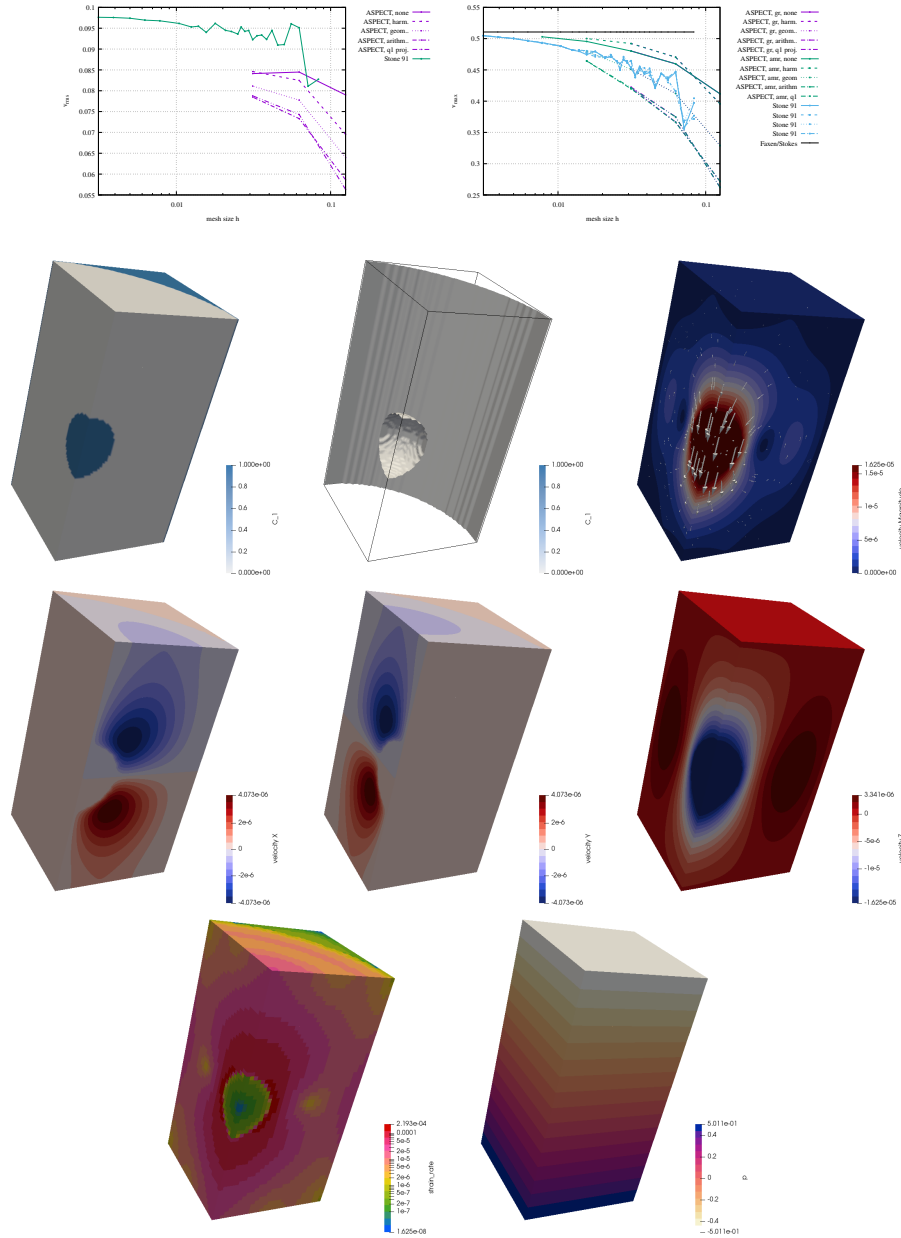
Measurements obtained with ASPECT and STONE 10 for various averaging schemes.



4+4 mesh

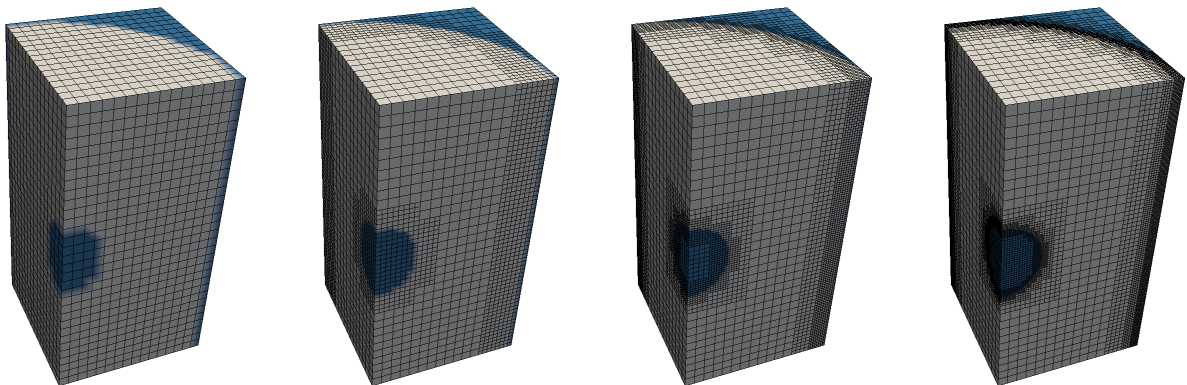
## CYL results .

note that sphere and walls are only 1000 times more viscous than fluid



level 5 mesh, 32x32x64 elements.

I have proven in [STONE 92](#) that the sphere should probably be  $10^6$  times more viscous than the fluid and the box should be 1.5 in height to recover the Habermann/Faxen velocities.



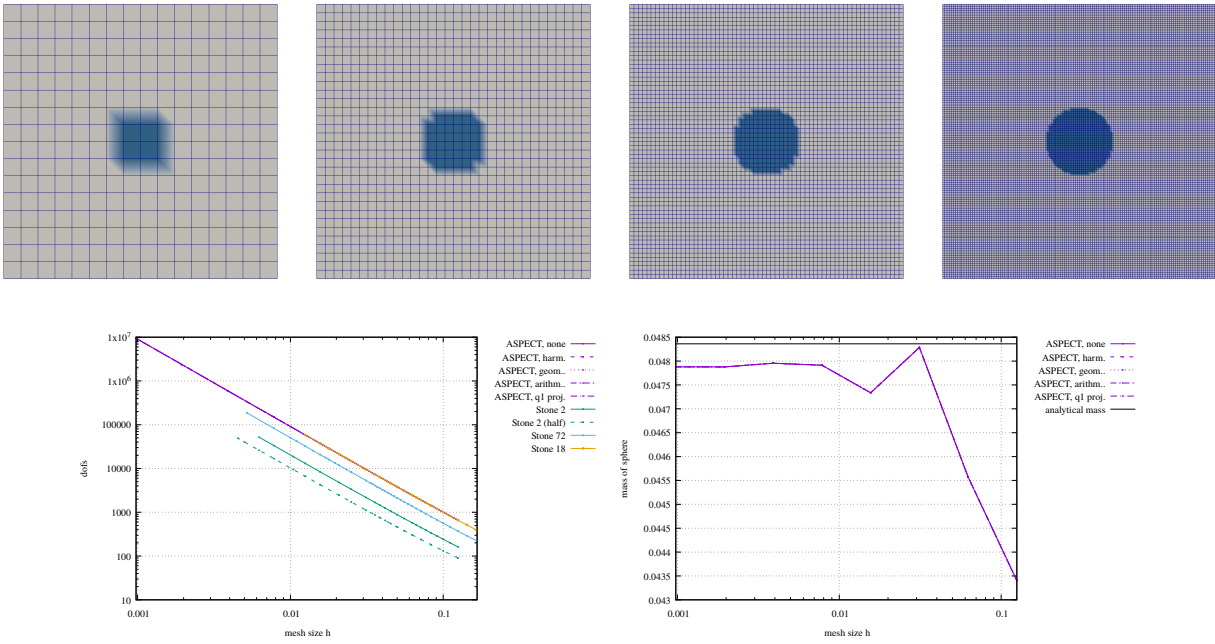
## 12.2.22 (Instantaneous) Stokes sphere (2D)

benchmark\_stokes\_sphere\_2D.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/stokes\\_sphere2D](https://github.com/cedrict/fieldstone/tree/master/images/stokes_sphere2D)

This is the same experiment as in the 3D case but in 2D. When using ASPECT, we simply start with regular meshes ranging from  $8^2$  to  $512^2$  elements and we use the default  $Q_2 \times Q_1$  element. This corresponds to

```
subsection Mesh refinement
set Initial adaptive refinement = 0
set Initial global refinement = 3->9
set Refinement fraction = 0.9
set Coarsening fraction = 0
set Strategy = composition
end
```



Left: total number of dofs in the Stokes problem; Right: mass of composition 1 as measured in ASPECT.

The total mass of the system is

$$M = (L_x L_y - \pi R^2) \rho_{fluid} + \pi R^2 \rho_{sphere} \quad (12.274)$$

$$= (L_x L_y - \pi R^2) \rho_{fluid} + \pi R^2 (\rho_{fluid} + \delta \rho) \quad (12.275)$$

$$= L_x L_y \rho_{fluid} + \pi R^2 \delta \rho \quad (12.276)$$

$$= 1 + \pi \cdot 0.123456789^2 \cdot 0.01 \quad (12.277)$$

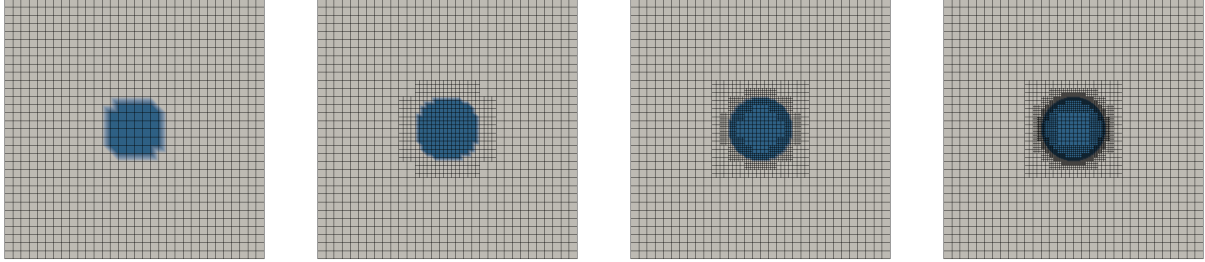
$$\simeq 1.00047882831 \quad (12.278)$$

The Stokes velocity can be obtained as follows: on p61 of Landau & Lifschitz, it is reported that the drag force on a disk moving in its plane is  $F = 32\eta_f R \mathbf{v}_s / 3$ . The buoyancy force is  $F = \pi R^2 \delta \rho g$ , so the velocity is then

$$\mathbf{v}_s = \frac{3\pi}{32} \frac{\delta \rho}{\eta_f} R g \simeq 0.00036361025$$

Given the dimensions, this is obviously given per meter of the infinite cylinder. This is substantially smaller than what we recover, so I keep the 3D velocity as reference for now.

In a second time, we make use of the mesh refinement capabilities of the code, as shown here:

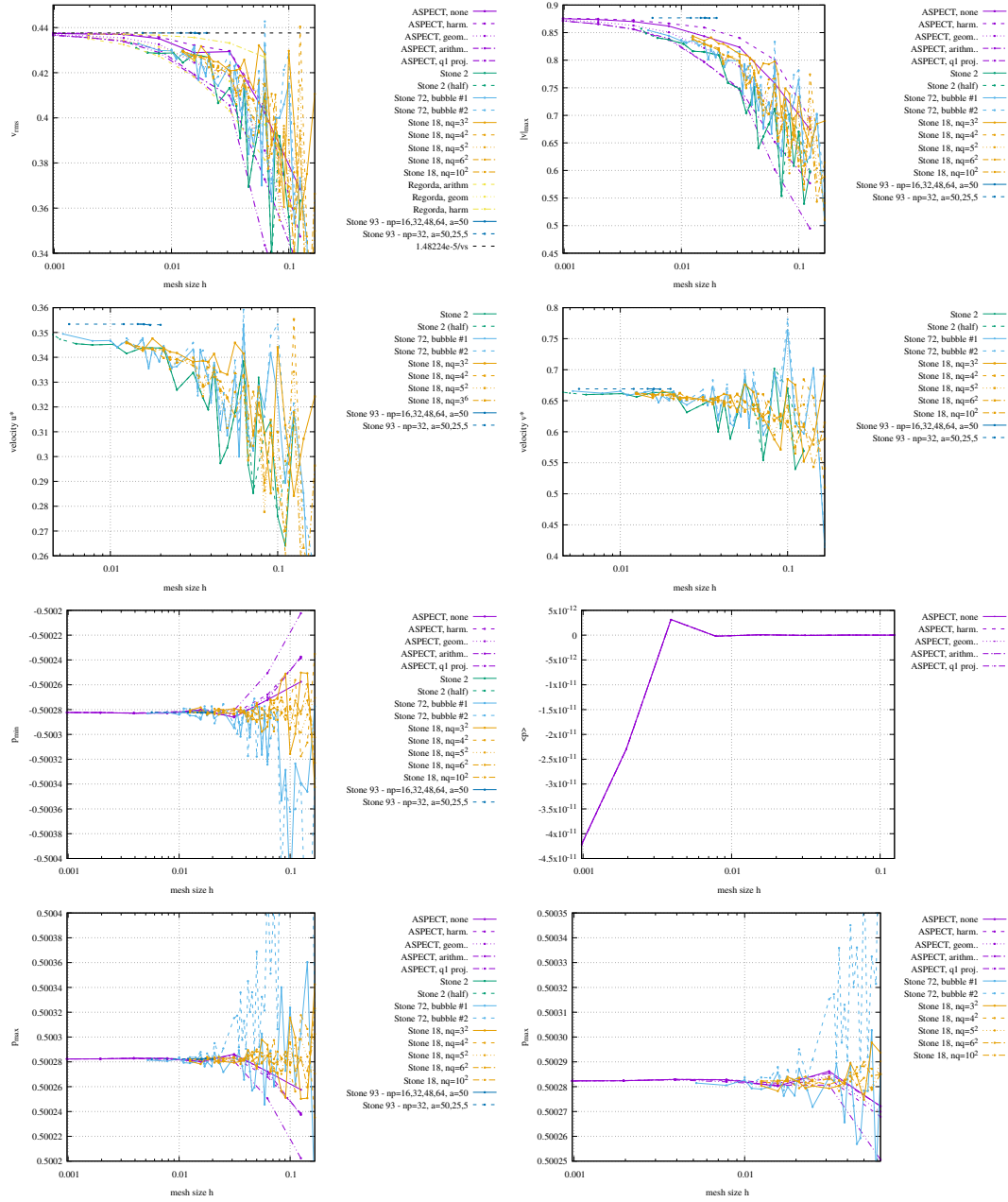


Octree-based ASPECT meshes. The background mesh is  $16^3$  and refinement is allowed to take place 4 times. Note that a special output is automatically generated in the code which subdivides all elements in 8 for visualisation purposes only. Initial refinements 0,1,2,3,4.

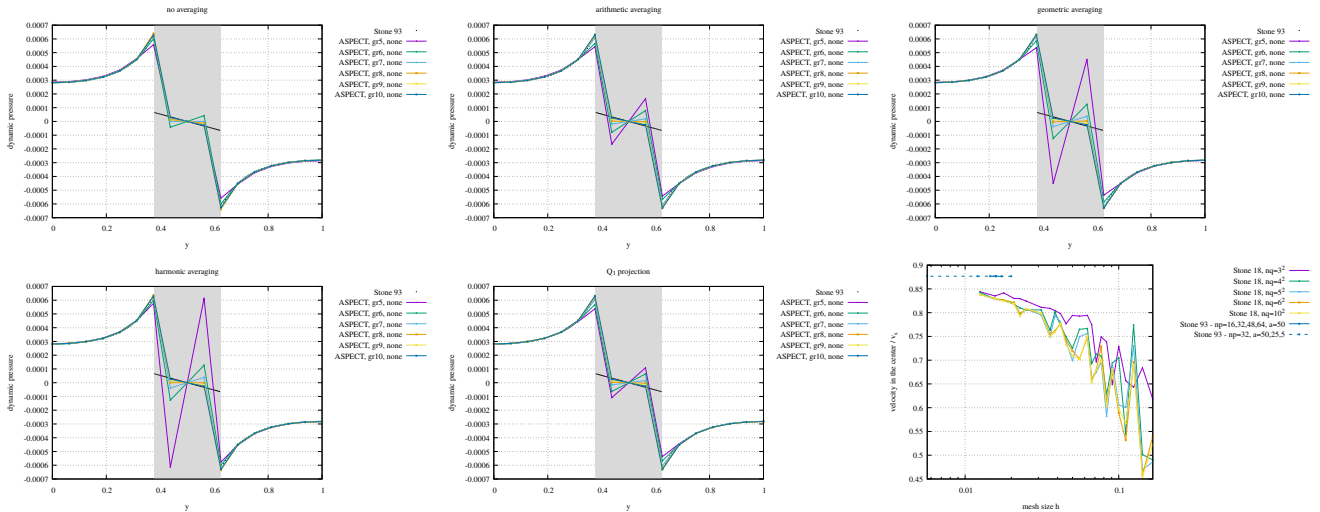
This corresponds to

```
subsection Mesh refinement
set Initial adaptive refinement = 0 -> 9
set Initial global refinement = 4
set Refinement fraction = 0.9
set Coarsening fraction = 0
set Strategy = composition
end
```

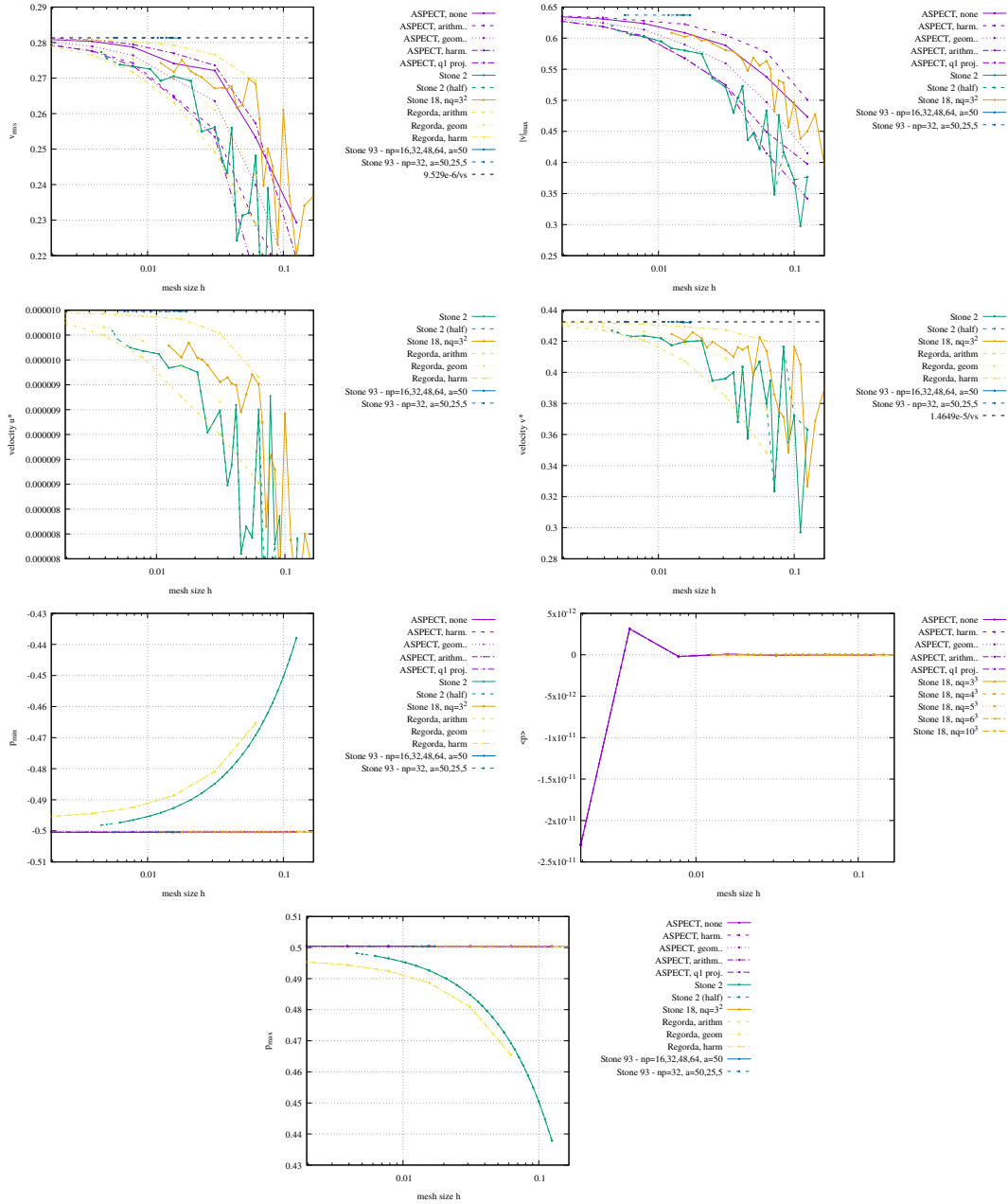
# FS results



Measurements obtained with ASPECT for various averaging schemes and with different stones.

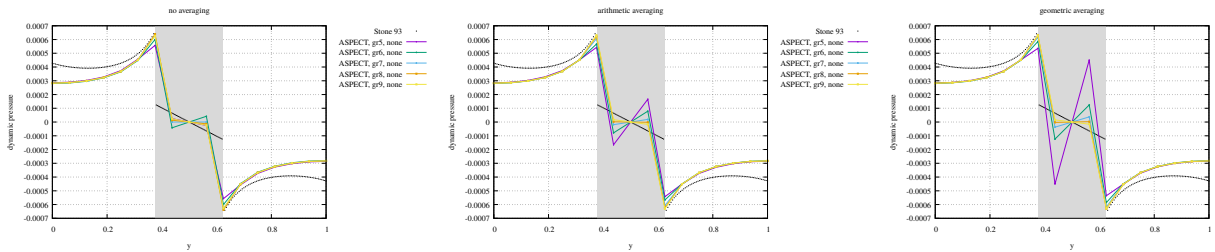


## NS results

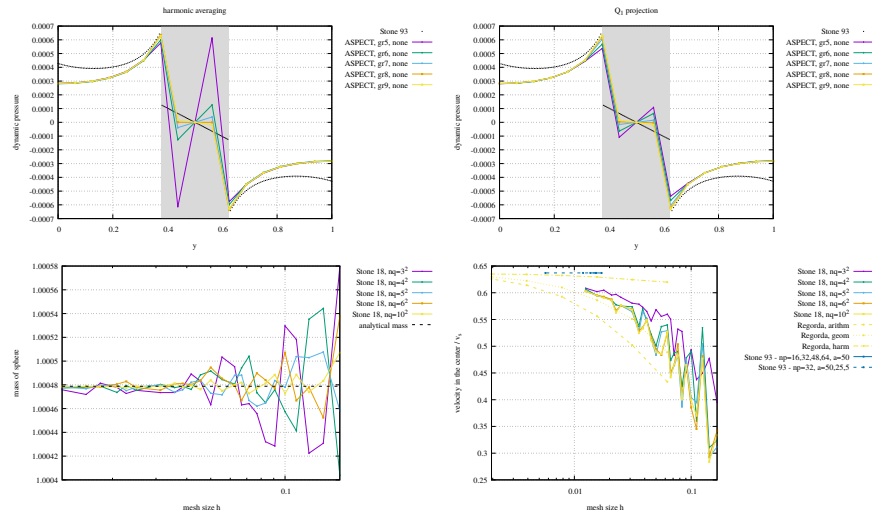


Measurements obtained with ASPECT for various averaging schemes and with different stones.

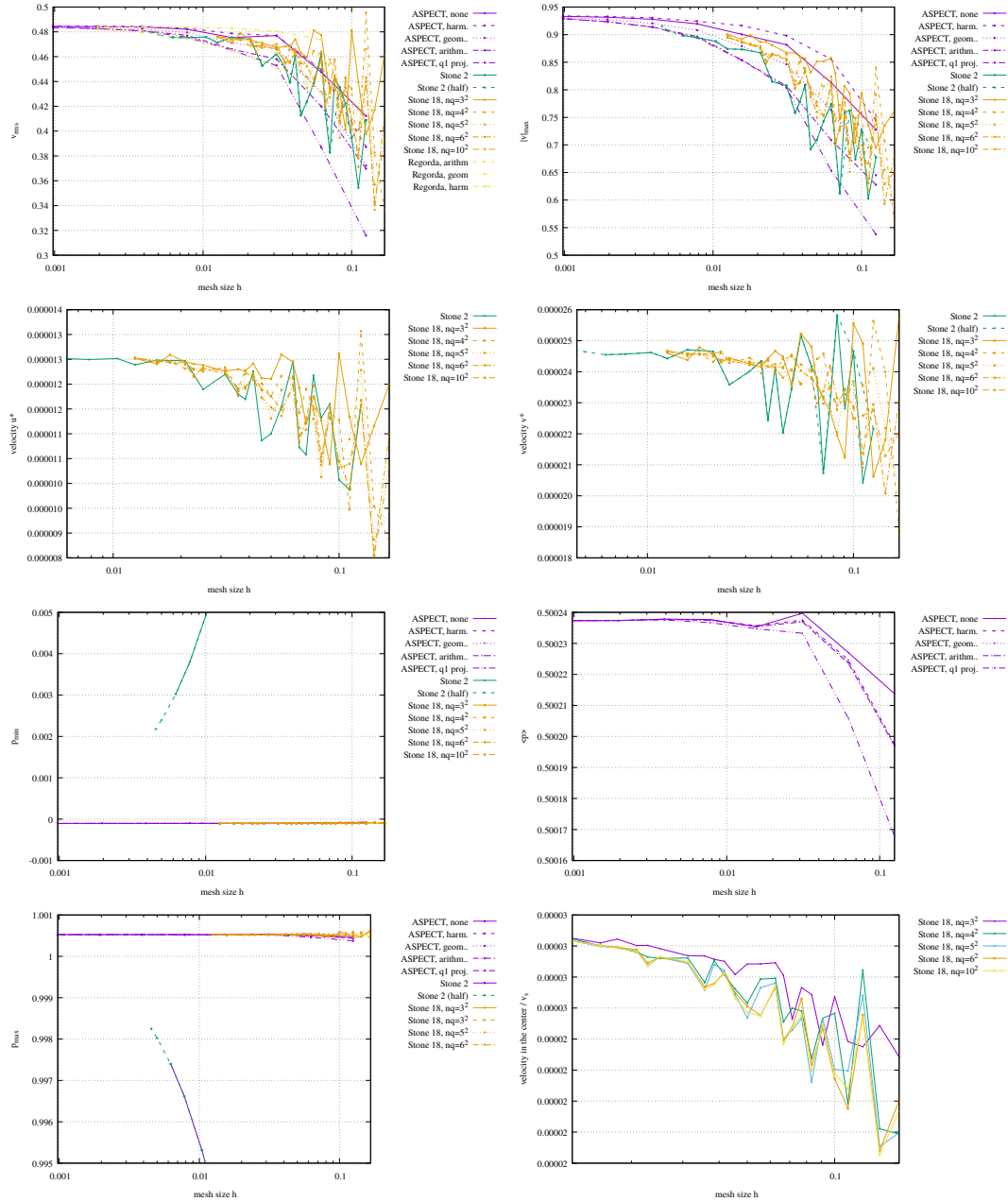
I have also retrieved the pressure at 16 equidistant locations on the  $x = 0.5$  line for all five averagings. Because the signal is dominated by the lithostatic pressure I have subtracted it from the data, so that I hereunder plot the dynamic pressure as a function of the  $y$  coordinate:



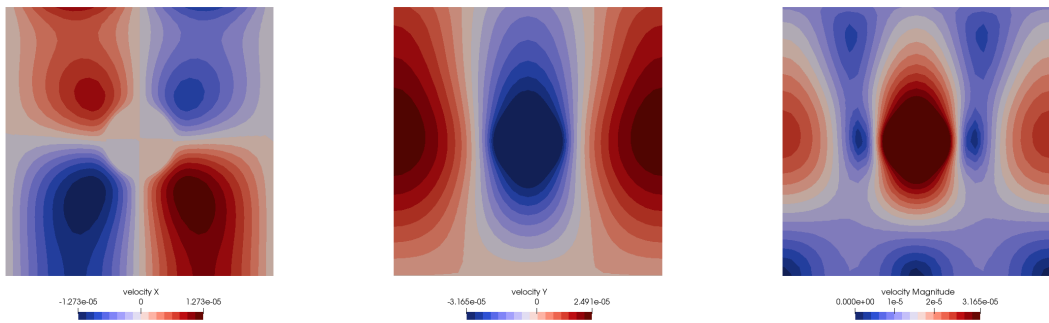




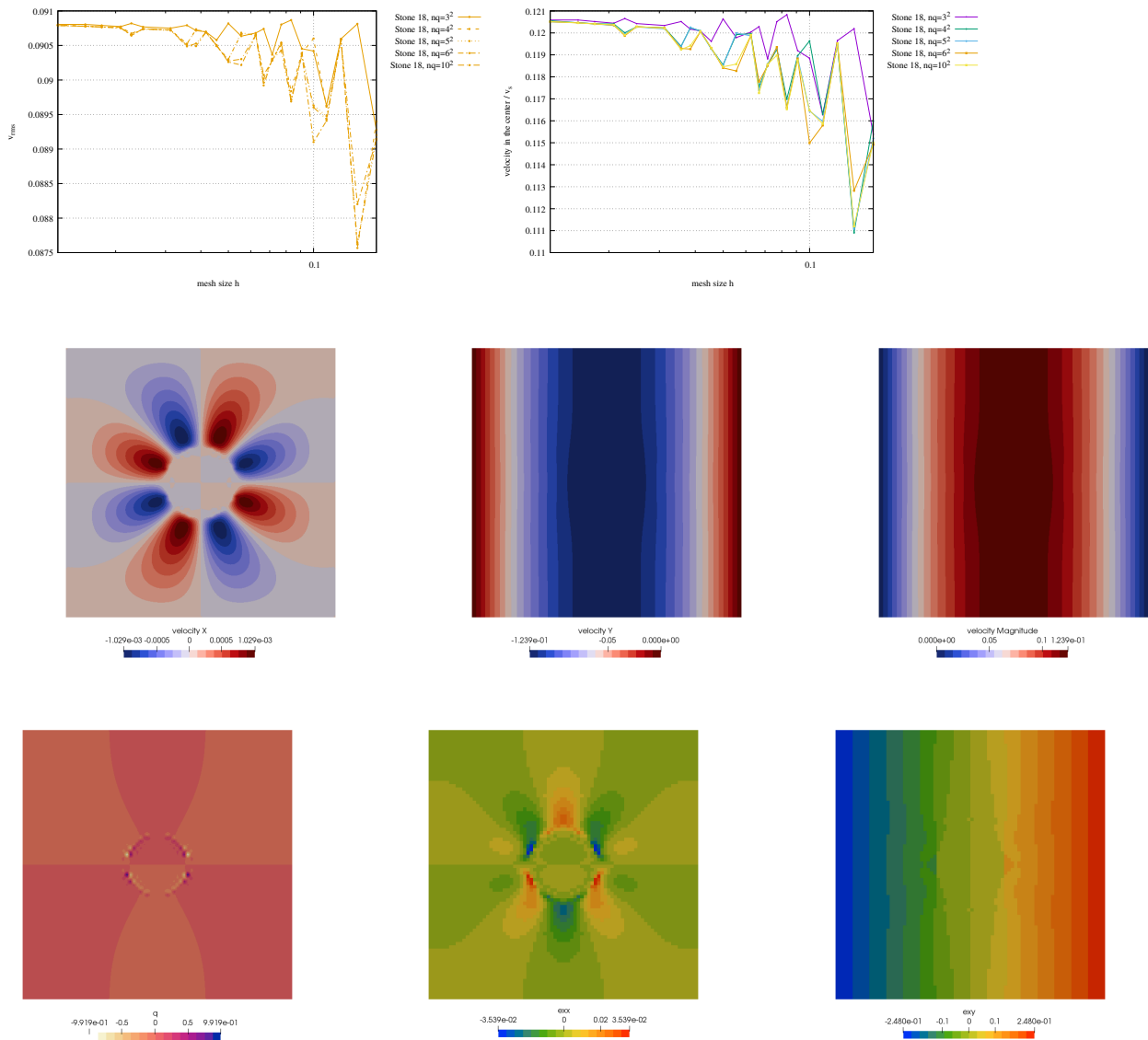
## OT results



Measurements obtained with ASPECT for various averaging schemes and with different stones.



BO results .



### 12.2.23 Stokes sphere (2D) in fluid with deformable free surface

benchmark\_stokes\_sphere\_fs\_2D.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/stokes\\_sphere\\_fs2D/](https://github.com/cedrict/fieldstone/tree/master/images/stokes_sphere_fs2D/)

The domain is a  $1 \times 0.75$  box. If sticky air is used, then its thickness should be 0.25 so that the domain is a unit square. The fluid is characterised by  $\rho_f = 1$  and  $\eta_f = 1$ . The sphere is characterised by  $\rho_s = 2$  and  $\eta_s = 10^3$ . The air is characterised by  $\rho_a = 0$  and  $\eta_a = 10^{-3}$ . Gravity is vertical with  $\vec{g} = -\vec{e}_y$ . The sphere has a radius  $R_s = 0.123456789$  and its center is at position  $\vec{r}_c = (0.5, 0.6)$ . Boundary conditions are free slip on all sides (unless a true free surface is used). Pressure is normalised so that its average is zero on the top (if no free surface is used). The model is run for 200s. The CFL number is set to 0.25 with a maximum time step of 0.5.

We wish to keep track of the following quantities as a function of time:

- the position and velocity of the sphere center,
- the minimum and maximum topography,
- the volume of fluid  $V_f(t)$ , sphere  $V_s(t)$  and air  $V_a(t)$
- root mean square velocity  $\mathbf{v}_{rms}$  for the whole domain, as well as for the air, fluid and sphere separately, and for the fluid+sphere,
- the maximum velocity and pressure in the domain,
- the time step value  $\delta t$ ,
- the average density<sup>11</sup> and viscosity in the domain:

$$\begin{aligned}\langle \rho \rangle(t) &= \frac{1}{L_x L_y} \iint \rho(x, y, t) dx dy = \frac{1}{L_x L_y} (V_a(t) \rho_a + V_f(t) \rho_f + V_s(t) \rho_s) \\ \langle \eta \rangle(t) &= \frac{1}{L_x L_y} \iint \eta(x, y, t) dx dy = \frac{1}{L_x L_y} (V_a(t) \eta_a + V_f(t) \eta_f + V_s(t) \eta_s)\end{aligned}\quad (12.279)$$

Initial values are

$$\langle \rho \rangle(0) = \frac{1}{L_x L_y} (V_a(0) \rho_a + V_f(0) \rho_f + V_s(0) \rho_s) = 0.25 * 0 + (0.75 - \pi R_s^2) * 1 + \pi R_s^2 * 2 = 0.75 + \pi R_s^2 \simeq 0.79788283$$

$$\langle \eta \rangle(0) = \frac{1}{L_x L_y} (V_a(0) \eta_a + V_f(0) \eta_f + V_s(0) \eta_s) = 0.25 * 10^{-3} + (0.75 - \pi R_s^2) * 1 + \pi R_s^2 * 10^3 \simeq 48.5851989989$$

- the min/max of the compositional fields when these are used;
- the velocity, pressure and material at position (0.5, 0.6);
- the pressure at position (0.5, 0).

---

<sup>11</sup>Because  $L_x L_y = 1$ , also equal to the total mass of the system

Participating codes:

- ASPECT uses  $Q_2 \times Q_1$  elements by default.  $Q_2 \times P_{-1}$  elements can be used by setting

```
subsection Discretization
 set Use locally conservative discretization = true
end
```

The default stabilisation method when compositional fields are used is the entropy viscosity method, but SUPG has also been implemented and can be triggered with

```
subsection Discretization
 subsection Stabilization parameters
 set Stabilization method = SUPG
 end
end
```

The default mesh settings are as follows:

```
subsection Mesh refinement
 set Initial adaptive refinement = 1
 set Initial global refinement = 6
 set Refinement fraction = 0.9
 set Strategy = composition
 set Coarsening fraction = 0.1
end
```

In the results hereafter when the combination 6-0 or 7-0 are mentioned, this means that the coarsening fraction has been set to zero and these correspond then to regular meshes with 64x64 and 128x128 elements respectively.

Active particles can also replace the compositional fields and this is how it is triggered from the input file:

```
subsection Compositional fields
 set Number of fields = 2
 set Names of fields = sphere , air
 set Compositional field methods = particles , particles
 set Mapped particle properties = sphere:initial sphere , air: initial air
end

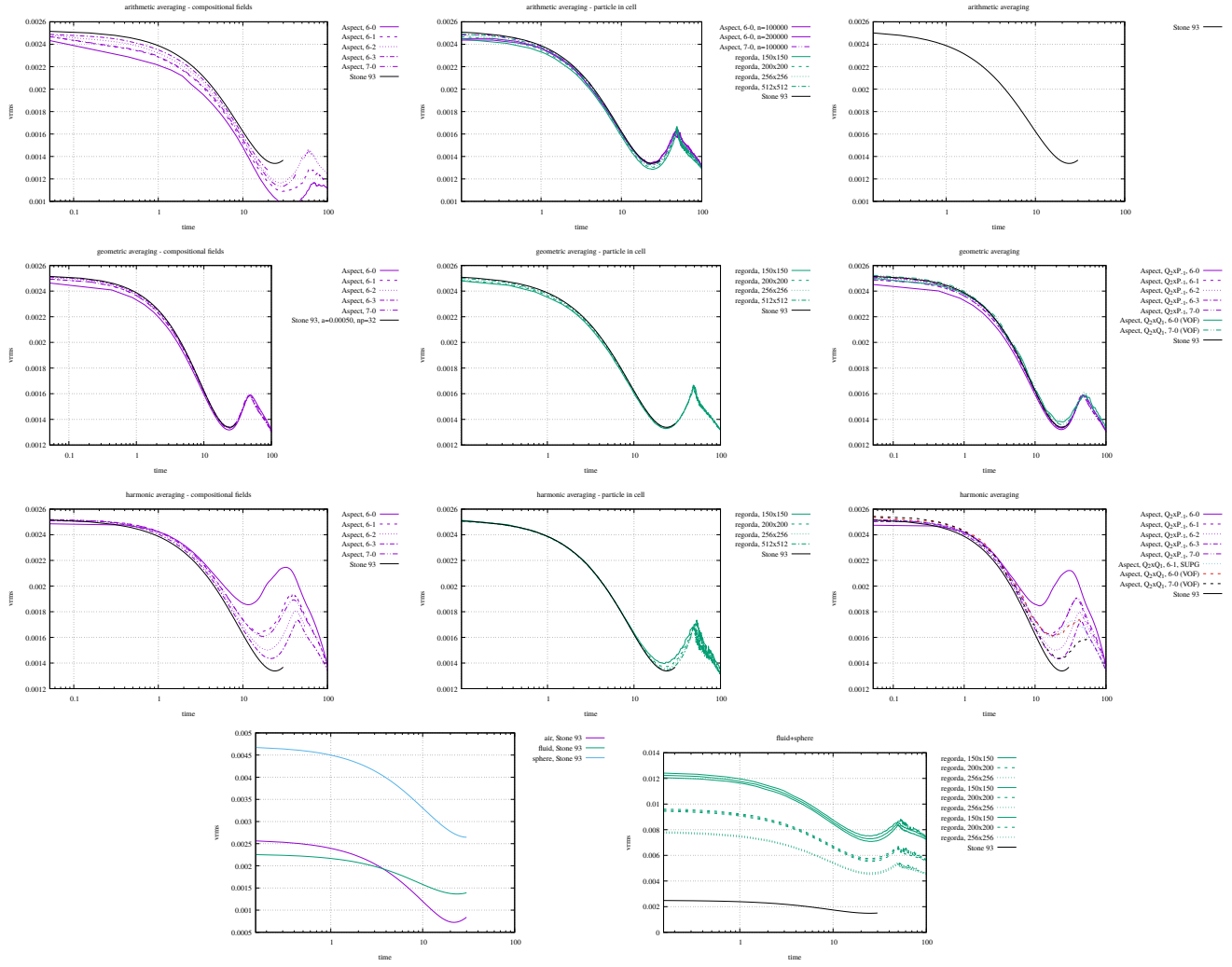
subsection Postprocess
 set List of postprocessors = visualization , ... , particles
 subsection Visualization
 set List of output variables = density , viscosity , strain rate
 set Time between graphical output = 1
 end
 subsection Particles
 set Number of particles = 100000
 set Time between data output = 0
 set Data output format = vtu
 set List of particle properties = velocity , initial composition , initial position #, integrated strain
 set Interpolation scheme = cell average
 set Update ghost particles = true
 set Particle generator name = random uniform
 end
end
```

- Stone 93. Code based on unstructured mesh of Crouzeix-Raviart triangular elements. The resolution is controlled by the minimum area of the triangles as passed as argument to the triangle mesher, and the parameter  $np$  which controls the number of points on the hull ( $np$  on each side), the surface ( $5 * np$ ) and the sphere ( $5 * np$ );
- Alessandro Regorda's code: The number of markers is fixed per element with random distribution. At the beginning of the simulation there are: 562500 for 150x150 grid with 25 markers per element, 600,000 for 200x200 grid with 15 markers per element, 655,360 for 256x256 grid with 10 markers per element, 2,621,440 for 512x512 grid with 10 markers per element.

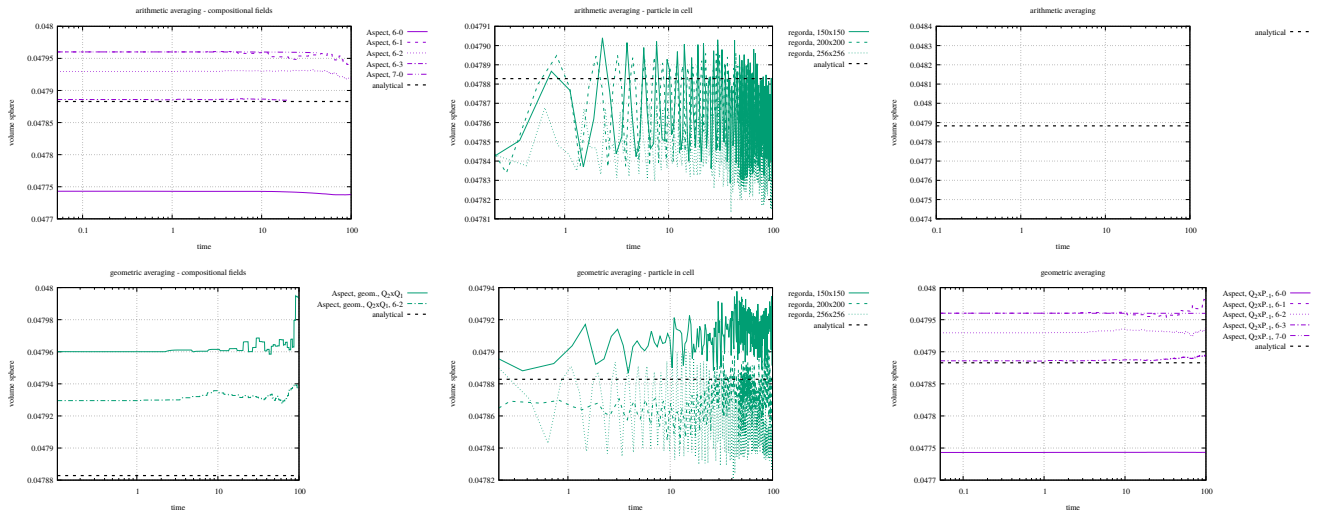
The advection is RK4. The code maintains the number of markers per element between half of the initial number and the initial number plus half (e.g. in 512x512 markers are between 5 and 15). When in an element there are less markers than the minimum it adds random markers to reach the minimum, while if the number is higher than the maximum some of them are deleted. In this way elements are never empty. When new markers are added they assume the type of the nearest marker. Averaging only applies to viscosity, density is always arithmetically averaged.

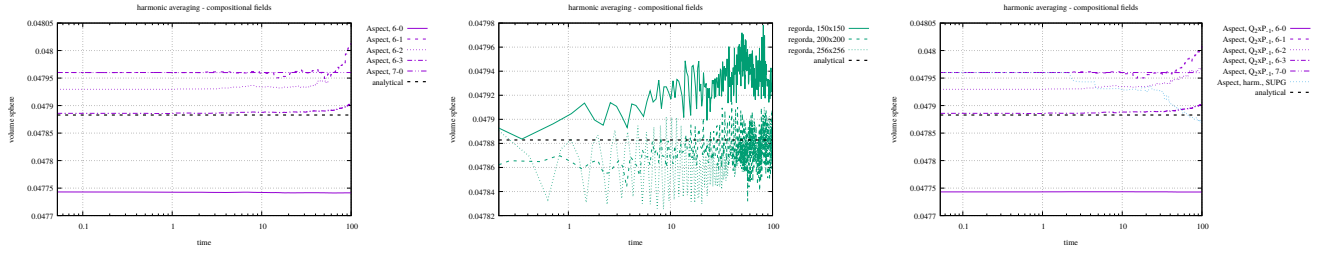
In what follows 6-0 and 7-0 correspond to regular grids (no coarsening, no refinement)

# Root mean square velocity

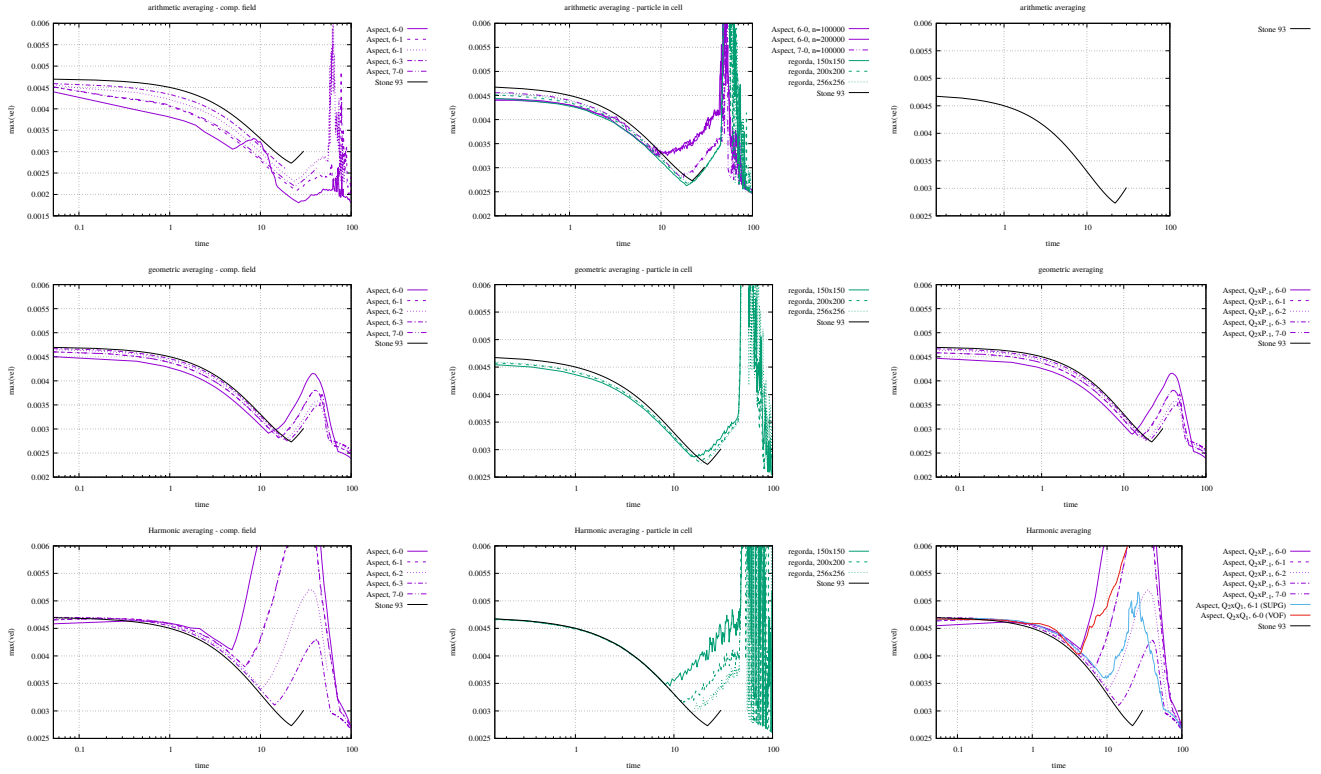


# Volumes

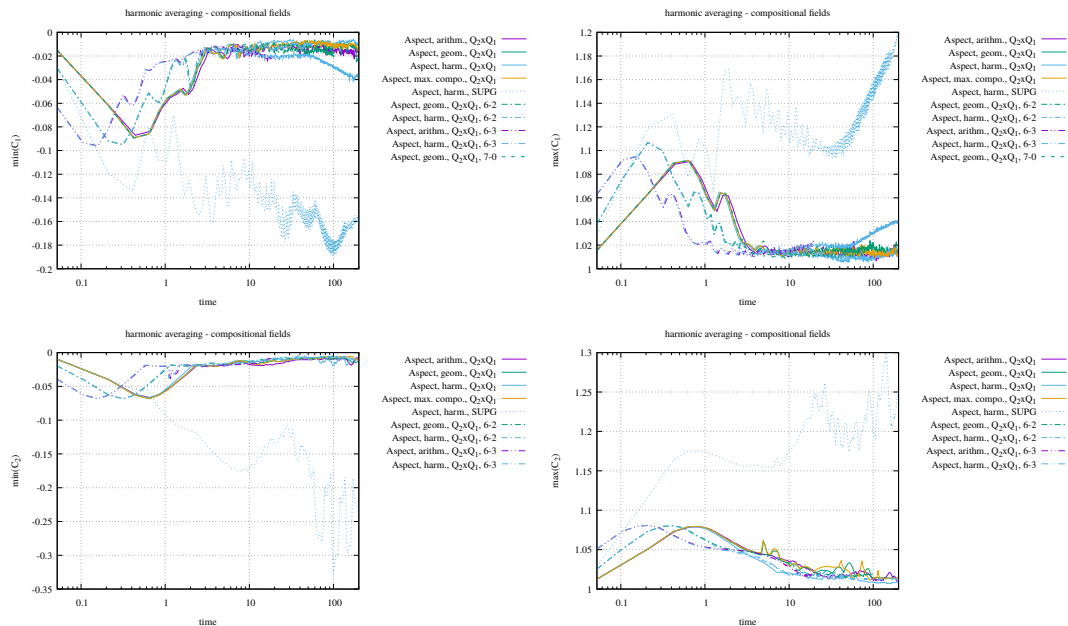




## Maximum velocity

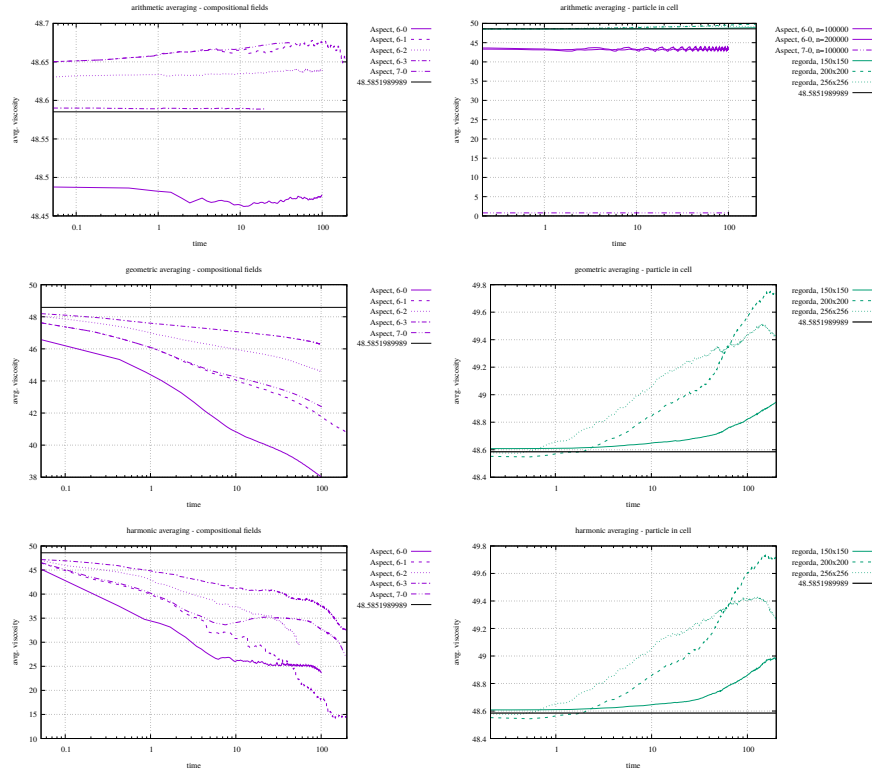


## Compositions min/max

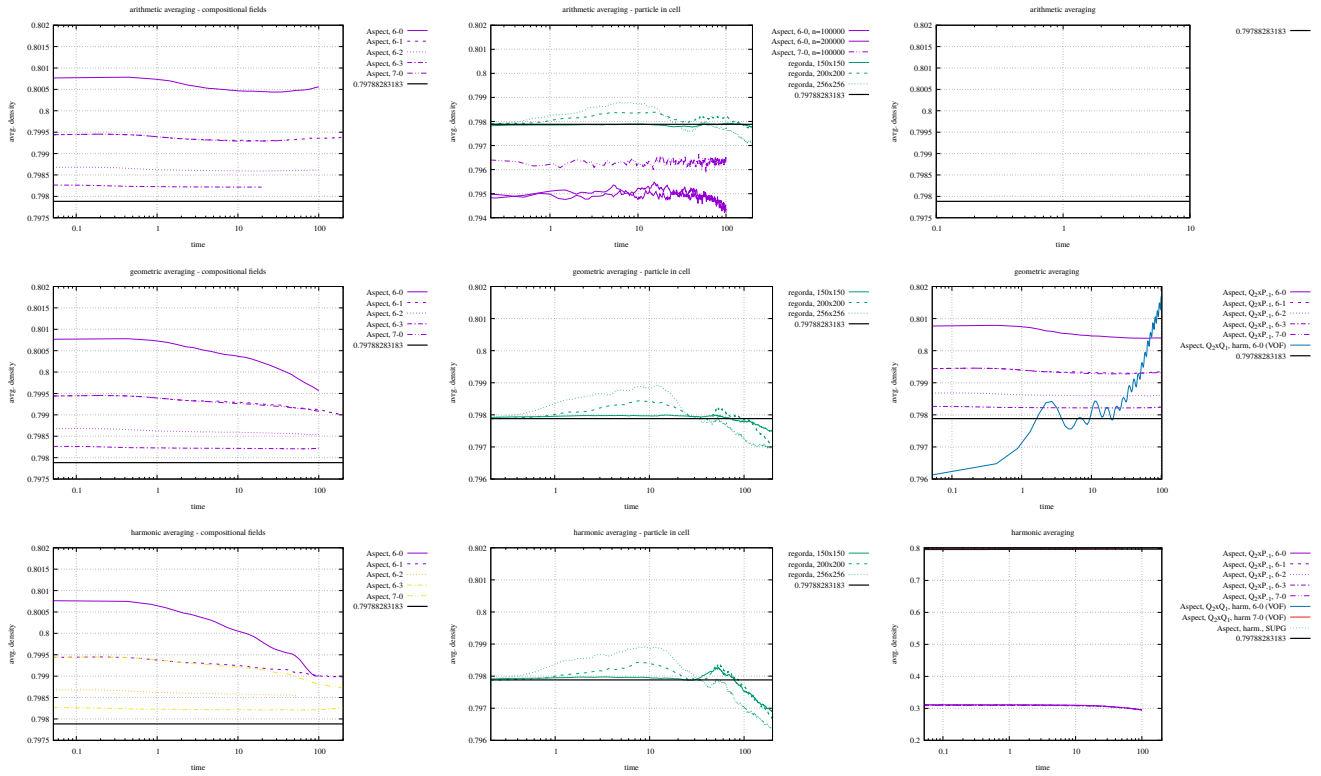




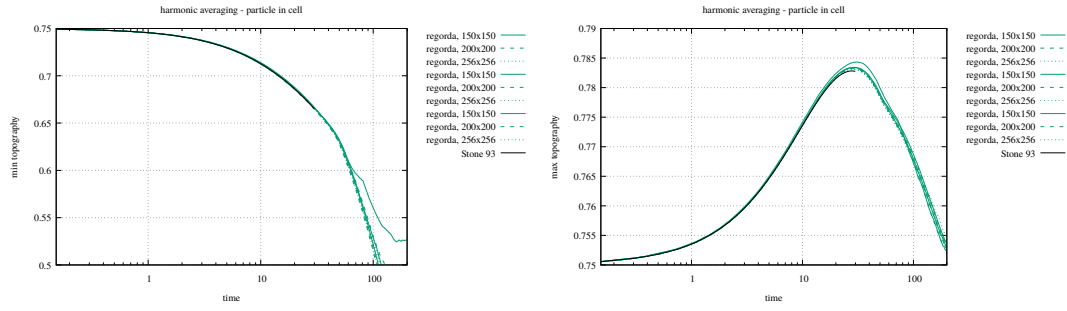
## Viscosity volume average



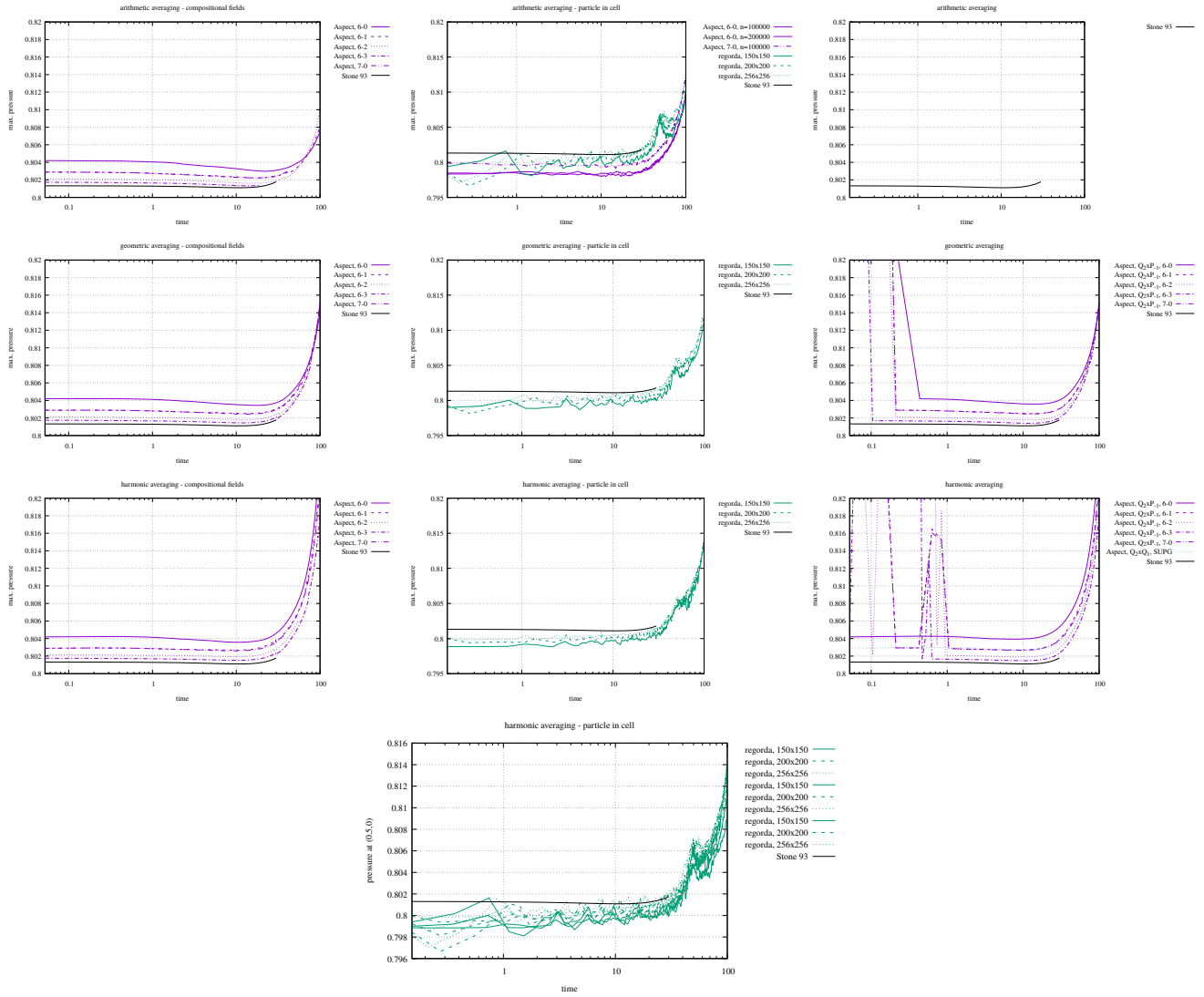
## Density volume average



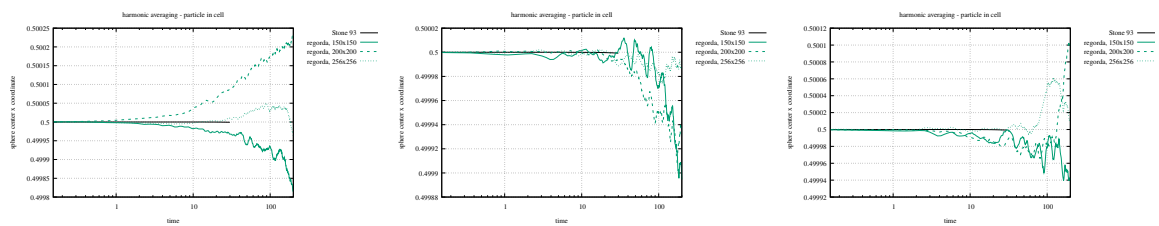
## Topography

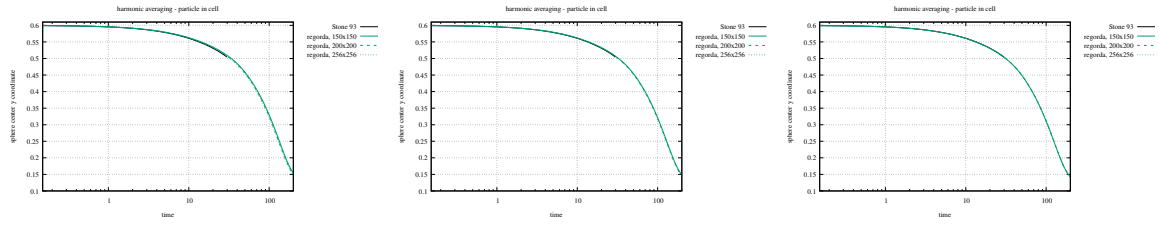


## Pressure measurements

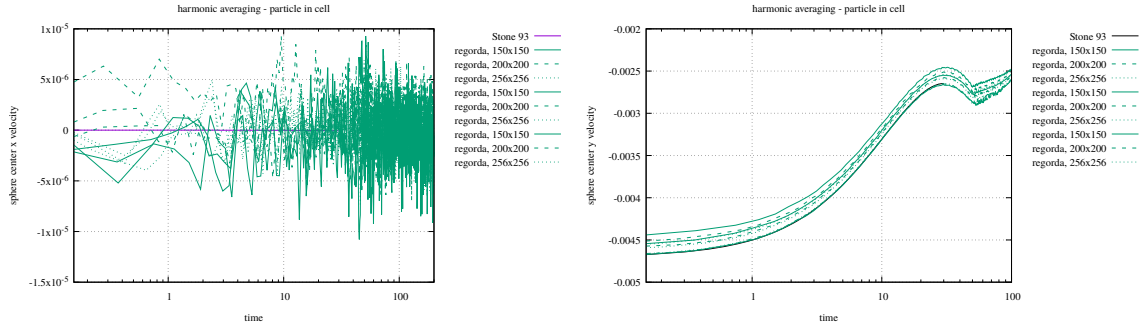


## Position of the sphere center

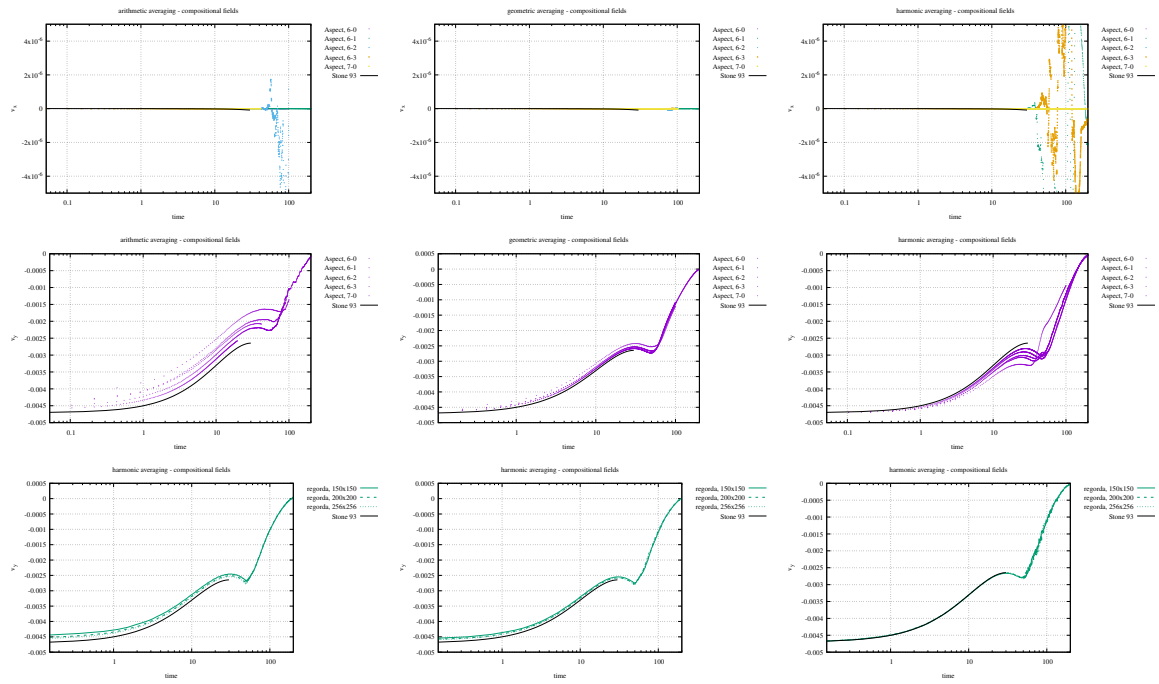




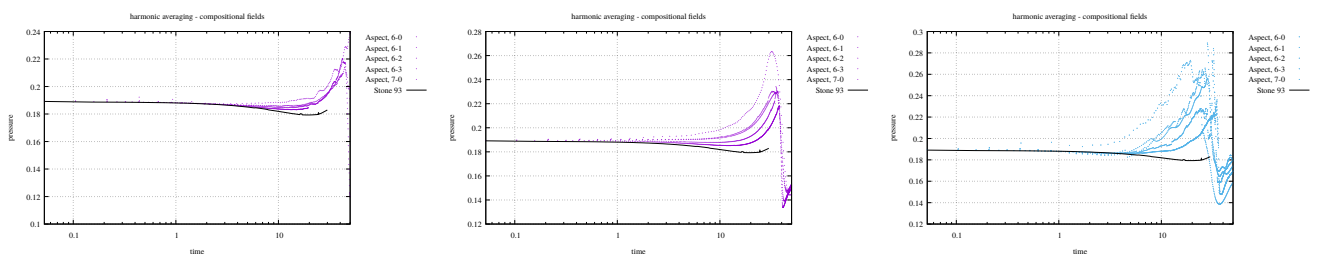
## Velocity of the sphere center .

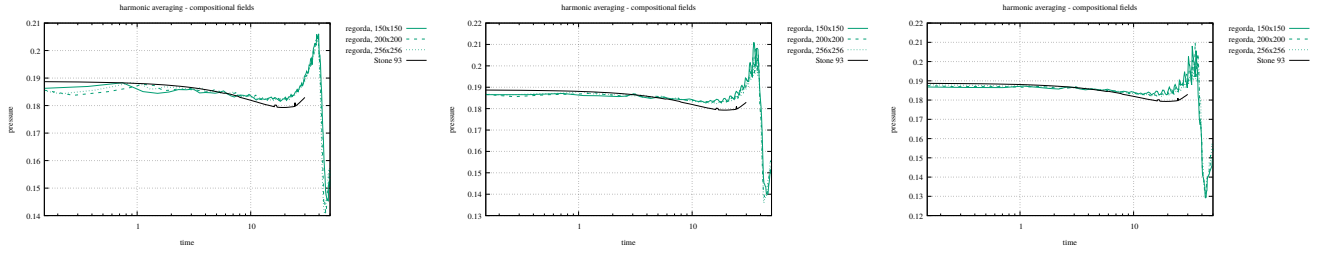


## Velocity at (0.5,0.6) location .

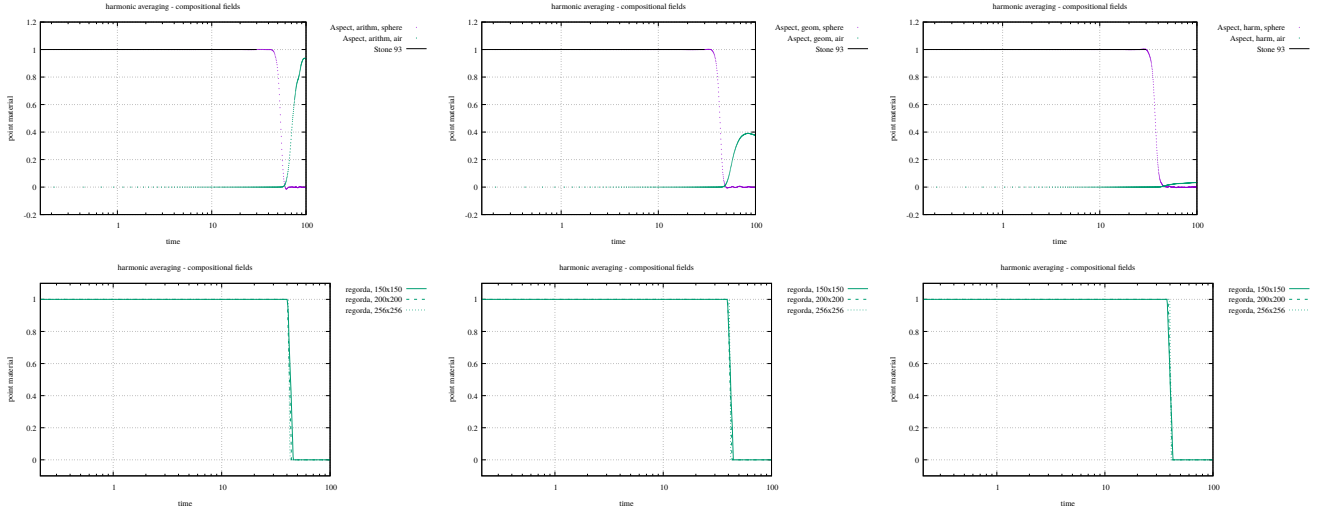


## Pressure at (0.5,0.6) location .

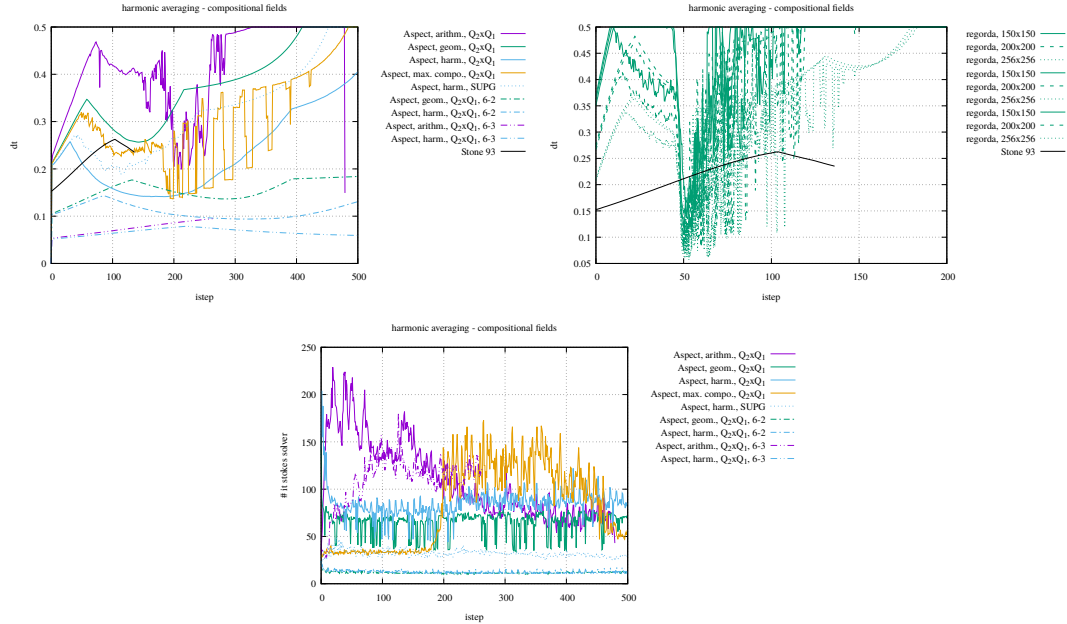




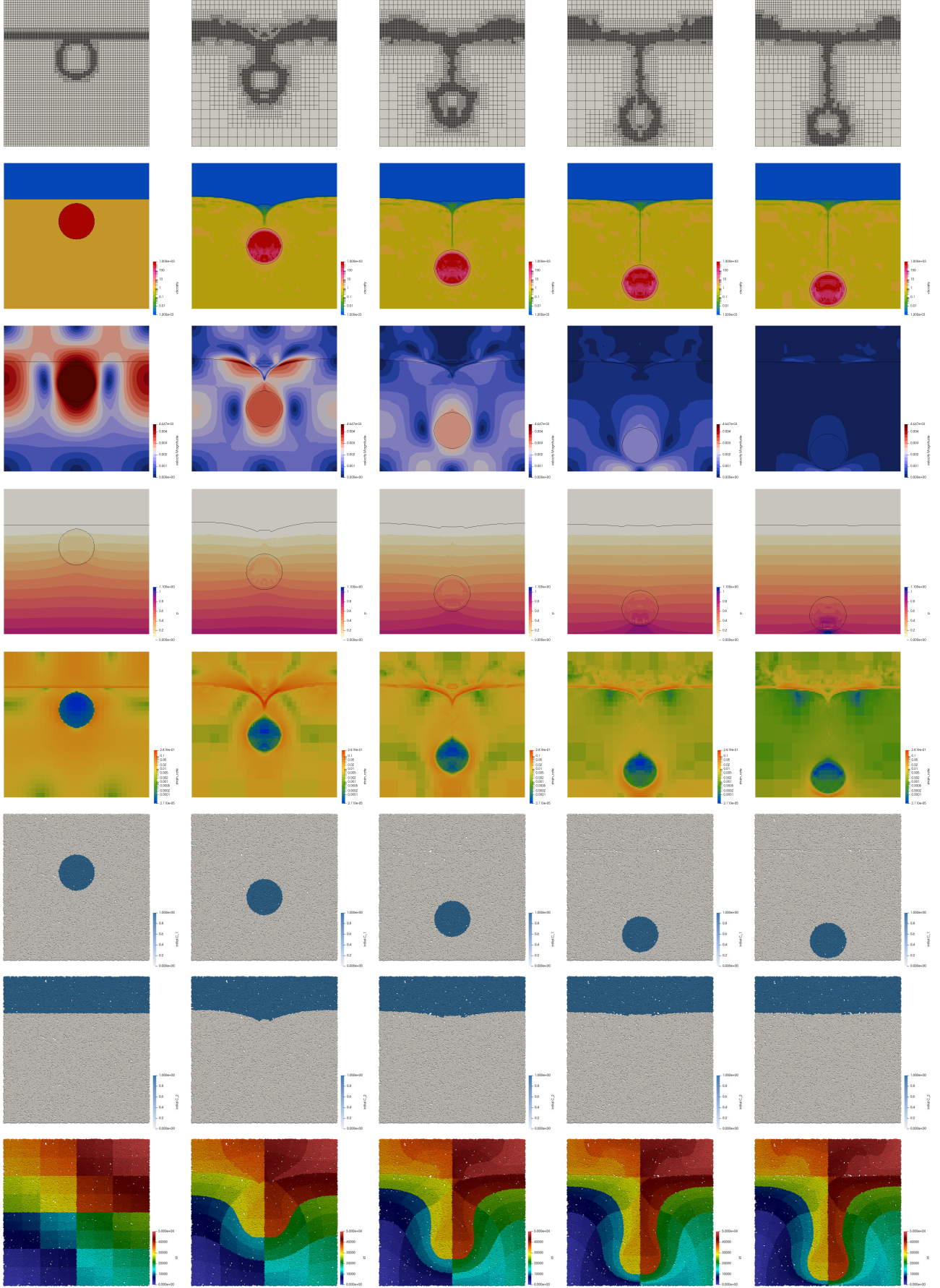
## Material at (0.5,0.6) location .



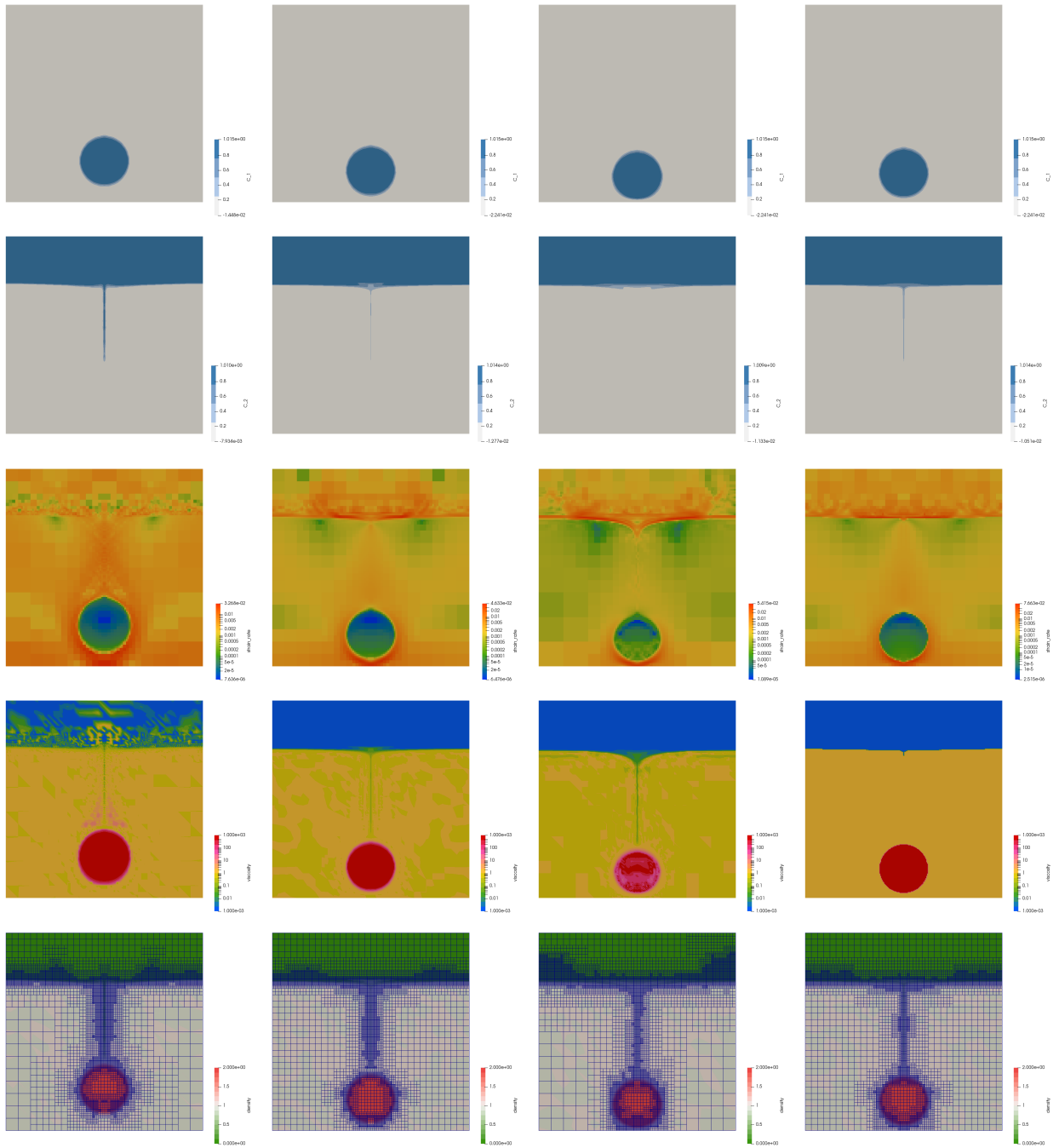
## Timestep and solver convergence



Stone 93 results seem to be most influenced by the resolution on the sphere and surface than the resolution in the fluid. Some conclusions: arithmetic yields very high inner iteration counts. Discontinuous pressure also better on that topic. Geometric averaging yields very good agreement for vrms. Funny enough, geometric does not correspond to a physical arrangement of viscous dampers... Arithm and harm do ultimately converge towards geom but at very high resolution. Using no amr does not change things that much.



ASPECT results: fields and passive markers time evolution. Last row is the particle id, between 0 and 50000. System at times 0,50,100,150,200. Harmonic averaging.



Obtained with Aspect. From left to right: Arithmetic, geometric, harmonic, maximum composition, all at  $t = 200$ .



This is the ASPECT input file for this benchmark:

```

set Dimension = 2
set Start time = 0
set End time = 200
set Use years in output instead of seconds = false
set CFL number = 0.25
set Output directory = output-stokes
set Maximum time step = 0.5
set Pressure normalization = surface

subsection Solver parameters
 subsection Stokes solver parameters
 set Number of cheap Stokes solver steps = 0
 end
end

subsection Geometry model
 set Model name = box
 subsection Box
 set X extent = 1
 set Y extent = 1
 end
end

subsection Boundary velocity model
 set Tangential velocity boundary indicators = left , right , bottom , top
end

subsection Material model
 set Model name = multicomponent
 subsection Multicomponent
 set Densities = 1, 2, 0
 set Viscosities = 1, 1000, 0.001
 #set Viscosity averaging scheme = maximum composition
 #set Viscosity averaging scheme = arithmetic
 #set Viscosity averaging scheme = geometric
 set Viscosity averaging scheme = harmonic
 set Thermal expansivities = 0
 end
end

subsection Gravity model
 set Model name = vertical
 subsection Vertical
 set Magnitude = 1
 end
end

subsection Boundary temperature model
 set Fixed temperature boundary indicators = bottom , top
 set List of model names = box
end

subsection Initial temperature model
 set Model name = function

 subsection Function
 set Function expression = 0
 end
end

subsection Compositional fields
 set Number of fields = 2
end

subsection Initial composition model
 set Model name = function

 subsection Function
 set Variable names = x,y
 set Function constants = r=0.123456789, xc=0.5, yc=0.6
 set Function expression = if(sqrt((x-xc)*(x-xc)+(y-yc)*(y-yc)) < r, 1, 0) ; if (y>0.75,1,0)
 end
end

subsection Mesh refinement
 set Initial adaptive refinement = 1
 set Initial global refinement = 6
 set Refinement fraction = 0.9
 set Strategy = composition
 set Coarsening fraction = 0.1
end

subsection Postprocess
 set List of postprocessors = visualization , velocity statistics , composition statistics , pressure statistics ,
 material statistics , global statistics , point values , particles

 subsection Point values
 set Evaluation points = 0.5,0.6
 end

 subsection Particles
 set Number of particles = 50000
 set Time between data output = 1
 set Data output format = vtu
 set List of particle properties = initial composition , initial position
 end

 subsection Visualization
 set List of output variables = density , viscosity , strain rate
 set Time between graphical output = 1
 end
end

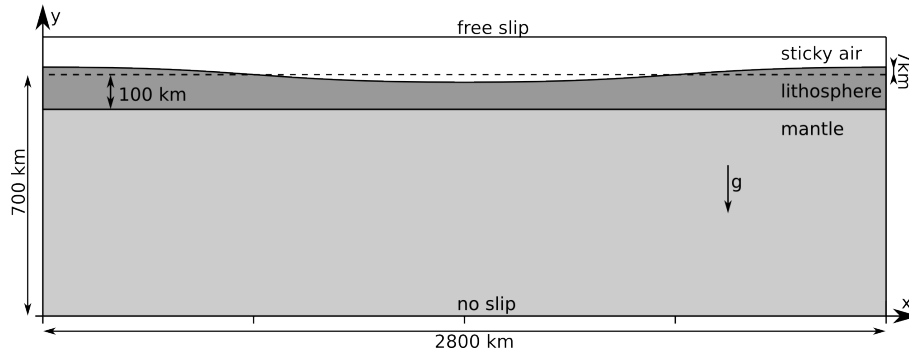
```

### 12.2.24 Relaxation of topography (Crameri *et al.*, 2012)

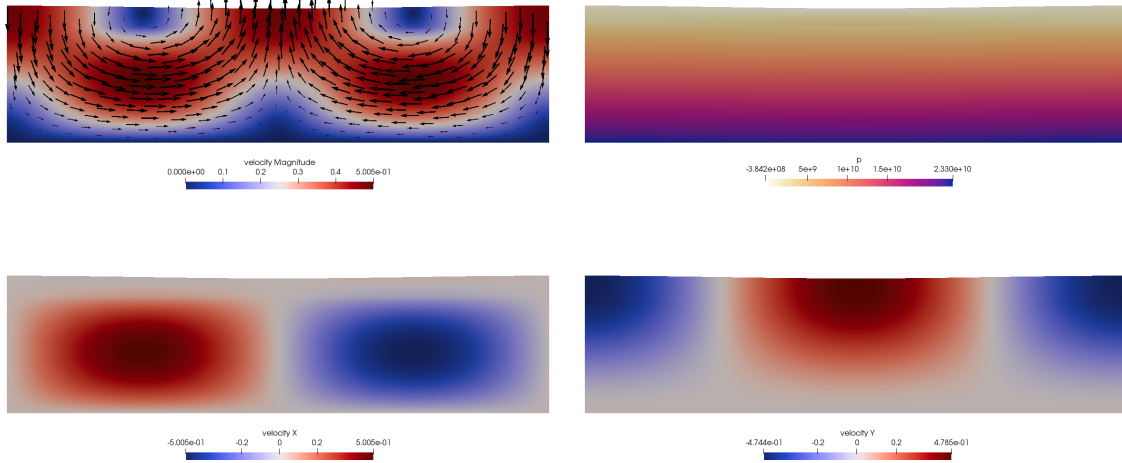
benchmark\_crsg12.tex

This benchmark was first presented in Crameri *et al.* (2012) [285] and is also presented in Hillebrand *et al.* (2014) [571]. It is designed to test the accuracy of the free surface representation in geodynamics code.

The model box spans 2800km by 700–1100km (greater model height is necessary when employing sticky air on top). The initial condition is specified by a mantle of 600km thickness, overlain by a cosine shaped, 93 – 107km-thick lithosphere:



The sticky air layer has a thickness varying between 10 and 400km. The lithosphere is a highly viscous, dense medium ( $\rho_L = 3300\text{kg m}^{-3}$ ,  $\mu_L = 10^{23}\text{Pa s}$ ). The underlying ambient mantle has a density  $\rho_M = 3300\text{kg m}^{-3}$  and a viscosity  $\mu_M = 10^{21}\text{Pa s}$ .



Results at  $t = 0$  obtained with ASPECT by running the included cookbook.

The sticky air layer on the top has a density  $\rho_{air} = 0\text{kg m}^{-3}$  and a viscosity  $\mu_{air} = 10^{18} - 10^{20}\text{Pa s}$  and is bordered by a free-slip top boundary condition. Free slip is also imposed at the sides while the bottom boundary is set to no slip condition.

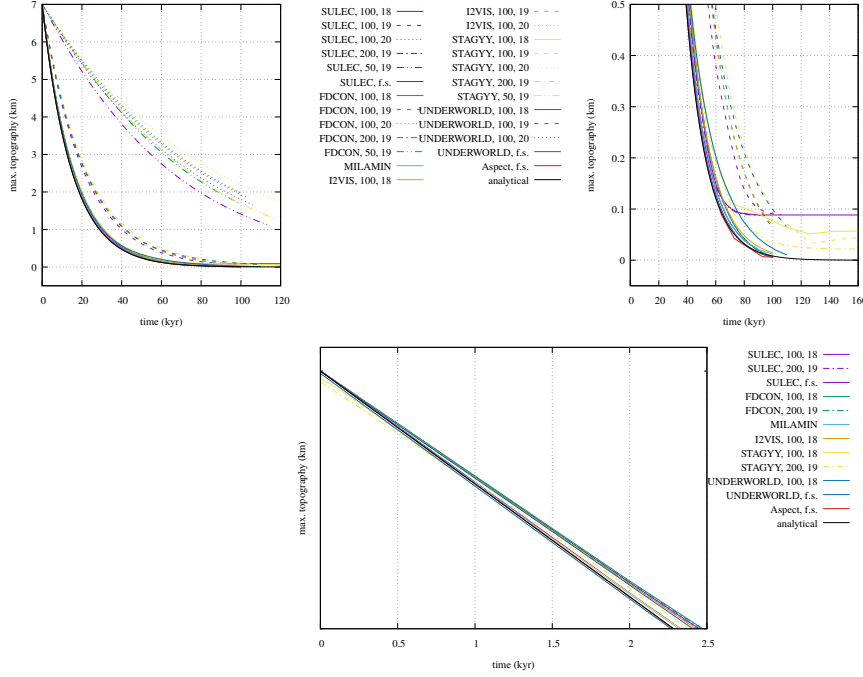
The setup for the real free-surface model is identical to the setup described above, but the weak surface layer is removed and replaced by zero normal stress boundary conditions.

An analytical solution is presented by [1035]: the maximum topography at time  $t$  can be derived analytically using the relaxation rate  $\gamma$  and from the initial maximum topography  $h_{init}$ :

$$h_{analytic} = h_{init} \exp(\gamma t) \quad (12.280)$$



where  $t = 14.825\text{kyr}$  is the characteristic relaxation time and  $\gamma = -0.2139 \times 10^{-11}\text{s}^{-1}$  is the characteristic relaxation rate of the three-layer case at a given wavelength of  $2800\text{km}$ . It should be noted that these values are valid for infinitesimal amplitudes, whereas deviations are to be expected for small but finite amplitudes. In particular, keeping the interface between the middle and lower layer flat and assuming a finite amplitude of the interface between the upper and middle layer implies that the thickness of the highly viscous middle layer varies laterally by  $\pm 7\%$  (in the case of an initial maximum topography of  $7\text{km}$ ). This variation increases the effective viscous flexural rigidity and leads to a slightly longer relaxation time. The system is let to relax over time (typically  $200\text{kyrs}$ ) and the position of the free surface at  $x = 0$  is recorded over time.



Data from the original paper. Aspect data are obtained by running the available example in the code.

## 12.2.25 3D spherical shell convection benchmark

benchmark\_sscb3D.tex

The governing equations are those of an incompressible fluid with constant viscosity whose density depends on temperature. The fluid convects in a three-dimensional hollow sphere. The inner radius is  $11/9$  while the outer radius is  $20/9$ , so that the depth of the mantle is exactly 1. Boundary conditions are free-slip at the top and bottom boundaries and isothermal with nondimensional temperatures of 0 and 1 at the top and bottom boundaries, respectively.

The dynamics of the system are governed by the Rayleigh number:

$$\text{Ra} = \frac{\rho_0 \alpha g \Delta T \Delta R^3}{\kappa \eta}$$

Initial conditions for temperature are given as a function of coordinates with perturbations at some given spherical harmonics superimposed on a conductive temperature profile:

$$T(r, \theta, \phi) = \frac{R_i(r - R_o)}{r(R_i - R_o)} + \sum_m [\epsilon_{c,m} \cos m\phi + \epsilon_{s,m} \sin m\phi] p_{lm}(\theta) \sin \left( \pi \frac{r - R_i}{R_o - R_i} \right)$$

The first term represents a purely conductive temperature profile, while the second term is a perturbation to this profile, determining the final patterns of polyhedral symmetry.  $l$  and  $m$  are the

spherical harmonic degree and order, respectively.  $\epsilon_{c,m}$  and  $\epsilon_{s,m}$  are the magnitudes of the individual spherical harmonic constituents.  $p_{lm}$  is a normalized associated Legendre polynomial that is related to the associated Legendre polynomial  $P_{lm}$  as<sup>12</sup>:

$$p_{lm}(\theta) = \sqrt{\frac{(2l+1)(l-m)!}{2\pi(1+\delta_{m0})(l+m)!}} P_l^m(\theta)$$

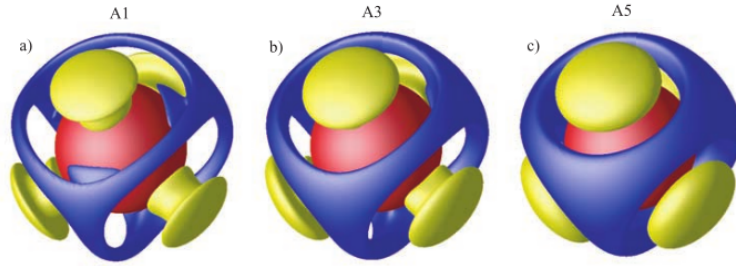
where  $P_l^m$  are the (unnormalized) associated Legendre functions and  $\delta_{m0}$  is the Kronecker delta.

Note that there are somewhat subtle notation differences between the papers reporting on this benchmark with regards to the spherical harmonic parts and the normalisation term. Also note that the  $(1+\delta_{m0})$  term is often omitted from spherical harmonics libraries (see ASPECT documentation).


In each case, we compute, as a function of time, Nusselt numbers for both the top and bottom boundaries,  $Nu_t$  and  $Nu_b$ , averaged temperature for the whole mantle,  $\langle T \rangle$  and averaged RMS velocity  $v_{rms}$  for the whole mantle with

$$Nu_t = \frac{R_o(R_o - R_i)}{R_i} Q_t \quad Nu_b = \frac{R_i(R_o - R_i)}{R_o} Q_b$$

where  $Q_t$  and  $Q_b$  are the surface and bottom heat fluxes.



Taken from Zhong *et al.* (2008) [1412]. Representative steady state residual temperature  $\delta T = T(r, \theta, \phi) - \langle T(r) \rangle$  for cases a, b, c.

 **Relevant Literature:** Zhong *et al.* [1412] (CITCOMS ), Arrial *et al.* [29] (CITCOMS vs. radial basis function), Shahnas *et al.* [1151] (own code), Liu & King [801] (ASPECT ).

## 12.2.26 2D convection benchmark ('Blankenbach *et al.* benchmark')

benchmark\_blb89.tex

The abstract of the original publication by Blankenbach *et al.* (1989) [95] reads:

*We have carried out a comparison study for a set of benchmark problems which are relevant for convection in the Earth's mantle. The cases comprise steady isoviscous convection, variable viscosity convection and time-dependent convection with internal heating. We compare Nusselt numbers, velocity, temperature, heat-flow, topography and geoid data. Among the applied codes are finite-difference, finite-element and spectral methods. In a synthesis we give best estimates of the 'true' solutions and ranges of uncertainty. We recommend these data for the validation of convection codes in the future.*

The temperature is fixed to zero on top and to  $\Delta T$  at the bottom, with reflecting symmetry at the sidewalls (i.e.  $\partial_x T = 0$ ) and there are no internal heat sources. Free-slip conditions are implemented on all boundaries.

---

<sup>12</sup>There is a typo in [1151]

The Rayleigh number is given by

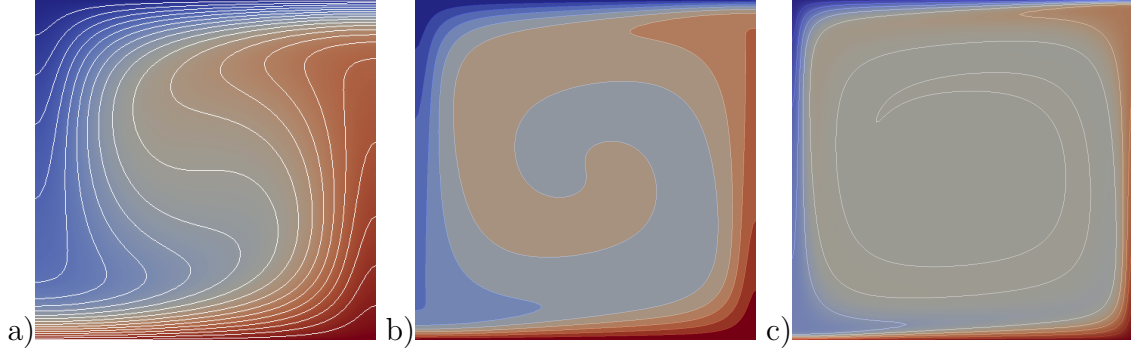
$$\text{Ra} = \frac{\alpha g_y \Delta T h^3}{\kappa \nu} = \frac{\alpha g_y \Delta T h^3 \rho^2 c_p}{k \mu}$$

The initial temperature field is given by

$$T(x, y) = (1 - y) - 0.01 \cos(\pi x) \sin(\pi x)$$

The perturbation in the initial temperature fields leads to a perturbation of the density field and sets the fluid in motion.

Depending on the initial Rayleigh number, the system ultimately reaches a steady state after some time.



Temperature fields at steady-state for  $\text{Ra} = 10^4$  (a),  $\text{Ra} = 10^5$  (b),  $\text{Ra} = 10^6$  (c). Obtained with ELEFANT code [1257].

a)

| Mesh size               | $\frac{1}{16}$ | $\frac{1}{32}$ | $\frac{1}{64}$ | $\frac{1}{128}$   | $\frac{1}{256}$   | Reference             |
|-------------------------|----------------|----------------|----------------|-------------------|-------------------|-----------------------|
| # DoFs                  | 5 276          | 20 532         | 80 996         | $3.2 \times 10^5$ | $1.3 \times 10^6$ | —                     |
| Period                  | 0.048231       | 0.048051       | 0.048031       | 0.048030          | 0.048029          | $0.04803 \pm 0.00003$ |
| $\text{Nu}^{\max}$      | 7.4065         | 7.3822         | 7.3789         | 7.3788            | 7.3788            | $7.379 \pm 0.005$     |
| $\text{Nu}^{\min}$      | 6.5062         | 6.4717         | 6.4691         | 6.4691            | 6.4692            | $6.468 \pm 0.005$     |
| $\text{Nu}^{\max}$      | 7.2637         | 7.2047         | 7.1969         | 7.1960            | 7.1960            | $7.196 \pm 0.005$     |
| $\text{Nu}^{\min}$      | 6.7878         | 6.7949         | 6.7961         | 6.7965            | 6.7966            | $6.796 \pm 0.005$     |
| $v_{\text{rms}}^{\max}$ | 60.726         | 60.398         | 60.361         | 60.359            | 60.360            | $60.367 \pm 0.015$    |
| $v_{\text{rms}}^{\min}$ | 31.829         | 31.965         | 31.981         | 31.981            | 31.982            | $31.981 \pm 0.02$     |
| $v_{\text{rms}}^{\max}$ | 58.225         | 57.517         | 57.442         | 57.437            | 57.436            | $57.43 \pm 0.05$      |
| $v_{\text{rms}}^{\min}$ | 30.392         | 30.330         | 30.324         | 30.323            | 30.322            | $30.32 \pm 0.03$      |


b)

| Finest mesh size        | $\frac{1}{64}$              | $\frac{1}{128}$             | $\frac{1}{256}$             | Reference             |
|-------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------|
| # DoFs                  | $4.5 \dots 6.0 \times 10^4$ | $1.6 \dots 2.2 \times 10^5$ | $5.6 \dots 8.0 \times 10^5$ | —                     |
| Period                  | 0.048029                    | 0.048030                    | 0.048030                    | $0.04803 \pm 0.00003$ |
| $\text{Nu}^{\max}$      | 7.3809                      | 7.3792                      | 7.3788                      | $7.379 \pm 0.005$     |
| $\text{Nu}^{\min}$      | 6.4718                      | 6.4695                      | 6.4691                      | $6.468 \pm 0.005$     |
| $\text{Nu}^{\max}$      | 7.1996                      | 7.1967                      | 7.1960                      | $7.196 \pm 0.005$     |
| $\text{Nu}^{\min}$      | 6.7986                      | 6.7969                      | 6.7965                      | $6.796 \pm 0.005$     |
| $v_{\text{rms}}^{\max}$ | 60.366                      | 60.361                      | 60.360                      | $60.367 \pm 0.015$    |
| $v_{\text{rms}}^{\min}$ | 31.980                      | 31.981                      | 31.981                      | $31.981 \pm 0.02$     |
| $v_{\text{rms}}^{\max}$ | 57.449                      | 57.434                      | 57.435                      | $57.43 \pm 0.05$      |
| $v_{\text{rms}}^{\min}$ | 30.322                      | 30.322                      | 30.322                      | $30.32 \pm 0.03$      |

a) Results for the 2-D benchmark problem with uniform mesh refinement. # DoFs indicates the number of degrees of freedom. Reference results from Blankenbach *et al.* (1989). b) Results with adaptive mesh refinement. The number of degrees of freedom (# DoFs) for each finest mesh size  $h$  varies between time steps; the indicated numbers provide a typical range.

|                                       |             | $V_{rms}$                | Nu                       |
|---------------------------------------|-------------|--------------------------|--------------------------|
| Blankenbach <i>et al.</i> (1989) [95] | $Ra = 10^4$ | $42.864947 \pm 0.000020$ | $4.884409 \pm 0.000010$  |
|                                       | $Ra = 10^5$ | $193.21454 \pm 0.00010$  | $10.534095 \pm 0.000010$ |
|                                       | $Ra = 10^6$ | $833.98977 \pm 0.00020$  | $21.972465 \pm 0.000020$ |
| Tackley (1994) [1227]                 | $Ra = 10^4$ | 42.775                   | 4.878                    |
|                                       | $Ra = 10^5$ | 193.11                   | 10.531                   |
|                                       | $Ra = 10^6$ | 833.55                   | 21.998                   |
| King (2009) [701]                     | $Ra = 10^4$ | 42.867                   | 4.885                    |
|                                       | $Ra = 10^5$ | 193.248                  | 10.536                   |
|                                       | $Ra = 10^6$ | 834.353                  | 21.981                   |
| Thieulot (2014) [1257]                | $Ra = 10^4$ | 42.867                   | 4.882                    |
|                                       | $Ra = 10^5$ | 193.255                  | 10.507                   |
|                                       | $Ra = 10^6$ | 834.712                  | 21.695                   |
| ASPECT [44]                           | $Ra = 10^4$ |                          |                          |
|                                       | $Ra = 10^5$ |                          |                          |
|                                       | $Ra = 10^6$ |                          |                          |

Steady state Nusselt number and  $V_{rms}$  measurements as reported in the literature and obtained with ELEFANT on a  $200 \times 200$  grid. King (2009) results on 200x200 grid with ConMan.

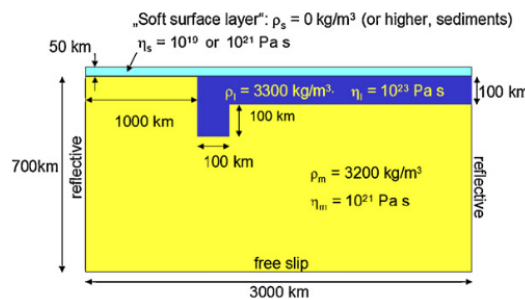
 **Relevant Literature:** Travis *et al.* [1278] (1990), Ogawa [953] (1993), Trompert and Hansen [1283] (1998), Auth and Harder [33] (1999), Christon, Gresho, and Sutton [256] (2002), Chiu-Webster, Hinch, and Lister [234] (2008), Kameyama, Kageyama, and Sato [668] (2005), King [701] (2009), Beuchert and Podladchikov [86] (2010), Leng and Zhong [769] (2011), Davies, Wilson, and Kramer [309] (2011), Vynnytska, Rognes, and Clark [1333] (2013), Trim, Butler, and Spiteri [1282] (2021), Davies, Kramer, Ghelichkhan, and Gibson [306] (2002), Sime and Wilson [1169] (2020), [STONE 3](#).

### 12.2.27 Subduction 'benchmark' of Schmeling *et al.* (2008)

benchmark\_scbe08.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/benchmark\\_scbe08](https://github.com/cedrict/fieldstone/tree/master/images/benchmark_scbe08)

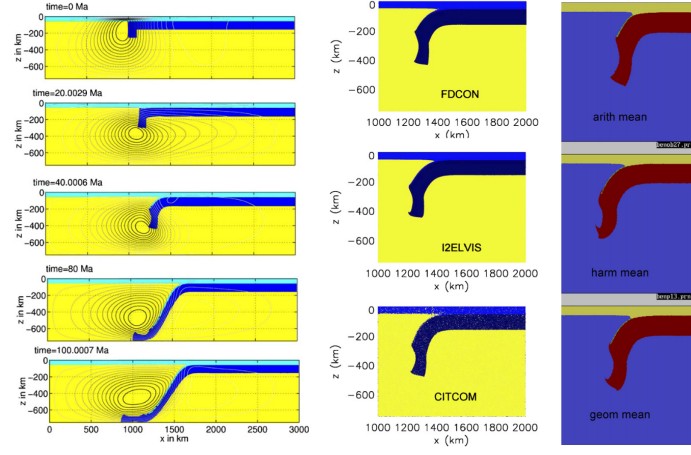
The setup is as follows :



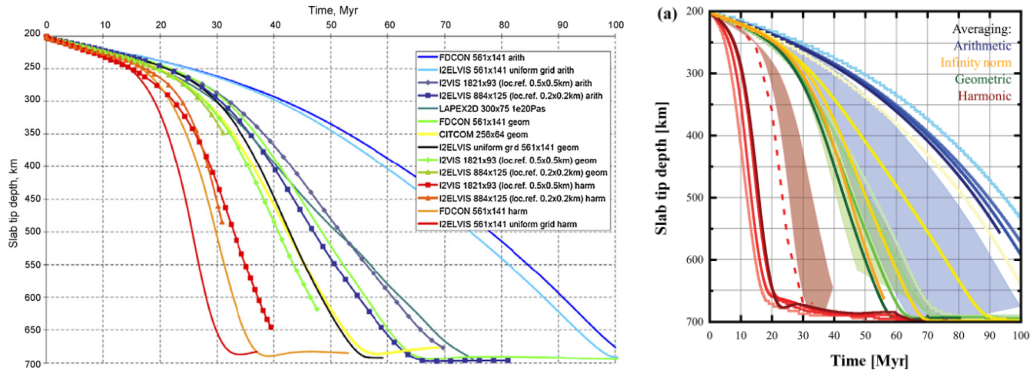
Taken from Schmeling *et al.* (2008) [1124].

Materials are linear viscous, initial geometry is simple, boundary conditions are simple. On paper it sounds like a good idea. See [STONE 67](#) for a discussion on why this setup was doomed from the beginning as a benchmark.

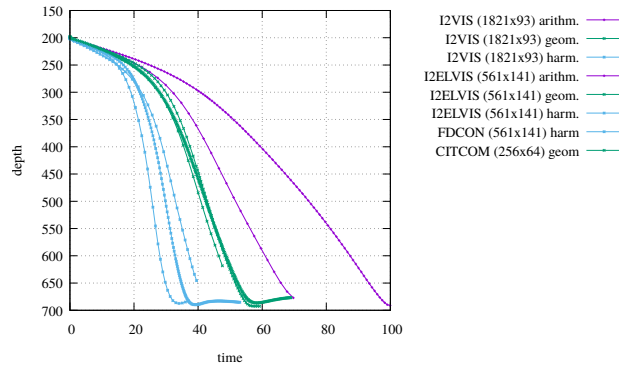
Experiments have been conducted by a handful of codes, investigating the effect of averaging and mesh resolution:



Taken from Schmeling *et al.* (2008) [1124]. Left: case 1 model (here FDCON-4 is shown). Streamlines are also shown. Middle: Shapes of different case 1 models at similar stages: FDCON: 40 Myears, I2ELVIS: 34.7 Myears, CITCOM: 38.1 Myears. Viscosity averaging: geometric mean in all cases. Right: Comparison of the shapes of the slabs for different viscosity averaging methods using I2VIS. Note that the snapshots are taken at different times (59.6, 24.4, 37.8 Myears from top to bottom), so that the slab tips have reached comparable levels.



Temporal behaviour of case 1 modelled by different codes with highest resolutions each. Each curve shows the position of the deepest part of the slab (slab tip) as a function of time below the initial surface of the lithosphere. Left: taken from Schmeling *et al.* (2008) [1124]; Right: taken from Glerum *et al.* (2018) [467].



Data courtesy of Prof. H. Schmeling.

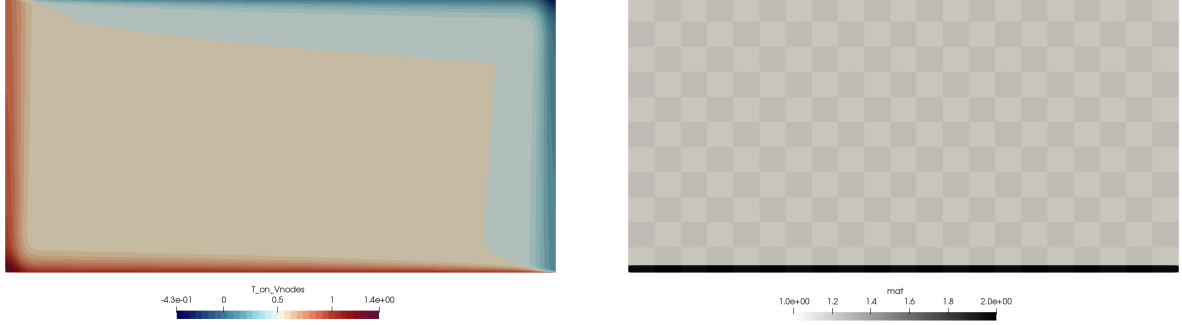
## 12.2.28 Thin layer entrainment

benchmark\_thin\_layer\_entrapment.tex

Data pertaining to this section are to be found at:  
[https://github.com/cedrict/fieldstone/tree/master/images/benchmark\\_thinlayer](https://github.com/cedrict/fieldstone/tree/master/images/benchmark_thinlayer)

The problem is a simulation to study the amount of entrainment by thermal convection of a dense, thin layer at the bottom of the model [1309]. To the author's knowledge only two other publications (Tackley & King (2003) [1229], van Thienen (2007) [1255]) have presented results pertaining to this benchmark. The results shown here after are obtained with my ELEFANT code using the particle-in-cell technique and originate in the ELEFANT paper [1257].

The box is  $2 \times 1$ , and contains two fluids:



Fluid 1 has a density  $\rho_1 = 1$  and a viscosity  $\eta = 1$ . Fluid 2 is heavier ( $\rho_2 = \rho_1 + \Delta\rho$ ) but has the same viscosity. Both fluids have a thermal expansion coefficient  $\alpha = 10^{-10}$ , a thermal conductivity  $k = 1$ , and a heat capacity coefficient  $c_p = 1$ . Fluid 2 is placed at the bottom of the box ( $0 \leq y \leq 0.025$ ).

This experiment is parameterised by the thermal Rayleigh number  $\text{Ra} = 300,000$  and the compositional Rayleigh number  $\text{Ra}_c = 450,000$  which are defined as follows:

$$\text{Ra}_T = \frac{\alpha \rho g \Delta T L_y^3}{\kappa \eta} = \frac{\alpha \rho^2 g \Delta T L_y^3 c_p}{k \eta} = \alpha g \quad (12.281)$$

$$\text{Ra}_c = \frac{\Delta \rho g L_y^3}{\kappa \eta} = \frac{\rho \Delta \rho g L_y^3 c_p}{k \eta} = \Delta \rho g \quad (12.282)$$

where I have used the relationship  $\kappa = k/\rho c_p$ .  $B$  is defined as  $B = \text{Ra}_T/\text{Ra}_c$  so The gravity acceleration is therefore set to  $g = \text{Ra}/\alpha$  and this yields  $\Delta\rho = \text{Ra}_c/g = B\text{Ra}_T/g = B \times \alpha$ .

Free-slip boundary conditions are imposed on all sides of the domain. Temperature boundary conditions are  $T(x, y = 0) = 1$  and  $T(x, y = 1) = 0$ . The analytical initial temperature field is given by

$$T(x, y) = T_u(x, y) + T_l(x, y) + T_r(x, y) + T_s(x, y) - \frac{3}{2} \quad (12.283)$$

where

$$\begin{aligned} T_u(x, y) &= \frac{1}{2} \text{erf} \left( \frac{1-y}{2} \sqrt{\frac{u_0}{x}} \right) \\ T_l(x, y) &= 1 - \frac{1}{2} \text{erf} \left( \frac{y}{2} \sqrt{\frac{u_0}{L_x - x}} \right) \\ T_r(x, y) &= \frac{1}{2} + \frac{Q}{2\sqrt{\pi}} \sqrt{\frac{u_0}{y+1}} \exp \left( -\frac{x^2 u_0}{4y+4} \right) \\ T_s(x, y) &= \frac{1}{2} - \frac{Q}{2\sqrt{\pi}} \sqrt{\frac{u_0}{2-y}} \exp \left( -\frac{(L_x - x)^2 u_0}{8 - 4y} \right) \end{aligned} \quad (12.284)$$

with

$$u_0 = \frac{L_x^{7/3}}{(1 + L_x^4)^{2/3}} \left( \frac{\text{Ra}}{2\sqrt{\pi}} \right)^{2/3} \quad Q = 2\sqrt{\frac{L_x}{\pi u_0}} \quad (12.285)$$

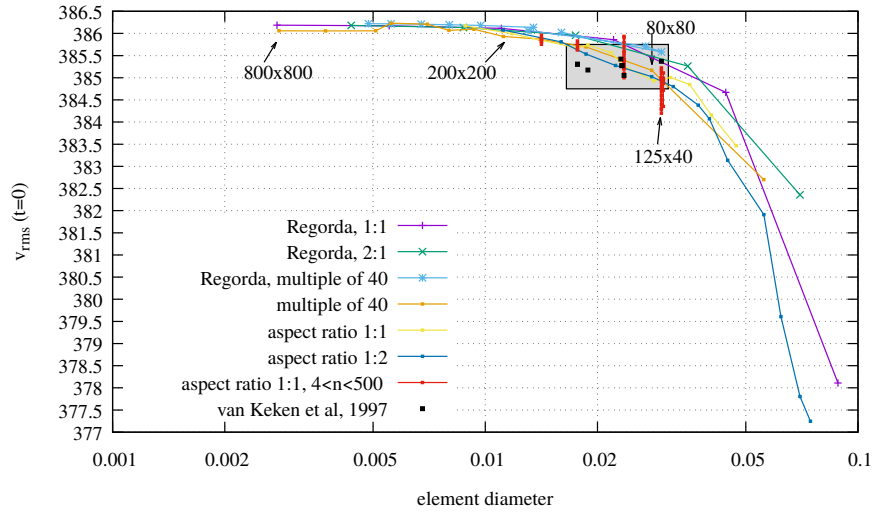


Using  $L_x = 2$ ,  $Ra = 3 \times 10^5$ , one gets  $u_0 \simeq 1469.315$  and  $Q \simeq 0.0416305$ .

Given the small thickness of the bottom layer, it seems quite legitimate to investigate the influence of grid resolution on the simulation. I have therefore looked at the initial root mean square velocity measurement as a function of the element diagonal value (a proxy for the average resolution in the case where elements are not square).

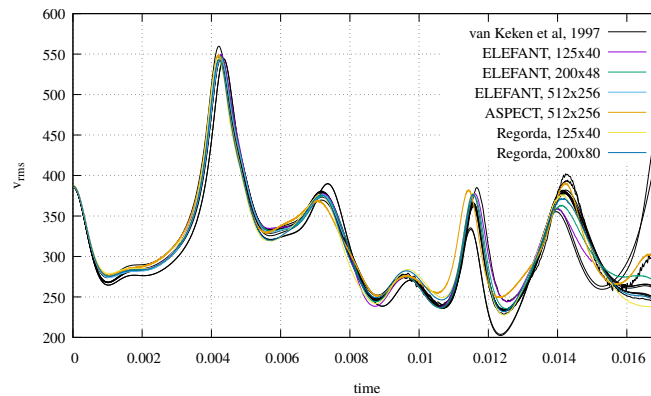
Results are confirm that the element size plays a non negligible role at startup on the dynamics of the system. Superimposed on the figure are the measurements provided by Prof. van Keken (black squares in the gray box). They agree well with my measurements but also indicate that none of the authors in the original study ran the experiment at a high-enough resolution to start with (their results were therefore most likely resolution dependent).

We see that the number of markers per element at startup is critical at (very) low resolution but that it does not lead to significant velocity variations at high resolution.



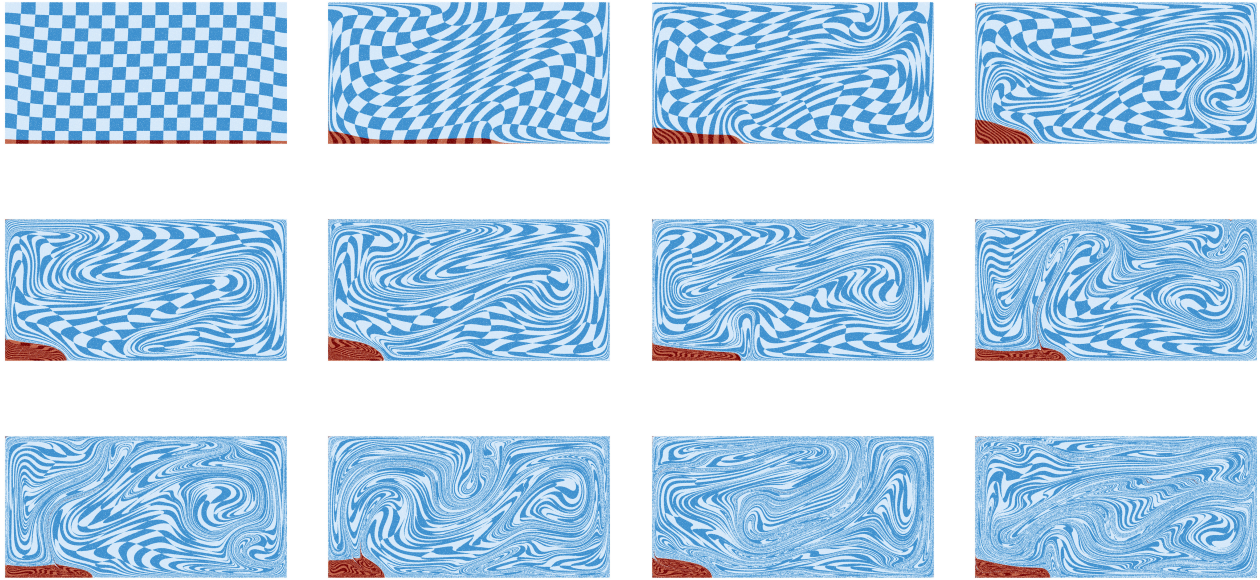
Thin layer entrainment experiment: root mean square velocity measurements at  $t = 0$  as a function of the element diagonal size. The red square points correspond to resolutions where the number of elements in each direction is a multiple of 40 (i.e.  $L_y/d$ ), so that no element would contain a mix of fluids 1 and 2. Pink points correspond to cases where the number of markers within each element was varied between 4 and 500 (random spatial distribution). Taken from ELEFANT paper [1257]

Looking at the root mean square velocity measurements, we see that the measurements done with ELEFANT agree nicely with those presented in van Keken *et al.* [1309]. Past  $t \sim 0.015$ , the curves diverge clearly across all codes and authors, so I only need to focus the comparison for times  $t < 0.015$ . For the three tested resolutions, measurements agree well and fall within the grey curves representing all results of van Keken *et al.*. Additional tests have been carried out concerning the value of the Courant number (0.1 to 0.25) and the initial number of markers per element (100 or 200) and these parameters led to extremely similar results.



Thin layer entrainment experiment. Root mean square velocity as a function of time. All results presented in van Keken *et al.* (1997) are collapsed in black dashed lines. All simulations were run with an initial marker density of 100 markers per element and with a Courant number of 0.25. Taken from ELEFANT paper [1257].

As observed in van Keken *et al.* , the dense layer is first swept into the lower left corner. Thermal instabilities then further develop in an asymmetrical way and entrain the dense material. Past  $t \simeq 0.015$  the system becomes more and more chaotic with markers being randomly mixed in the system in a non-orderly fashion.



Marker distribution as obtained with ELEFANT for grid 240x120, init\_marker\_density=7, random distribution, CFL=0.25, rkmethod=2, m.to.q=2. (unpublished).

 Relevant Literature: Trim *et al.* (2020) [1281].



# Chapter 13

## Vorticity-stream function approach

chapter\_streamfunction.tex

The Stream function (commonly denoted by  $\Psi$ ) approach is a useful approach in fluid dynamics as it can provide relatively quick solutions to 2D incompressible flow problems. Using a stream function formulation is numerically convenient because velocity information is contained in a single scalar equation and pressure vanishes from the solution process.

The stream function is a function of coordinates and time of an inviscid liquid. It allows to determine the components of velocity by differentiating the stream function with respect to the space coordinates. A family of curves  $\Psi = \text{constant}$  represent **streamlines**, i.e. the stream function remains constant along a streamline. Although also valid in 3D, this approach is mostly used in 2D because of its relative simplicity.

Glaisner and Tezduyar [464] state:

“The main advantages of the vorticity stream-function formulation are the simple form of the equations in two-dimensions and the in-built satisfaction of the incompressibility constraint.

In two dimensions, the vorticity transport equation and the Poisson equation for the stream function are scalar and there are only two degree of freedom in the problem. Moreover, the vorticity stream-function form of the Navier-Stokes equations allows equal order of interpolation for the vorticity and the stream function. In fact the bilinear interpolations which have been used here for both the unknowns are sufficient which is an asset from the point of view of implementation.

On the other hand, a flow field obtained by the solution of the vorticity stream-function equations is by definition divergence-free and the initial condition of the problem do not need to satisfy the Incompressibility constraint. This allows the use of an initial flow field noncontinuous at the boundary of the domain. ”



|                 |                   |               |
|-----------------|-------------------|---------------|
| Stream function | $\leftrightarrow$ | Stroomfunctie |
| Vorticity       | $\leftrightarrow$ | vorticiteit   |

### 13.1 Vorticity-stream function formulation of the isoviscous Navier-Stokes equation

What follows is adapted from Glaisner and Tezduyar [464] (1987). The vorticity transport equation can be obtained by taking the curl of the (isoviscous) incompressible momentum equation,

$$\vec{\nabla} \times \left[ \frac{\partial \vec{v}}{\partial t} + (\vec{v} \cdot \vec{\nabla}) \vec{v} \right] = \vec{\nabla} \times \left[ -\frac{1}{\rho} \vec{\nabla} p + \nu \vec{\nabla}^2 \vec{v} + \vec{g} \right] \quad (13.1)$$

with  $\nu = \eta/\rho$ . **i need to propagate gravity in what follows** Computing term-by-term, we have

$$\vec{\nabla} \times \frac{\partial \vec{v}}{\partial t} = \frac{\partial}{\partial t}(\vec{\nabla} \times \vec{v}) = \frac{\partial \vec{\omega}}{\partial t}$$

where the vorticity vector  $\vec{\omega}$  is defined as

$$\vec{\omega} = \vec{\nabla} \times \vec{v} = \begin{vmatrix} \partial_x & u \\ \partial_y & v \\ \partial_z & w \end{vmatrix} \quad (13.2)$$

Because the curl of a gradient is zero (**assuming rho constant?!**),

$$\vec{\nabla} \times \left( -\frac{1}{\rho} \vec{\nabla} p \right) = 0$$

Also, (**assuming nu constant?!**)

$$\vec{\nabla} \times (\nu \vec{\nabla} |\vec{v}|^2) = \nu \vec{\nabla}^2 (\vec{\nabla} \times \vec{v}) = \nu \vec{\nabla}^2 \vec{\omega}$$

Now, we use the vector identity

$$(\vec{v} \cdot \vec{\nabla}) \vec{v} = \frac{1}{2} \vec{\nabla} |\vec{v}|^2 - \vec{v} \times (\vec{\nabla} \times \vec{v}) = \frac{1}{2} \vec{\nabla} |\vec{v}|^2 - \vec{v} \times \vec{\omega} \quad (13.3)$$

Taking the curl of both sides and making use of the curl of a gradient equals zero and  $\vec{\omega} = \vec{\nabla} \times \vec{v}$ , results in

$$\vec{\nabla} \times [(\vec{v} \cdot \vec{\nabla}) \vec{v}] = -\vec{\nabla} \times (\vec{v} \times \vec{\omega}) = \vec{\nabla} \times (\vec{\omega} \times \vec{v})$$

Combining all the above terms, we have thus obtained the vorticity equation

$$\boxed{\frac{\partial \vec{\omega}}{\partial t} + \vec{\nabla} \times (\vec{\omega} \times \vec{v}) = \nu \vec{\nabla}^2 \vec{\omega}} \quad (13.4)$$

Also, we have the identity

$$\vec{\nabla} \times (\vec{\omega} \times \vec{v}) = (\vec{v} \cdot \vec{\nabla}) \vec{\omega} - (\vec{\omega} \cdot \vec{\nabla}) \vec{v}$$

so that in the end:

$$\boxed{\frac{\partial \vec{\omega}}{\partial t} + (\vec{v} \cdot \vec{\nabla}) \vec{\omega} = (\vec{\omega} \cdot \vec{\nabla}) \vec{v} + \nu \vec{\nabla}^2 \vec{\omega}} \quad (13.5)$$

In the case of a two-dimensional flow, Eq. (13.5) simplifies to<sup>1</sup>

$$\frac{\partial \vec{\omega}}{\partial t} + \vec{v} \cdot \vec{\nabla} \vec{\omega} = \nu \vec{\nabla}^2 \vec{\omega}$$

The velocity  $\vec{v}$  being solenoidal, there exists a vector field  $\vec{\Psi} = (0, 0, \Psi)$  such that

$$\vec{v} = \vec{\nabla} \times \vec{\Psi}$$

The nonzero component of  $\Psi$  is called the stream function. A relation between the vorticity and the stream function is obtained as

$$\vec{\omega} = \vec{\nabla} \times \vec{\nabla} \times \vec{\Psi} = -\vec{\nabla}^2 \vec{\Psi}$$

---

<sup>1</sup>In the case of a two-dimensional flow, the vorticity-stretching term  $(\vec{v} \cdot \vec{\nabla}) \vec{v}$  is zero and the vorticity transport equation contains only one nonlinear term. This, together with the scalar nature of the equations, is the reason why the vorticity stream-function formulation is very convenient in two dimensions.

verify  $\text{rot rot} = -\text{laplace}$

For 2D flows in the  $x - y$  plane, since  $w = 0$  and  $\partial_z = 0$  then it has only one non-zero component  $\vec{\omega} = \omega \vec{e}_z$  with

$$\omega = \frac{\partial v}{\partial x} - \frac{\partial u}{\partial y} = \frac{\partial}{\partial x} \left( -\frac{\partial \Psi}{\partial x} \right) - \frac{\partial}{\partial y} \left( \frac{\partial \Psi}{\partial y} \right) = -\frac{\partial^2 \Psi}{\partial x^2} - \frac{\partial^2 \Psi}{\partial y^2} = -\vec{\nabla}^2 \Psi$$

The vorticity stream-function form of the Navier-Stokes equations for two-dimensional flows is the set of coupled scalar equations given below :

$$\frac{\partial \omega}{\partial t} + \vec{v} \cdot \vec{\nabla} \omega = \nu \vec{\nabla}^2 \omega \quad \text{on } \Omega \times [0, T] \quad (13.6)$$

$$\vec{\nabla}^2 \Psi = -\omega \quad \text{on } \Omega \times [0, T] \quad (13.7)$$

where  $\Omega$  is a domain of  $\mathbb{R}^2$ ,  $\vec{v}(\vec{x}, t)$  is the velocity,  $\omega(\vec{x}, t)$  is the vorticity,  $\Psi(\vec{x}, t)$  is the stream function, and  $\nu$  is the kinematic viscosity. The vorticity and the stream function are related to the velocity field through

$$\omega = \frac{\partial v}{\partial x} - \frac{\partial u}{\partial y} \quad (13.8)$$

$$u = \frac{\partial \Psi}{\partial y} \quad (13.9)$$

$$v = -\frac{\partial \Psi}{\partial x} \quad (13.10)$$

For flow domains having solid boundaries ( e.g., walls or obstacles ) the no-slip and no-penetration conditions yield constraints on the derivatives of  $\Psi$  , i.e,

$$\vec{n} \cdot \vec{\nabla} \Psi = v_t \quad (13.11)$$

$$\vec{t} \cdot \vec{\nabla} \Psi = 0 \quad (13.12)$$

where  $\vec{t}$  and  $\vec{n}$  are the tangential and normal unit vectors at the surface whereas  $v_t$  is the tangential component of the velocity at the wall. Moreover, the stream function assumes constant values along solid boundaries. However, it is difficult to derive any boundary condition for vorticity at solid surfaces. The boundary values of the vorticity at solid surfaces are related to the local boundary-layer profiles and are thus time dependent. On the other hand, they do not arise naturally from the no-slip and no-penetration conditions.

## 13.2 Vorticity-stream function formulation of the isoviscous Stokes equation in 2d

In two dimensions the velocity is obtained as follows:

$$\vec{v} = (u, v) = \left( \frac{\partial \Psi}{\partial y}, -\frac{\partial \Psi}{\partial x} \right) \quad (13.13)$$

Provided the function  $\Psi$  is a smooth enough function, this automatically insures that the flow is incompressible:

$$\vec{\nabla} \cdot \vec{v} = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = \frac{\partial^2 \Psi}{\partial x \partial y} - \frac{\partial^2 \Psi}{\partial x \partial y} = 0 \quad (13.14)$$

Also, we have

$$\vec{v} \cdot \vec{\nabla} \Psi = \vec{v} \cdot \left( \frac{\partial \Psi}{\partial x}, \frac{\partial \Psi}{\partial y} \right) = (u, v) \cdot (-v, u) = 0$$

Assuming constant viscosity, the Stokes equation writes:

$$-\vec{\nabla} p + \eta_0 \Delta \vec{v} + \rho \vec{g} = \vec{0} \quad (13.15)$$

or, in each dimension:

$$0 = -\frac{\partial p}{\partial x} + \eta_0 \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial x^2} \right) + \rho g_x \quad (13.16)$$

$$0 = -\frac{\partial p}{\partial y} + \eta_0 \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) + \rho g_y \quad (13.17)$$

Using Eq. (13.13), we can write

$$\begin{aligned} 0 &= -\frac{\partial p}{\partial x} + \eta_0 \left( \frac{\partial^3 \Psi}{\partial y^3} + \frac{\partial^3 \Psi}{\partial x^2 \partial y} \right) + \rho g_x \\ 0 &= -\frac{\partial p}{\partial y} - \eta_0 \left( \frac{\partial^3 \Psi}{\partial x^3} + \frac{\partial^3 \Psi}{\partial y^2 \partial x} \right) + \rho g_y \end{aligned} \quad (13.18)$$

The pressure terms in both equations can be removed by first differentiating the first line with regards to  $y$  and the second line with regards to  $x$ ,

$$\begin{aligned} 0 &= -\frac{\partial^2 p}{\partial x \partial y} + \eta_0 \left( \frac{\partial^4 \Psi}{\partial y^4} + \frac{\partial^4 \Psi}{\partial x^2 \partial y^2} \right) + \frac{\partial(\rho g_x)}{\partial y} \\ 0 &= -\frac{\partial^2 p}{\partial y \partial x} - \eta_0 \left( \frac{\partial^4 \Psi}{\partial x^4} + \frac{\partial^4 \Psi}{\partial y^2 \partial x^2} \right) + \frac{\partial(\rho g_y)}{\partial x} \end{aligned} \quad (13.19)$$

and next by subtracting the resulting equations, leading to:

$$0 = \eta_0 \left( \frac{\partial^4 \Psi}{\partial x^4} + 2 \frac{\partial^4 \Psi}{\partial x^2 \partial y^2} + \frac{\partial^4 \Psi}{\partial y^4} \right) + \frac{\partial(\rho g_x)}{\partial y} - \frac{\partial(\rho g_y)}{\partial x} \quad (13.20)$$

or

$$\vec{\nabla}^4 \Psi = \frac{1}{\eta_0} \left( -\frac{\partial(\rho g_x)}{\partial y} + \frac{\partial(\rho g_y)}{\partial x} \right) \quad (13.21)$$

where

$$\vec{\nabla}^4 = \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right)$$

which we find for example in According to Gerya's book2, page 70.

Note that  $\vec{\nabla}^2 \vec{\nabla}^2 = \vec{\nabla}^4$  is known as the **Biharmonic operator**. These equations are also to be found in the geodynamics literature, see Ismail-Zadeh and Tackley [626, eq. 1.43] or Gerya [455, p 70-71].

From  $\vec{\nabla}^2 \Psi = -\omega$ , we can also write

$$\vec{\nabla}^2 \omega = -\frac{1}{\eta_0} \left( -\frac{\partial(\rho g_x)}{\partial y} + \frac{\partial(\rho g_y)}{\partial x} \right) \quad (13.22)$$

### 13.3 Vorticity-stream function formulation of the non-isoviscous Stokes equation in 2d

In this case, we can no longer write the Stokes equation as follows

$$-\vec{\nabla} p + \eta_0 \Delta \vec{v} + \rho \vec{g} = \vec{0} \quad (13.23)$$

Indeed, it should be

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\epsilon}) + \rho \vec{g} = \vec{0} \quad (13.24)$$

We have

$$2\eta \dot{\epsilon} = 2\eta \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{1}{2}(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}) \\ \frac{1}{2}(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}) & \frac{\partial v}{\partial y} \end{pmatrix}$$

so that

$$W_x = -\frac{\partial p}{\partial x} + \frac{\partial}{\partial x} \left( 2\eta \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \eta \left( \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) \right) + \rho g_x \quad (13.25)$$

$$W_y = -\frac{\partial p}{\partial y} + \frac{\partial}{\partial x} \left( \eta \left( \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) \right) + \frac{\partial}{\partial y} \left( 2\eta \frac{\partial v}{\partial y} \right) + \rho g_y \quad (13.26)$$

Using Eq. (13.13) again, we can write

$$W_x = -\frac{\partial p}{\partial x} + \frac{\partial}{\partial x} \left( 2\eta \frac{\partial^2 \Psi}{\partial x \partial y} \right) + \frac{\partial}{\partial y} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \rho g_x = 0 \quad (13.27)$$

$$W_y = -\frac{\partial p}{\partial y} + \frac{\partial}{\partial x} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) - \frac{\partial}{\partial y} \left( 2\eta \frac{\partial^2 \Psi}{\partial y \partial x} \right) + \rho g_y = 0 \quad (13.28)$$

The pressure terms in both equations can be removed by first differentiating the first line with regards to  $y$  and the second line with regards to  $x$ ,

$$\frac{\partial W_x}{\partial y} = -\frac{\partial^2 p}{\partial x \partial y} + \frac{\partial^2}{\partial x \partial y} \left( 2\eta \frac{\partial^2 \Psi}{\partial x \partial y} \right) + \frac{\partial^2}{\partial y^2} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \frac{\partial \rho g_x}{\partial y} \quad (13.29)$$

$$\frac{\partial W_y}{\partial x} = -\frac{\partial^2 p}{\partial y \partial x} + \frac{\partial^2}{\partial x^2} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) - \frac{\partial^2}{\partial x \partial y} \left( 2\eta \frac{\partial^2 \Psi}{\partial y \partial x} \right) + \frac{\partial \rho g_y}{\partial x} \quad (13.30)$$

and next by subtracting the resulting equations, leading to:

$$\begin{aligned} 0 &= \frac{\partial^2}{\partial x \partial y} \left( 2\eta \frac{\partial^2 \Psi}{\partial x \partial y} \right) + \frac{\partial^2}{\partial y^2} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) - \frac{\partial^2}{\partial x^2} \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \frac{\partial^2}{\partial x \partial y} \left( 2\eta \frac{\partial^2 \Psi}{\partial y \partial x} \right) + \frac{\partial \rho g_x}{\partial y} - \frac{\partial \rho g_y}{\partial x} \\ 0 &= \frac{\partial^2}{\partial x \partial y} \left( 2\eta \left( \frac{\partial^2 \Psi}{\partial x \partial y} + \frac{\partial^2 \Psi}{\partial y \partial x} \right) \right) + \left( \frac{\partial^2}{\partial y^2} - \frac{\partial^2}{\partial x^2} \right) \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \frac{\partial \rho g_x}{\partial y} - \frac{\partial \rho g_y}{\partial x} \\ 0 &= 4 \frac{\partial^2}{\partial x \partial y} \left( \eta \frac{\partial^2 \Psi}{\partial x \partial y} \right) + \left( \frac{\partial^2}{\partial y^2} - \frac{\partial^2}{\partial x^2} \right) \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \frac{\partial \rho g_x}{\partial y} - \frac{\partial \rho g_y}{\partial x} \end{aligned} \quad (13.31)$$

i.e.

$$\boxed{4 \frac{\partial^2}{\partial x \partial y} \left[ \eta \frac{\partial^2 \Psi}{\partial x \partial y} \right] + \left( \frac{\partial^2}{\partial y^2} - \frac{\partial^2}{\partial x^2} \right) \left[ \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right] = -\frac{\partial \rho g_x}{\partial y} + \frac{\partial \rho g_y}{\partial x}}$$

This is eq 5.36c of Gerya. This expression (esp. the lhs) is to be found in its dimensionless form in Christensen and Yuen [245] (1984) or in Schmeling and Jacoby [1123] (1981) for example:

$$4 \frac{\partial^2}{\partial x \partial y} \left( \eta \frac{\partial^2 \Psi}{\partial x \partial y} \right) + \left( \frac{\partial^2}{\partial y^2} - \frac{\partial^2}{\partial x^2} \right) \left( \eta \left( \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \right) \right) + \frac{\partial \rho g_x}{\partial y} - \frac{\partial \rho g_y}{\partial x} = Ra \frac{\partial T}{\partial x} - Rb \frac{\partial \Gamma}{\partial x}$$

In the presence of temperature variations and multiple compositions, Trim, Lowman, and Butler [1281] (2020) use the following identical nondimensional equation:

$$\left( \frac{\partial^2}{\partial x^2} - \frac{\partial^2}{\partial y^2} \right) \left[ \eta \left( \frac{\partial^2 \Psi}{\partial x^2} - \frac{\partial^2 \Psi}{\partial y^2} \right) \right] + 4 \frac{\partial^2}{\partial x y} \left[ \eta \frac{\partial^2 \Psi}{\partial x y} \right] = Ra_T \frac{\partial T}{\partial x} - Ra_C \frac{\partial C}{\partial x}$$

PB: what if C is a discontinuous field? then I guess C needs to be carried on the nodes.

Introducing the vorticity again, we can rewrite the above equation as follows: ??? can it ???

## 13.4 Boundary conditions

When we use a velocity-pressure formulation in a 2D Cartesian domain Dirichlet boundary conditions translate into zeroing either  $u$ ,  $v$ , or both, or assigning  $u$  and/or  $v$  a value on (part of) the boundary. In this case our unknowns are the  $\Psi$  function and the vorticity  $\omega$  so we need to think about boundary conditions a bit more.

The solution of vorticity transport equation and stream function equation requires that appropriate vorticity and stream function boundary conditions are specified at the boundaries. The specification of these boundary conditions is extremely important since it directly affects the stability and accuracy of the solution. However, neither vorticity nor its derivatives at the boundary are usually known in advance. Therefore a set of boundary conditions must be constructed.

For example, since the flow is parallel to a solid boundary, solid boundaries and symmetry planes are surfaces of constant stream function. In other words, since flow is parallel to the walls of the cavity, walls may be treated as streamline. Thus, the stream function value on the wall streamline is set as a constant (often taken to be zero for simplicity). Since the stream function is a constant along a wall, all the derivatives of stream function along the wall also vanish.

In the end we have to account for three physical types of boundary conditions: free slip, no slip, and prescribed velocity (yes, technically the last one encompasses the second one).

Check Napolitano, Pascazio, and Quartapelle [929] (1999) for a review of vorticity conditions in the numerical solution of the vorticity-stream function equations.

### 13.4.1 Free slip, 'stress free' boundary conditions

The free slip boundary condition states that at the interface between a moving fluid and a stationary wall, 1) the normal component of the fluid velocity field is equal to zero, but the tangential component is unrestricted. 2) the tangential stress is set to zero (which is why it is often called 'stress free' b.c.).

Free slip at the top and bottom then implies:

1.  $v = 0$  which translates into  $\frac{\partial \Psi}{\partial x} = 0$ . Actually stating that  $\frac{\partial \Psi}{\partial x} = 0$  along an horizontal edge means that  $\Psi$  is then constant along that edge.
2.  $\tau_{xy} = 2\eta\dot{\epsilon}_{xy} = \eta(\partial_y u + \partial_x v) = 0$  on the boundary. Since  $v = 0$ , then it does not change in the  $x$  direction on the boundary, i.e.  $\partial_x v = 0$ , leaving the condition  $\partial_y u = 0$  which translates into  $\frac{\partial^2 \Psi}{\partial y^2} = 0$ . Since  $\Psi$  is constant on that edge (see point above) then we also have  $\frac{\partial^2 \Psi}{\partial x^2} = 0$  so that we can write that  $\omega = -\vec{\nabla}^2 \Psi = 0$  on that boundary.

The same reasoning on the sides means:  $u = 0$  and  $\partial_x v = 0$  which translates into  $\frac{\partial \Psi}{\partial y} = 0$  and  $\frac{\partial^2 \Psi}{\partial x^2} = 0$ . In general the Poisson equation  $\vec{\nabla}^2 \Psi = -\omega$  becomes on a wall:

$$\left. \frac{\partial^2 \Psi}{\partial n^2} \right|_{wall} = -\omega_{wall}$$

where  $n$  is the normal direction. Indeed, in Hsui [597] (1978) we read:

“Besides the initial condition, boundary conditions are also required in order to form a well-posed system. Equation 13 [ i.e  $\omega = \text{curl } \Psi$ ] requires the specification of stream function at all boundaries. Since only the derivatives of stream functions which determine the velocities are important to the problem, stream function can thus be specified to within a constant. Consequently, they are chosen, for convenience, to be zero at all boundaries.”

This is also coherent with for example:

$$\begin{array}{llll}
\psi = 0 & \partial^2 \psi / \partial z^2 = 0 & \theta = 0 & \text{on } z = 1 \\
\psi = 0 & \partial^2 \psi / \partial z^2 = 0 & \partial \theta / \partial z = 1 & \text{on } z = 0 \\
\psi = 0 & \partial^2 \psi / \partial x^2 = 0 & \partial \theta / \partial x = 0 & \text{on } x = 0, l
\end{array}$$

(“free stress b.c.”) Richter and McKenzie [1073]

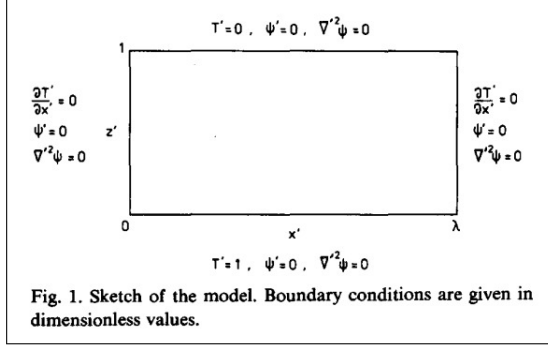
(1983)

### 2.3 Boundary conditions

The boundaries of the calculation region are free-slip and stress-free; therefore boundary conditions of this situation are expressed as follows,

- (1)  $\psi' = 0$  at all boundaries,
- (2)  $\partial^2 \psi' / \partial x'^2 = 0$  at both side boundaries,
- (3)  $\partial^2 \psi' / \partial z'^2 = 0$  at top and bottom boundaries.

Matsumoto and Tomoda [841]



Hansen and Ebel [533]

For all experiments we used boundary conditions according to

$$\begin{array}{llll}
z = 0, & 0 < x < \lambda: \psi = 0; & \nabla^2 \psi = 0, & T = 1. \\
z = 1, & 0 < x < \lambda: \psi = 0; & \nabla^2 \psi = 0, & T = 0. \\
0 < z < 1, & x = 0: \psi = 0; & \nabla^2 \psi = 0, & \frac{\partial T}{\partial x} = 0. \\
0 < z < 1, & x = \lambda: \psi = 0; & \nabla^2 \psi = 0, & \frac{\partial T}{\partial x} = 0.
\end{array} \quad (5)$$

Hansen and Ebel [540]

The boundary conditions are stress-free and impermeable at the top and bottom boundaries, and reflecting along the sides. The dimensionless temperature,  $T$ , is 0 at the top and 1 at the bottom. The boundary conditions are

$$\omega(x, z) = \psi(x, z) = 0 \begin{cases} \text{for } x = 0 \text{ and } x = x_0, \\ \text{for } z = 0 \text{ and } z = 1, \end{cases} \quad (4a)$$

$$T(x, z) = \begin{cases} 1 & \text{for } z = 0, \\ 0 & \text{for } z = 1, \end{cases} \quad (4b)$$

$$\frac{\partial T}{\partial x} = 0 \text{ at } x = 0 \text{ and } x = x_0, \quad (4c)$$

where  $x_0$  defines the aspect-ratio of the box.

Larsen, Yuen, Moser, and Fornberg [751]

## 13.4.2 No slip

Rather counter intuitively no-slip boundary conditions prove to be substantially more difficult to implement than free-slip boundary conditions. For example we find in DeMarco, DeAndrade, and Zapparoli [327] (2003):

“According to Layton [755] (1999), there are at least three natural ways of imposing zero tangential velocities value along the boundary: (i) Lagrange multiplier of tangential velocity component equal to zero as a constraint; (ii) Penalty term imposing tangential velocity component equal to zero approximately; (iii) Replacing no-slip with slip with friction. ”

DeMarco, DeAndrade, and Zapparoli [327] (2023) present “an analysis of penalty method application to imposing no-slip boundary condition. The approach used herein consists to express the vorticity boundary condition through the natural boundary condition depending on the solid wall tangential velocity component.”

Let us consider the bottom boundary. No-slip means  $u = v = 0$ , i.e.  $\frac{\partial \Psi}{\partial y} = \frac{\partial \Psi}{\partial x} = 0$ . We have seen previously that  $v = 0$  yields  $\Psi = \text{constant}$  on the boundary, leaving  $\frac{\partial \Psi}{\partial y} = 0$ . This is problematic because if one solves the biharmonic equation as two Poisson equations for  $\omega$  and  $\Psi$  it does not translate to a boundary condition for  $\omega$ !

In general zeroing the tangential component of the velocity on a boundary will write  $\partial \Psi / \partial n = 0$  where  $n$  is the normal to the boundary.

$$\left(\frac{\partial^2}{\partial z^2} - \frac{\partial^2}{\partial x^2}\right) [\eta(\psi_{zz} - \psi_{xx})] + 4 \frac{\partial^2}{\partial x \partial z} (\eta \psi_{xz}) = \text{Ra} \frac{\partial T}{\partial x}$$

$\psi = 0$  and  $\psi_{xx} = 0$  at  $x = 0, x = l$ ;  
 $\psi = 0$  and  $\psi_{zz} = 0$  (free slip) or  $\psi_z = 0$  (no slip) at  $z = 0, z = 1$ .

Christensen [243] (1984)

For a rectangular domain (the box containing the two fluids) the boundary conditions are:

$$\begin{aligned} \Psi = 0, \quad \frac{\partial^2 \Psi}{\partial n^2} = 0 & \quad \text{free slip} \\ \Psi = 0, \quad \frac{\partial \Psi}{\partial n} = 0 & \quad \text{fixed} \end{aligned} \quad (5b)$$

Woidt [1366] (1978)

Here  $u$  is prescribed at the top, but not zero:

The boundary conditions of the model are

All boundaries:  $\Psi = 0$

Sides:  $\omega = 0$   
 $\partial T / \partial x = 0$

Top:  $\partial \Psi / \partial y = u_D U^U(x, t)$  (8)  
 $T = 0$

Bottom:  $\omega = 0$   
 $T = 1$ ,  
or  $\partial T / \partial y = 0$

where  $\omega$  is the vorticity and the stream function,  $\Psi$ , is defined as

$$u_x = \frac{\partial \Psi}{\partial y} \quad u_y = -\frac{\partial \Psi}{\partial x} \quad (9)$$

Gurnis and Davies [514]

In Lux, Davies, and Thomas [817] (1979) we find this great summary of the problem and the sketch of its solution:

“The numerical methods we have used are well established and easy to apply. If we introduce the vorticity  $\omega$ , then the fourth-order stream function equation can be separated into a pair of second-order Poisson equations,

$$\omega = -\vec{\nabla}^2 \Psi \quad \vec{\nabla}^2 \omega = -\frac{\partial T}{\partial x},$$

for which efficient matrix reduction methods are available. For prescribed stress boundary conditions, the Poisson equations can simply be solved successively as done by McKenzie, Roberts, and Weiss [857] (1974). A complication arises for a prescribed boundary velocity because the boundary conditions on vorticity are then not explicit, see Richter [1071] (1973). On the boundary the vorticity is proportional to the shear stress and, since the shear stress is not known until the flow field is found, we lack a boundary condition on  $\omega$ . The Poisson equation for the stream function has an over-prescribed boundary condition, since both  $\Psi$ , and  $\partial \Psi / \partial y$  are given. This dilemma is resolved by using the condition on  $\Psi$  for the stream function equation and the condition on  $\partial \Psi / \partial y$  in a Taylor series expansion for  $\omega$  at the boundary for the vorticity equation. To ensure numerical stability with this boundary condition, it is necessary to iterate between the equations; an efficient way of optimizing this iteration has been given by Ehrlich and Gupta [363] (1975). The iteration technique for the coupled Poisson equations was tested against analytical solutions of the biharmonic equation given by Davies [314] (1977) and shown to be very accurate. The solution of the convection problem described here evolves with time, so the boundary condition on  $\omega$  must be found at each time step by the iteration method. This procedure is time consuming, so to increase the rate of convergence of this iteration a boundary condition on  $\omega$  accurate to first order (Taylor series) is used.”

In Comini, Manzan, and Nonino [275] (1994) we read:



In the analysis of two-dimensional incompressible flows the streamfunction-vorticity formulation of the Navier-Stokes equations allows the elimination of pressure from the problem and automatically satisfies the continuity constraint. On the other hand, the value of the vorticity at no-slip boundaries is difficult to specify and a poor evaluation of this boundary condition leads, almost invariably, to serious difficulties in obtaining a converged solution. [...] A guideline for a correct specification of boundary conditions at no-slip walls has been given by Roache.' His numerical recipe can be summarized as follows: at no-slip walls first specify the streamfunction and then, in the procurement of the wall vorticity, utilize the additional information on the normal component of the streamfunction. In this way the boundary conditions for the streamfunction are not overspecified and the wall vorticity is correctly related to the tangential component of the velocity. Obviously, at stationary walls the tangential component of the velocity is zero, but at moving walls serious errors may result if the gradient of the streamfunction is not properly taken into account.

### 13.4.3 Stress b.c.

To Do

### 13.4.4 Line of symmetry

When the flow is truly symmetrical, the axis of symmetry can be considered a streamline. Therefore the value of streamfunction along this boundary can be specified. Obviously, the velocity component normal to the the symmetry boundary would be zero, whereas the streamwise component is extrapolated from the interior solution.

## 13.5 Pressure Poisson Equation for the isoviscous N-S equation

This section is inspired by Salih [1098].

The PPE can be derived by taking the divergence of vector form of momentum equation

$$\vec{\nabla} \cdot \left[ \frac{\partial \vec{v}}{\partial t} + (\vec{v} \cdot \vec{\nabla}) \vec{v} \right] = \vec{\nabla} \cdot \left[ -\frac{1}{\rho} \vec{\nabla} p + \nu \vec{\nabla}^2 \vec{v} \right] \quad (13.32)$$

or,

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} = -\frac{1}{\rho} \frac{\partial p}{\partial x} + \nu \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) \quad (13.33)$$

$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} = -\frac{1}{\rho} \frac{\partial p}{\partial y} + \nu \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \quad (13.34)$$

Differentiating these equations with respect to  $x$  and  $y$  respectively, and assuming  $\rho$  constant in space yields

$$\frac{\partial^2 u}{\partial x \partial t} + \frac{\partial u}{\partial x} \frac{\partial u}{\partial x} + u \frac{\partial^2 u}{\partial x^2} + \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} + v \frac{\partial^2 u}{\partial x \partial y} = -\frac{1}{\rho} \frac{\partial^2 p}{\partial x^2} + \nu \frac{\partial}{\partial x} \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) \quad (13.35)$$

$$\frac{\partial^2 v}{\partial y \partial t} + \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + u \frac{\partial^2 v}{\partial y \partial x} + \frac{\partial v}{\partial y} \frac{\partial v}{\partial y} + v \frac{\partial^2 v}{\partial y^2} = -\frac{1}{\rho} \frac{\partial^2 p}{\partial y^2} + \nu \frac{\partial}{\partial y} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \quad (13.36)$$

Adding both equations together, i.e. taking the divergence of the first equation, yields

$$\begin{aligned} & \frac{\partial^2 u}{\partial x \partial t} + \frac{\partial u}{\partial x} \frac{\partial u}{\partial x} + u \frac{\partial^2 u}{\partial x^2} + \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} + v \frac{\partial^2 u}{\partial x \partial y} \frac{\partial^2 v}{\partial y \partial t} + \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + u \frac{\partial^2 v}{\partial y \partial x} + \frac{\partial v}{\partial y} \frac{\partial v}{\partial y} + v \frac{\partial^2 v}{\partial y^2} \\ = & -\frac{1}{\rho} \frac{\partial^2 p}{\partial x^2} + \nu \frac{\partial}{\partial x} \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) - \frac{1}{\rho} \frac{\partial^2 p}{\partial y^2} + \nu \frac{\partial}{\partial y} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \end{aligned} \quad (13.37)$$

$$\begin{aligned} & \partial_t \underbrace{\left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)}_{=0} + \frac{\partial u}{\partial x} \frac{\partial u}{\partial x} + \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} + \frac{\partial v}{\partial y} \frac{\partial v}{\partial y} + u \underbrace{\left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 v}{\partial y \partial x} \right)}_A + v \underbrace{\left( \frac{\partial^2 v}{\partial y^2} + \frac{\partial^2 u}{\partial x \partial y} \right)}_B \\ = & -\frac{1}{\rho} \left( \frac{\partial^2 p}{\partial x^2} + \frac{\partial^2 p}{\partial y^2} \right) + \nu \frac{\partial}{\partial x} \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + \nu \frac{\partial}{\partial y} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \end{aligned} \quad (13.38)$$

We have

$$\begin{aligned} A &= \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 v}{\partial y \partial x} = \frac{\partial u}{\partial x} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) = 0 \\ B &= \frac{\partial^2 v}{\partial y^2} + \frac{\partial^2 u}{\partial x \partial y} = \frac{\partial v}{\partial y} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) = 0 \end{aligned}$$

so that

$$\left( \frac{\partial u}{\partial x} \right)^2 + 2 \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \left( \frac{\partial v}{\partial y} \right)^2 = -\frac{1}{\rho} \left( \frac{\partial^2 p}{\partial x^2} + \frac{\partial^2 p}{\partial y^2} \right) + \nu \frac{\partial}{\partial x} \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + \nu \frac{\partial}{\partial y} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \quad (13.39)$$

The viscosity terms in the rhs can be rearranged as

$$\frac{\partial}{\partial x} \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + \frac{\partial}{\partial y} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) = \frac{\partial^2}{\partial x^2} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) + \frac{\partial^2}{\partial y^2} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) = 0$$

We are left with

$$\left( \frac{\partial u}{\partial x} \right)^2 + 2 \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \left( \frac{\partial v}{\partial y} \right)^2 = -\frac{1}{\rho} \left( \frac{\partial^2 p}{\partial x^2} + \frac{\partial^2 p}{\partial y^2} \right) \quad (13.40)$$

Now, the left-hand side can be further reduced as follows

$$\left( \frac{\partial u}{\partial x} \right)^2 + 2 \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \left( \frac{\partial v}{\partial y} \right)^2 = \underbrace{\left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)^2}_{=0} - 2 \frac{\partial u}{\partial x} \frac{\partial v}{\partial y} + 2 \frac{\partial u}{\partial y} \frac{\partial v}{\partial x}$$

In the end

$$\vec{\nabla}^2 p = 2\rho \left( \frac{\partial u}{\partial x} \frac{\partial v}{\partial y} - \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} \right)$$

The PPE can also be written in terms of stream function:

$$\boxed{\vec{\nabla}^2 p = 2\rho \left[ \frac{\partial^2 \Psi}{\partial x^2} \frac{\partial^2 \Psi}{\partial y^2} - \left( \frac{\partial^2 \Psi}{\partial x \partial y} \right)^2 \right]}$$

The Poisson equation for pressure is an elliptic equation, showing the elliptic nature of pressure in incompressible flows. For a steady flow problem, the PPE is solved only once, i.e., after the steady state values of  $\omega$  and  $\Psi$  have been computed.

Having established the PDE we must now turn to the boundary conditions. On a solid boundary, boundary values of pressure are obtained using the tangential momentum equation of the fluid

adjacent to the wall surface. For a wall located at  $y = 0$  in Cartesian coordinate system (i.e. the bottom boundary), the tangential momentum equation, that is the  $x$ -momentum equation, reduces to

$$\left. \frac{\partial p}{\partial x} \right|_{wall} = \eta \left. \frac{\partial^2 u}{\partial y^2} \right|_{wall} = -\eta \left. \frac{\partial \omega}{\partial y} \right|_{wall}$$

since  $v = \partial_x v = 0$  on that boundary.

In Reddy and Gartling [1051] we find at page 163: The pressure, if required, may be computed from a Poisson equation of the form:

$$\Delta P = \vec{\nabla} \cdot (\vec{f} - \nu \Delta \vec{v} - (\vec{v} \cdot \vec{\nabla}) \vec{v})$$

Check Gunzburger and Peterson [508] (1988) which “describe an alternate method of recovering the pressure which does not encounter any of the above difficulties; in particular, absolutely no boundary conditions on the pressure are needed at solid walls.”

## 13.6 Vorticity-stream function formulation in 3D

The vorticity-streamfunction approach has seen considerable use for two-dimensional incompressible flows. It has become less popular in recent years because its extension to three-dimensional flows is difficult. Both the vorticity and streamfunction become three-component vectors in three dimensions so one has a system of six partial differential equations in place of the four that are necessary in a velocity-pressure formulation. It also inherits the difficulties in dealing with variable fluid properties, compressibility, and boundary conditions that were described above for two dimensional flows.

In Reddy and Gartling [1051] we find at page 163:

“ A 3D version of the vorticity transport equation is possible, although it has seen relatively little use in computation due to the complexity of the vorticity boundary conditions. ”

Glaisner and Tezduyar [464] (1987) state: “

The extension of the codes to three dimensions is not easy for the vorticity stream-function formulation. In three dimensions, the vorticity- stretching term does not vanish and the equation system has two nonlinear terms. ”

In Zhong, Yuen, Moresi, and Knepley [1415] (2012) we read:

“For 3-D mantle convection, we can likewise employ the poloidal potential  $\Phi$  and a vorticity-like scalar function  $\Omega$  (Busse, 1989; Chandrasekhar, 1961; Travis et al., 1990). It is to be noted that this potential  $\Phi$  is not the same as the stream function  $\Psi$  in 2-D (see Chapter 7.04). From the general representation of an arbitrary solenoidal vector field (Busse, 1989), we can write a 3-D velocity vector as

$$\vec{v} = \vec{\nabla} \times \vec{\nabla} \times (\Phi \vec{e}_z) + \vec{\nabla} \times (\Theta \vec{e}_z)$$

where  $\vec{e}_z$  is the unit vector in the vertical,  $z$ -direction, pointing upward.  $\Theta$  is the toroidal potential and is present in problems with lateral variations of viscosity (Christensen and Harder, 1991; Gable et al., 1991; Zhang and Christensen, 1993). Thus, for constant viscosity,  $\Theta$  is zero unless driven by a boundary condition (e.g., Gable et al., 1991), and the velocity vector  $\vec{v} = (u, v, w)$  involves higher-order derivatives of  $\Phi$  in this formulation:

$$u = \frac{\partial^2 \Phi}{\partial y \partial z} \quad v = \frac{\partial^2 \Phi}{\partial x \partial z} \quad w = -\left( \frac{\partial^2 \Phi}{\partial x^2} + \frac{\partial^2 \Phi}{\partial y^2} \right)$$

The 3-D momentum equation for constant properties can be written as a system of coupled Poisson equations in 3-D:

$$\vec{\nabla}^2 \Phi = \Omega \quad (13.41)$$

$$\vec{\nabla}^2 \Omega = \text{Ra}T \quad (13.42)$$

where all differential operators are 3-D in character and  $\Omega$  is a scalar function playing a role analogous to the vorticity in the 2-D formulation.

We note that FD and FV methods using the primitive variables of velocity and pressure are currently predominant in models of 3-D mantle convection with variable viscosity. ”

ToDo:

1. scan 3D literature
2. find all formulations cartesian or not
3. document num methods used

### 13.6.1 Cartesian domain (1)

Sotin and Labrosse [1180] (1999) write:

“ The conservation equations for momentum and mass are transformed into four Poisson equations by introducing a vector potential stream function  $\Psi$  and vorticity  $\omega$  [1280]:

$$\Delta \Psi_x + \omega_x = 0 \quad (13.43)$$

$$\Delta \omega_x - \text{Ra}_T \frac{\partial T}{\partial y} = 0 \quad (13.44)$$

$$\Delta \Psi_y + \omega_y = 0 \quad (13.45)$$

$$\Delta \omega_y + \text{Ra}_T \frac{\partial T}{\partial x} = 0 \quad (13.46)$$

Since the buoyancy force acts only in the  $z$ -direction,  $\Psi_z$  and  $\omega_z$  vanish identically. Shear-stress free boundaries on top and bottom yields the following boundary conditions:

$$\omega_x = \omega_y = \Psi_x = \Psi_y = 0$$

at the top and bottom. Periodicity is assumed on all vertical boundaries.

The four Poisson equations are solved using a multigrid iterative method described by Parmentier, Sotin, and Travis [976] (1994) and Sotin et al. 1995. ”

### 13.6.2 Cartesian domain (2)

Let us now turn to Larsen, Yuen, Malevsky, and Smedsmo [750] (1996):

“For the spatial discretization fourth-order correct bi-cubic splines are employed [826]. [...] The mesh is uniform in the horizontal direction, and non-uniform in the vertical direction with mesh-refinement near the boundary layers. [...] The mechanical boundary conditions are stress-free and impermeable at the top and bottom boundaries, and reflecting along the sides. The dimensionless temperature  $T$  is zero at the top and unity at the bottom ( $z = 1$ ), and there is zero heat flux along the sides.”

In Larsen, Yuen, Moser, and Fornberg [751] (1997) the authors start with a 2D formulation:

“Instead of the biharmonic equation for  $\Psi$ , (e.g., Christensen [243], 1984), the conservation equations for the mass and momentum are given by two second-order partial differential equations (1) and (2). The horizontal and vertical coordinates are, respectively,  $x$  and  $z$  with  $z$  pointing upwards. Time,  $t$ , has been non-dimensionalized by the thermal diffusion time across the depth of the layer. A more detailed description of the scheme for non-dimensionalization is given in Weinstein, Olson, and Yuen [1347] (1989). The boundary conditions are stress-free and impermeable at the top and bottom boundaries, and reflecting along the sides. The dimensionless temperature,  $T$ , is 0 at the top and 1 at the bottom. The boundary conditions are:  $\Psi = \omega = 0$  on all sides.”

They go further and explain the 3D approach:

“ We employ the poloidal potential,  $\Psi$ , and a vorticity-like scalar function  $\Omega$  for solving the momentum equation with constant properties (Busse [194], 1989; Travis, Olson, and Schubert [1280], 1990b). We note that the potential  $\Psi$  is not the same as the streamfunction in two dimensions. From the general representation of an arbitrary solenoidal vector field (Busse [194], 1989), we can write a three-dimensional velocity vector,

$$\vec{v} = \vec{\nabla} \times (\vec{\nabla} \times \vec{e}_z \Psi)$$

with the  $\vec{e}_z$  axis pointing upwards. Thus the velocity vector involves higher-order derivatives of  $\Psi$  in this formulation:

$$u = \frac{\partial^2 \Psi}{\partial x \partial z} \quad v = \frac{\partial^2 \Psi}{\partial y \partial z} \quad \omega = - \left( \frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} \right)$$

with  $\vec{v} = (u, v, w)$ . The three-dimensional momentum equation can be written as a system of coupled Poisson equations:

$$\vec{\nabla}^2 \Psi = \Omega \quad \vec{\nabla}^2 \Omega = \text{Ra} T$$

Here  $\Omega$  is a scalar function playing an analogous role to the vorticity in the two-dimensional formulation.

The numerical method employed in this paper is based on a recursive algorithm for generating the weights in the finite-difference approximation.”

### 13.6.3 Cartesian domain (3)

See for example Houseman [593] (1990) and refs therein:

“As the flow is incompressible, velocity may be expressed as the curl of a solenoidal vector potential  $\vec{\Psi}$  [e.g. Richardson & Cornish (1977)]:

$$\vec{v} = \vec{\nabla} \times \vec{\Psi}$$

For 2-D flow,  $\vec{\Psi}$  is a vector everywhere normal to the plane of the flow. It may then be treated as a scalar, and is given the name streamfunction.

The momentum equation may be restated as a biharmonic equation:

$$\vec{\nabla}^4 \vec{\Psi} = \vec{f}$$

where  $\vec{f}$  is minus the curl of the buoyancy force divided by viscosity, and is hereafter treated as a known quantity. It is also useful to define the vorticity, which is simply related to the potential function

$$\vec{\omega} = \vec{\nabla} \times \vec{v} = -\vec{\nabla}^2 \Psi$$

We now consider a finite difference technique for solving equation (3) in a region spanned by a regular isotropic array of nodes.”

Boundary conditions of all kinds are extensively discussed!

### 13.6.4 Spherical shell

Zebib, Schubert, and Straus [1402] (1980) present axisymmetric steady convective solutions

Likewise Solheim and Peltier [1178] and Solheim and Peltier [1177]: “The model is spherical but restricted in generality to the analysis of axisymmetric solutions.”

## 13.7 Vorticity-stream function formulation in polar/cylindrical coordinates

[https://en.wikipedia.org/wiki/Stokes\\_stream\\_function](https://en.wikipedia.org/wiki/Stokes_stream_function) uses very different definitions! **sort it out!** actually it is for three-dimensional incompressible flow with axisymmetry!

$$\begin{aligned} v_r &= \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \\ v_\theta &= -\frac{\partial \Psi}{\partial r} \\ \vec{\omega} &= \left( \frac{1}{r} \frac{\partial v_z}{\partial \theta} - \frac{\partial v_\theta}{\partial z} \right) \vec{e}_r + \left( \frac{\partial v_r}{\partial z} - \frac{\partial v_z}{\partial r} \right) \vec{e}_\theta + \left( \frac{1}{r} \frac{\partial(r v_\theta)}{\partial r} - \frac{1}{r} \frac{\partial v_r}{\partial \theta} \right) \vec{e}_z \\ \vec{\nabla}^4 \Psi &= \left( \frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} \right) \left( \frac{\partial^2 \Psi}{\partial r^2} + \frac{1}{r} \frac{\partial \Psi}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \Psi}{\partial \theta^2} \right) \end{aligned}$$

or

$$\vec{\nabla}^4 \Psi = \Psi_{,rrrr} + \frac{2}{r} \Psi_{,rrr} - \frac{1}{r^2} (\Psi_{,rr} - 2\Psi_{,rr\theta\theta}) + \frac{1}{r^3} (\Psi_{,r} - 2\Psi_{,r\theta\theta}) + \frac{1}{r^4} (4\Psi_{,\theta\theta} + 2\Psi_{,\theta\theta\theta\theta})$$

See Hsui [597]. [1277] [1279]

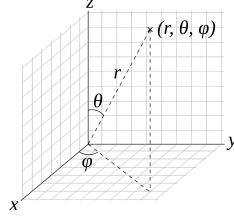
## 13.8 Vorticity-stream function formulation in spherical coordinates for three-dimensional incompressible flow with axisymmetry

What follows comes from [https://en.wikipedia.org/wiki/Stokes\\_stream\\_function](https://en.wikipedia.org/wiki/Stokes_stream_function).

The flow velocity components  $v_r$  and  $v_\theta$  are related to the Stokes stream function  $\Psi$  through:

$$v_r = \frac{1}{r^2 \sin \theta} \frac{\partial \Psi}{\partial \theta} \quad v_\theta = -\frac{1}{r \sin \theta} \frac{\partial \Psi}{\partial r}$$

The azimuthal velocity component  $v_\phi$  is not a function of the Stokes stream function  $\Psi$ .



The vorticity is defined as:

$$\vec{\omega} = \vec{\nabla} \times \vec{v} = \vec{\nabla} \times \vec{\nabla} \times \Psi \quad \text{with} \quad \Psi = -\frac{\Psi}{r \sin \theta} \vec{e}_\phi$$

From the definition of the curl in spherical coordinates<sup>2</sup>:

$$\omega_r = \frac{1}{r \sin \theta} \left( \frac{\partial}{\partial \theta} (v_\phi \sin \theta) - \frac{\partial v_\theta}{\partial \phi} \right) \vec{e}_r \quad (13.47)$$

$$\omega_\theta = \frac{1}{r} \left( \frac{1}{\sin \theta} \frac{\partial v_r}{\partial \phi} - \frac{\partial}{\partial r} (r v_\phi) \right) \vec{e}_\theta \quad (13.48)$$

$$\omega_\phi = \frac{1}{r} \left( \frac{\partial}{\partial r} (r v_\theta) - \frac{\partial v_r}{\partial \theta} \right) \vec{e}_\phi \quad (13.49)$$

First notice that the  $r$  and  $\theta$  components are equal to 0. **prove!** Secondly substitute  $v_r$  and  $v_\theta$  into  $\omega_\phi$ . The result is:

$$\omega_r = 0 \quad (13.50)$$

$$\omega_\theta = 0 \quad (13.51)$$

$$\omega_\phi = \frac{1}{r} \left[ \frac{\partial}{\partial r} \left( r \left( -\frac{1}{r \sin \theta} \frac{\partial \Psi}{\partial r} \right) \right) - \frac{\partial}{\partial \theta} \left( \frac{1}{r^2 \sin \theta} \frac{\partial \Psi}{\partial \theta} \right) \right] \quad (13.52)$$

After some algebra we arrive at

$$\vec{\omega} = \begin{pmatrix} 0 \\ 0 \\ -\frac{1}{r \sin \theta} \left( \frac{\partial^2 \Psi}{\partial r^2} + \frac{\sin \theta}{r^2} \frac{\partial}{\partial \theta} \left( \frac{1}{\sin \theta} \frac{\partial \Psi}{\partial \theta} \right) \right) \end{pmatrix}$$

Proof that velocity is perpendicular to gradient of stream function:

$$\vec{\nabla} \Psi \cdot \vec{v} = \frac{\partial \Psi}{\partial r} \cdot \frac{1}{r^2 \sin \theta} \frac{\partial \Psi}{\partial \theta} + \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \cdot \left( -\frac{1}{r \sin \theta} \frac{\partial \Psi}{\partial r} \right) = 0$$

## 13.9 Numerical approach

When it comes to solving the biharmonic equation, there are essentially two options: solving/discretising the biharmonic operator involving 4th order derivatives, or introducing the vorticity and solving coupled PDEs in vorticity-stream function.

Malevsky and Yuen [828] state:

“ The fourth-order elliptic equation does not contain time explicitly. One can obtain the stream-function for a given temperature field by solving this nonlinear equation at each instant. Equation (12) is a nonlinear time-dependent advection- diffusion equation, where the  $\Psi(T)$  dependence is given by the biharmonic eq. Numerical solution of the

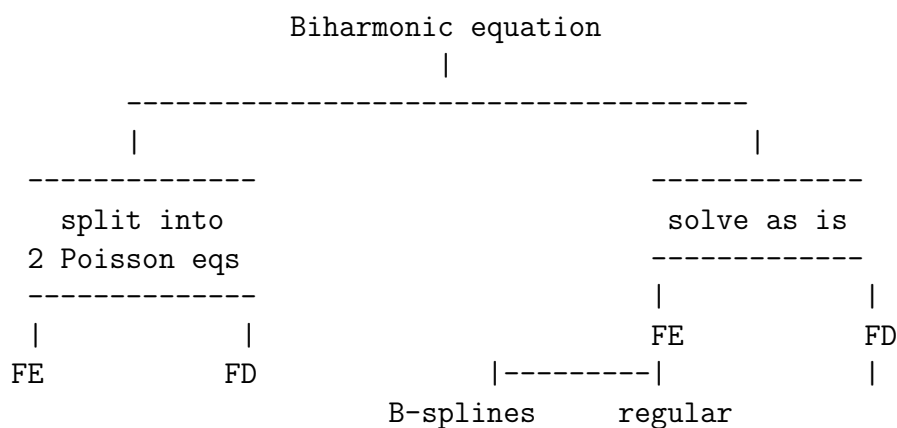
<sup>2</sup>[https://en.wikipedia.org/wiki/Del\\_in\\_cylindrical\\_and\\_spherical\\_coordinates](https://en.wikipedia.org/wiki/Del_in_cylindrical_and_spherical_coordinates)

advection-diffusion equation with advection dominating over diffusion is fraught with numerical difficulties. It is known that high-order finite-element and finite- difference schemes deteriorate from errors located near steep gradients of the advected field. Various upwinding schemes, designed to overcome this difficulty, introduce artificial smoothing (Cuvelier et al., 1988). A finite-element scheme based on the Lagrangian formulation of the total time derivative in the advection-diffusion equation [826] proves to be very efficient for the high Rayleigh number Newtonian convection. In stress-softening fluid, thermal advection can be very strong locally due to the decrease of the effective viscosity with growing stress. Therefore, we have chosen the Lagrangian scheme to solve the energy equation for the non-Newtonian case. A second-order predictor-corrector scheme with Lagrangian formulation was applied for time-stepping in the energy equation. The Lagrangian scheme requires interpolation to compute values of  $\Psi$  and  $T$  between nodal points. This scheme is very sensitive to the quality of spatial approximation. We used bicubic splines (Ahlberg et al., 1967) for spatial discretization of both the temperature and stream-function. ”

Rather interestingly Poliakov and Podlachikov [1008] (1992) state:

“ There is a class of methods based on the introduction of a stream function which satisfies this condition automatically and gives good results. However, the stress-free boundary condition causes significant difficulties because it requires the calculation of third-order derivatives of the stream function. These algorithms are also not applicable for irregular meshes. Thus, it is difficult to solve problems in regions with complicated geometry, and to refine the mesh in areas of special interest. A further disadvantage is that the approximation functions are not expressed explicitly in terms of the nodal variables and an additional system of linear equations must be solved. ”

At the moment (and before a thorough literature review) here is how I see things when it comes to numerical methods:



### 13.9.1 Finite differences

In Valera, Negredo, and Villaseñor [1302] we find

“To solve the equation of motion we have applied a second-order, central finite- difference scheme. We considered free slip (zero shear stress) boundary conditions at all boundaries. The solution of the equation on each node was then computed in terms of the 12 nearest nodes. The system matrix was square, symmetric and diagonally dominant, with only 12 non-null diagonals. This system could be solved by a very robust method based on a LU



triangular factorization by Gaussian elimination with partial pivoting, which is already implemented in the standard MATLAB code.”

In Zhong, Yuen, Moresi, and Knepley [1415] (2012) we read:

“We note that solving the biharmonic equation by FDs takes more time than solving two coupled Laplacian equations.”

ToDo:

1. scan literature. e.g. McKenzie, Roberts, and Weiss [857] (1974) explain their stencils
2. establish stencil(s) for equations
3. how to deal with boundaries
4. In gerya’s book it is solved with FD. see page 72, and example Streamfunction2D.m - translate to python?

### 13.9.2 Finite elements

We start from  $\vec{\nabla}^2 \vec{\nabla}^2 \Psi = f$ , which we rewrite as follows:

$$\vec{\nabla}^2 \Psi = -\omega \quad (13.53)$$

$$\vec{\nabla}^2 \omega = f \quad (13.54)$$

After establishing the weak form and discretising it, this will yield the following linear system:

$$\begin{pmatrix} K & M \\ 0 & K \end{pmatrix} \cdot \begin{pmatrix} \vec{\Psi} \\ \vec{\omega} \end{pmatrix} = \begin{pmatrix} \vec{0} \\ \vec{f} \end{pmatrix}$$

where  $K = \int B^T B dV$ . This is an ideal situation: one can first solve the second line, obtain  $\vec{\omega}$  and then solve the first line as  $K \cdot \vec{\Psi} = -M \cdot \vec{\omega}$ . Also it is the same matrix  $K$ , only different rhs!

Can i do this with Q1 elements only?!

Remarks:

- See Gresho and Sani [488] p523-525 for important remarks about establishing the weak forms of the vorticity-stream function equations. They refer to [1212].
- In Step 47 of deal.II<sup>3</sup> it is explained

“The fundamental issue with the equation is that [the biharmonic equation] takes four derivatives of the solution. In the case of the Laplace equation [...], and several other tutorial programs, one multiplies by a test function, integrates, integrates by parts, and ends up with only one derivative on both the test function and trial function - something one can do with functions that are continuous globally, but may have kinks at the interfaces between cells: The derivative may not be defined at the interfaces, but that is on a lower-dimensional manifold (and so doesn’t show up in the integrated value).

But for the biharmonic equation, if one followed the same procedure using integrals

<sup>3</sup>[https://www.dealii.org/current/doxygen/deal.II/step\\_47.html](https://www.dealii.org/current/doxygen/deal.II/step_47.html)

over the entire domain (i.e., the union of all cells), one would end up with two derivatives on the test functions and trial functions each. If one were to use the usual piecewise polynomial functions with their kinks on cell interfaces, the first derivative would yield a discontinuous gradient, and the second derivative with delta functions on the interfaces - but because both the second derivatives of the test functions and of the trial functions yield a delta function, we would try to integrate the product of two delta functions.”

Logically the biharmonic equation is then split into two Poisson equations.

- In van Keken’s phd thesis we read:

“The equation of motion is a 4th-order PDE and, for a FE approximation to be conform, the basis functions should be twice continuously differentiable in each element and least continuously differentiable throughout the computational domain  $\Omega$ . ”

van Keken then proceeds to use 2 different methods: a non-conforming type of element Hansen and Ebel [533] (1984), and a bicubic spline FE method Woidt [1366] (1978), Kopitzke et al. [722] (1979).

- check Comini, Manzan, and Nonino [275] (1994) for N-S implementation.
- check Gunzburger and Peterson [508] (1988) for 2D, 3D implementation

## the Spline FE method

Christensen [243] (1984) writes:

“The numerical grid consists of rectangular elements of variable size. All spatially varying quantities - i.e.  $\Psi$ ,  $T$  and  $\eta$  are represented by polynomial splines. These are polynomials of a certain degree  $n$  within each individual element which are continuous at the joints up to the  $k$ th derivative (  $k < n$ ). In a spline space at a rather compact base can be constructed, which means that each base function (B spline) is zero outside a small region consisting of few elements; this is an advantage over a polynomial or Fourier space. The shape of each B spline depends on the local grid structure. With the restriction that the size of the elements varies in both dimensions only by one of the factors  $1/\alpha$ ,  $1$  or  $\alpha$  compared to the neighbouring elements (with  $\alpha$  fixed) a restricted number of types of B splines are necessary which reduces the computational costs. Splines which are non-zero at the boundary of the mathematical domain are modified in such a way that the boundary conditions to [the biharmonic equation] are fulfilled for any function that can be constructed in the spline space. Compared to the normal finite element method with Lagrangian shape functions, higher requirements on the continuity can be easily fulfilled. On the other hand one is restricted to semi-regular grids with all element boundaries parallel or perpendicular to each other. [...] The integral contains fourth-order derivatives, however, by partial integration it can be transformed into one containing only second derivatives. The matrix built is band structured and a stable solution is obtained by a Cholesky transformation.”

Christensen and Yuen [246] (1989) states:

“The spline finite element method used to solve the set of equations (7) and (9) is described by Christensen [1984]. Bicubic splines are taken for  $\Psi$ , and biquadratic splines for  $T$ . The grid is nonuniform and allows higher resolution in the boundary layers”

In Malevsky and Yuen [828] we find:

“We can use the Galerkin method to approximate the solution. In this case a matrix  $\mathbf{A}$  is obtained as

$$a_{ij} = \iint \eta \left[ \left( \frac{\partial^2 B_i}{\partial y^2} - \frac{\partial^2 B_i}{\partial x^2} \right) \left( \frac{\partial^2 B_j}{\partial y^2} - \frac{\partial^2 B_j}{\partial x^2} \right) + 4 \frac{\partial B_i}{\partial y \partial x} \frac{\partial B_j}{\partial y \partial x} \right] dx dy$$

where  $B$  is a basic bicubic spline. The derivatives of a basic spline (a piecewise cubic polynomial) can be calculated analytically.”

ToDo:

- derive B spline equation, implement, python fet? check for example [722]
- derive equation above
- work out the 2 formulations of vKK phd thesis
- check book by Gunzburger & Peterson for splines + FEM

## 13.10 Algorithm for stream function-vorticity formulation

A solution algorithm for computing evolution of incompressible, two-dimensional flow using stream function-vorticity formulation is given as follows:

1. Initialize the velocity field and compute the associated vorticity field and streamfunction field using equations  $\omega_z = \partial_x v - \partial_y u$  and  $\Delta \Psi = -\omega$ .
2. Compute the boundary conditions for vorticity.
3. Solve the vorticity transport equation (13.6) to compute the vorticity at new time step; any standard time marching scheme may be used for this purpose.
4. Solve the Poisson equation for streamfunction  $\Delta \Psi = -\omega$  to compute the streamfunction field at new time step; any iterative scheme for elliptic equations may be used.
5. Compute the velocity field at new time step using the relations  $u = \partial_y \Psi$  and  $v = -\partial_x \Psi$ .
6. Return to step 2 and repeat the computation for another time step.

## 13.11 The nondimensional equations

### 13.11.1 Isoviscous case

We start from the equations (2.81),(2.82),(2.83):

$$-\vec{\nabla} p + \vec{\nabla} \cdot 2\dot{\epsilon} + \text{Ra} T \vec{e}_y = \vec{0} \quad (13.55)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (13.56)$$

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T \quad (13.57)$$

If the viscosity is constant then  $\eta = \eta_0 = \eta_{ref}$  so  $\eta = 1$  we can also write these as

$$-\vec{\nabla} p + \vec{\nabla}^2 \vec{v} + \text{Ra} T \vec{e}_y = \vec{0} \quad (13.58)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (13.59)$$

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T \quad (13.60)$$

Following the approach in Section 13.2, and omitting the colors to indicate dimensionless values: or, in each dimension:

$$0 = -\frac{\partial p}{\partial x} + \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial x^y} \right) \quad (13.61)$$

$$0 = -\frac{\partial p}{\partial y} + \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) + \text{Ra} T \quad (13.62)$$

Using Eq. (13.13), we can write

$$\begin{aligned} 0 &= -\frac{\partial p}{\partial x} + \left( \frac{\partial^3 \Psi}{\partial y^3} + \frac{\partial^3 \Psi}{\partial x^2 \partial y} \right) \\ 0 &= -\frac{\partial p}{\partial y} - \left( \frac{\partial^3 \Psi}{\partial x^3} + \frac{\partial^3 \Psi}{\partial y^2 \partial x} \right) + \text{Ra} T \end{aligned} \quad (13.63)$$

The pressure terms in both equations can be removed by first differentiating the first line with regards to  $y$  and the second line with regards to  $x$ ,

$$\begin{aligned} 0 &= -\frac{\partial^2 p}{\partial x \partial y} + \left( \frac{\partial^4 \Psi}{\partial y^4} + \frac{\partial^4 \Psi}{\partial x^2 \partial y^2} \right) \\ 0 &= -\frac{\partial^2 p}{\partial y \partial x} - \left( \frac{\partial^4 \Psi}{\partial x^4} + \frac{\partial^4 \Psi}{\partial y^2 \partial x^2} \right) + \frac{\partial(\text{Ra} T)}{\partial x} \end{aligned} \quad (13.64)$$

and next by subtracting the resulting equations, leading to:

$$0 = \left( \frac{\partial^4 \Psi}{\partial x^4} + 2 \frac{\partial^4 \Psi}{\partial x^2 \partial y^2} + \frac{\partial^4 \Psi}{\partial y^4} \right) - \text{Ra} \frac{\partial T}{\partial x} \quad (13.65)$$

or

$$\vec{\nabla}^4 \Psi = \text{Ra} \frac{\partial T}{\partial x} \quad (13.66)$$

Introducing again the vorticity  $\omega = -\vec{\nabla}^2 \Psi$ , then we must solve

$$\vec{\nabla}^2 \omega = -\text{Ra} \frac{\partial T}{\partial x} \quad \text{and} \quad \vec{\nabla}^2 \Psi = -\omega$$

check further

The vorticity equation is

$$\frac{\partial \omega}{\partial t} + \vec{v} \cdot \vec{\nabla} \omega = \text{Re}^{-1} \vec{\nabla}^2 \omega$$

ToDo:

- re-derive these equations , see for example Solheim and Peltier [1177] (1994) eq 8, or Kopitzke et al. [722] eq 1,2

## 13.12 Incorporation of phase changes

ToDo:

- scan literature for  $\Gamma$  fct
- isolate example

## 13.13 The energy equation

The (simple) energy equation

$$\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T = \kappa \Delta T$$

becomes in 2D:

$$\frac{\partial T}{\partial t} + \frac{\partial \Psi}{\partial y} \frac{\partial T}{\partial x} - \frac{\partial \Psi}{\partial x} \frac{\partial T}{\partial y} = \kappa \Delta T$$

rewrite with all coeffs

## 13.14 Remark, misc

Arie van den Berg (prov. comm.) writes:

“The stream function formulation is useful to obtain a conceptually simpler implementation with finite differences if you rewrite the 4th order biharmonic equation for the streamfunction in two coupled Poisson equations for streamfunction and vorticity respectively. That is for an isoviscous fluid. For variable viscosity it becomes more ugly.

It disappeared probably due its limitations for variable viscosity fluids and once the more power full but complex formulation alternatives became more well known. ”

In Gresho and Sani [488] (p. 941), the authors present an algorithm to compute  $\Psi$  from a velocity field obtained via the FEM.

## 13.15 Literature

- 1961** •Hsiao-Lan Kuo. “Solution of the non-linear equations of cellular convection and heat transport”. In: *Journal of Fluid Mechanics* 10.4 (1961), pp. 611–634. doi: 10.1017/S0022112061000408
- 1967** •DL Turcotte and ER Oxburgh. “Finite amplitude convective cells and continental drift”. In: *Journal of Fluid Mechanics* 28.1 (1967), pp. 29–42. doi: 10.1017/S0022112067001880
- 1969** •G Schubert, DL Turcotte, and ER Oxburgh. “Stability of planetary interiors”. In: *Geophysical Journal International* 18.5 (1969), pp. 441–460. doi: 10.1111/j.1365-246X.1969.tb03370.x
- 1970** •Hans Ramberg. “Folding of laterally compressed multilayers in the field of gravity, I”. in: *Physics of the Earth and Planetary Interiors* 2.4 (1970), pp. 203–232. doi: 10.1016/0031-9201(70)90010-5
- 1971** •KE Torrance and DL Turcotte. “Thermal convection with large viscosity variations”. In: *Journal of Fluid Mechanics* 47.1 (1971), pp. 113–125. doi: 10.1017/S002211207100096X  
•KE Torrance and DL Turcotte. “Structure of convection cells in the mantle”. In: *Journal of Geophysical Research* 76.5 (1971), pp. 1154–1161. doi: 10.1029/JB076i005p01154
- 1973** •Frank M Richter. “Convection and the large-scale circulation of the mantle”. In: *Journal of Geophysical Research* 78.35 (1973), pp. 8735–8745. doi: 10.1029/JB078i035p08735  
•Frank M Richter. “Dynamical models for sea floor spreading”. In: *Reviews of Geophysics* 11.2 (1973), pp. 223–287. doi: 10.1029/RG011i002p00223  
•Dan McKenzie, Jean Roberts, and Nigel Weiss. “Numerical models of convection in the earth’s mantle”. In: *Tectonophysics* 19.2 (1973), pp. 89–103. doi: 10.1016/0040-1951(73)90034-6

- 1974** ●Patrick Cassen and Ray T Reynolds. “Convection in the Moon: Effect of variable viscosity”. In: *Journal of Geophysical Research* 79.20 (1974), pp. 2937–2944. DOI: 10.1029/JB079i020p02937  
 ●Dan P McKenzie, Jean M Roberts, and Nigel O Weiss. “Convection in the Earth’s mantle: towards a numerical simulation”. In: *Journal of Fluid Mechanics* 62.3 (1974), pp. 465–538. DOI: 10.1017/S0022112074000784
- 1975** ●J.M. Hewitt, D.P. McKenzie, and N.O. Weiss. “Dissipative heating in convective flows”. In: *J. Fluid Mech.* 68.4 (1975), pp. 721–738. DOI: 10.1017/S002211207500119X  
 ●EM Parmentier, DL Turcotte, and KE Torrance. “Numerical experiments on the structure of mantle plumes”. In: *Journal of Geophysical Research* 80.32 (1975), pp. 4417–4424. DOI: 10.1029/JB080i032p04417
- 1976** ●EM Parmentier, DL Turcotte, and KE Torrance. “Studies of finite amplitude non-Newtonian thermal convection with application to convection in the Earth’s mantle”. In: *Journal of Geophysical Research* 81.11 (1976), pp. 1839–1846. DOI: 10.1029/JB081i011p01839
- 1977** ●Geoffrey F Davies. “Whole-mantle convection and plate tectonics”. In: *Geophysical Journal International* 49.2 (1977), pp. 459–486. DOI: 10.1111/j.1365-246X.1977.tb03717.x  
 ●Geoffrey F Davies. “Viscous mantle flow under moving lithospheric plates and under subduction zones”. In: *Geophysical Journal International* 49.3 (1977), pp. 557–563. DOI: 10.1111/j.1365-246X.1977.tb01303.x
- 1978** ●Peter Bird. “Finite element modeling of lithosphere deformation: the Zagros collision orogeny”. In: *Tectonophysics* 50.2-3 (1978), pp. 307–336  
 ●Albert T Hsui. “Numerical simulation of finite-amplitude thermal convection with large viscosity variation in axisymmetric spherical geometry: effect of mechanical boundary conditions”. In: *Tectonophysics* 50.2-3 (1978), pp. 147–162. DOI: 10.1016/0040-1951(78)90132-4  
 ●D.K. Gartling. *Nachos - A finite element computer program for incompressible flow problems*. Tech. rep. Sand77-1333. Sandia Laboratories, 1978  
 ●Aaron Tovish, Gerald Schubert, and Bruce P Luyendyk. “Mantle flow pressure and the angle of subduction: Non-Newtonian corner flows”. In: *Journal of Geophysical Research: Solid Earth* 83.B12 (1978), pp. 5892–5898  
 ●W.-D. Woitd. “Finite element calculations applied to salt dome analysis”. In: *Tectonophysics* 50 (1978), pp. 369–386. DOI: 10.1016/0040-1951(78)90143-9  
 ●F. Richter and D. McKenzie. “Simple plate models of mantle convection”. In: *J. Geophys.* 44 (1978), pp. 441–471. DOI: xxxx
- 1979** ●D.F. Griffiths. “Finite Elements for Incompressible Flow”. In: *Math. Meth. in the Appl. Sci.* 1 (1979), pp. 16–31  
 ●Dan McKenzie. “Finite deformation during fluid flow”. In: *Geophysical Journal International* 58.3 (1979), pp. 689–715. DOI: 10.1111/j.1365-246X.1979.tb04803.x  
 ●L. Bodri and B. Bodri. “Flow, stress and temperature in island arc areas”. In: *Geophysical & Astrophysical Fluid Dynamics* 13.1 (1979), pp. 95–105. doi: 10.1080/03091927908243763  
 ●U Kopitzke et al. “Finite element convection models: comparison of shallow and deep mantle convection, and temperatures in the mantle”. In: *Journal of Geophysics* 46.1 (1979), pp. 97–121. DOI: xxx  
 ●Richard A Lux, Geoffrey F Davies, and John H Thomas. “Moving lithospheric plates and mantle convection”. In: *Geophysical Journal International* 58.1 (1979), pp. 209–228
- 1980** ●Gary T. Jarvis and Dan P. McKenzie. “Convection in a compressible fluid with infinite Prandtl number”. In: *Journal of Fluid Mechanics* 96.3 (1980), pp. 515–583. DOI: 10.1017/S002211208000225X  
 ●Abdelfattah Zebib, Gerald Schubert, and Joe M Straus. “Infinite Prandtl number thermal convection in a spherical shell”. In: *Journal of Fluid Mechanics* 97.2 (1980), pp. 257–277. DOI: 10.1017/S0022112080002558  
 ●D.A. Yuen and W.R. Peltier. “Mantle plumes and the thermal stability of the D” layer”. In: *Geophysical Research Letters* 7.9 (1980), pp. 625–628. DOI: 10.1029/GL007i009p00625
- 1981** ●Wolfgang R Jacoby and Harro Schmeling. “Convection experiments and the driving mechanism”. In: *Geologische Rundschau* 70.1 (1981), pp. 207–230. DOI: 10.1007/BF01764323  
 ●H. Schmeling and W.R. Jacoby. “On modelling the lithosphere in mantle convection with non-linear rheology”. In: *Journal of Geophysics* 50 (1981), pp. 89–100  
 ●D.A. Yuen, W.R. Peltier, and G. Schubert. “On the existence of a second scale of convection in the upper mantle”. In: *Geophysical Journal of the Royal Astronomical Society* 65.1 (1981), pp. 171–190. DOI: 10.1111/j.1365-246X.1981.tb02707.x  
 ●Frank M Richter and Dan P McKenzie. “Parameterizations for the horizontally averaged temperature of infinite Prandtl number convection”. In: *Journal of Geophysical Research: Solid Earth* 86.B3 (1981), pp. 1738–1744. DOI: 10.1029/JB086iB03p01738  
 ●Greg A Houseman, D Po McKenzie, and Peter Molnar. “Convective instability of a thickened boundary layer and its relevance for the thermal evolution of continental convergent belts”. In: *Journal of Geophysical Research: Solid Earth* 86.B7 (1981), pp. 6115–6132. DOI: 10.1029/JB086iB07p06115
- 1982** ●P. England. “Some numerical investigations of large scale continental deformation”. In: *Mountain Building Processes*. Academic Press, 1982, pp. 129–189. DOI: xxxx  
 ●Gary T Jarvis and WR Peltier. “Mantle convection as a boundary layer phenomenon”. In: *Geophysical Journal International* 68.2 (1982), pp. 389–427. DOI: 10.1111/j.1365-246X.1982.tb04907.x  
 ●Frank J Spera, David A Yuen, and Stephen J Kirschvink. “Thermal boundary layer convection in silicic magma chambers: Effects of temperature-dependent rheology and implications for thermogravitational chemical fractionation”. In: *Journal of Geophysical Research: Solid Earth* 87.B10 (1982), pp. 8755–8767. DOI: 10.1029/JB087iB10p08755
- 1983** ●Takeshi Matsumoto and Yoshibumi Tomoda. “Numerical simulation of the initiation of subduction at the fracture zone”. In: *Journal of Physics of the Earth* 31.3 (1983), pp. 183–194. DOI: 10.4294/jpe1952.31.183  
 ●U Christensen. “Convection in a variable-viscosity fluid: Newtonian versus power-law rheology”. In: *Earth and Planetary Science Letters* 64.1 (1983), pp. 153–162. DOI: 10.1016/0012-821X(83)90060-2  
 ●Ulrich Christensen. “A numerical model of coupled subcontinental and oceanic convection”. In: *Tectonophysics* 95.1-2 (1983), pp. 1–23. DOI: 10.1016/0040-1951(83)90256-1

- 1984** ●U.R. Christensen and D.A. Yuen. "The interaction of a subducting lithospheric slab with a chemical or phase boundary". In: *J. Geophys. Res.* 89(B6) (1984), pp. 4389–4402. DOI: 10.1029/JB089iB06p04389
- U Christensen. "Convection with pressure- and temperature-dependent non-Newtonian rheology". In: *Geophysical Journal International* 77.2 (1984), pp. 343–384. DOI: 10.1111/j.1365-246X.1984.tb01939.x
- U Hansen and A Ebel. "Experiments with a numerical model related to mantle convection: boundary layer behaviour of small- and large scale flows". In: *Physics of the earth and planetary interiors* 36.3-4 (1984), pp. 374–390. DOI: 10.1016/0031-9201(84)90058-X
- 1985** ●NRA Hoffman and DP McKenzie. "The destruction of geochemical heterogeneities by differential fluid motions during mantle convection". In: *Geophysical Journal International* 82.2 (1985), pp. 163–206. DOI: 10.1111/j.1365-246X.1985.tb05134.x
- Ulrich R Christensen and David A Yuen. "Layered convection induced by phase transitions". In: *Journal of Geophysical Research: Solid Earth* 90.B12 (1985), pp. 10291–10300. DOI: 10.1029/JB090iB12p10291
- Satoru Honda. "Thermal structure beneath Tohoku, northeast Japan". In: *Tectonophysics* 112.1-4 (1985), pp. 69–102. DOI: 10.1016/0040-1951(85)90173-8
- 1986** ●Michael Gurnis and Geoffrey F Davies. "Mixing in numerical models of mantle convection incorporating plate kinematics". In: *Journal of Geophysical Research: Solid Earth* 91.B6 (1986), pp. 6375–6395. DOI: 10.1029/JB091iB06p06375
- Claire Harvey Craig and Dan McKenzie. "The existence of a thin low-viscosity layer beneath the lithosphere". In: *Earth and Planetary Science Letters* 78.4 (1986), pp. 420–426
- Geoffrey F Davies and Michael Gurnis. "Interaction of mantle dregs with convection: Lateral heterogeneity at the core-mantle boundary". In: *Geophysical Research Letters* 13.13 (1986), pp. 1517–1520
- Geoffrey F Davies. "Mantle convection under simulated plates: effects of heating modes and ridge and trench migration, and implications for the core-mantle boundary, bathymetry, the geoid and Benioff zones". In: *Geophys. J. R. astr. Soc* 84.1 (1986), pp. 153–183
- W Roger Buck. "Small-scale convection induced by passive rifting: the cause for uplift of rift shoulders". In: *Earth and Planetary Science Letters* 77.3-4 (1986), pp. 362–372. DOI: 10.1016/0012-821X(86)90146-9
- 1987** ●Ulrich R Christensen. "Time-dependent convection in elongated Rayleigh-Benard cells". In: *Geophysical Research Letters* 14.3 (1987), pp. 220–223. DOI: 10.1029/GL014i003p00220
- Harro Schmeling. "On the relation between initial conditions and late stages of Rayleigh-Taylor instabilities". In: *Tectonophysics* 133.1-2 (1987), pp. 65–80. DOI: 10.1016/0040-1951(87)90281-2
- J. Revenaugh and B. Parsons. "Dynamic topography and gravity anomalies for fluid layers whose viscosity varies exponentially with depth". In: *Geophysical Journal of the Royal Astronomical Society* 90.2 (1987), pp. 349–368. DOI: 10.1111/j.1365-246X.1987.tb00731.x
- Charles G Speziale. "On the advantages of the vorticity-velocity formulation of the equations of fluid dynamics". In: *J. Comp. Phys.* 73 (1987), pp. 476–480
- Marc Spiegelman and Dan McKenzie. "Simple 2-D models for melt extraction at mid-ocean ridges and island arcs". In: *Earth and Planetary Science Letters* 83.1-4 (1987), pp. 137–152
- U. Hansen and D.A. Yuen. "Evolutionary structures in double-diffusive convection in magma chambers". In: *Geophysical Research Letters* 14.11 (1987), pp. 1099–1102. DOI: 10.1029/GL014i011p01099
- 1988** ●MGG Foreman and AF Bennett. "On no-slip boundary conditions for the incompressible Navier-Stokes equations". In: *Dynamics of atmospheres and oceans* 12.1 (1988), pp. 47–70. DOI: 10.1016/0377-0265(88)90014-0
- A.P. Vincent and D.A. Yuen. "Thermal attractor in chaotic convection with high-Prandtl-number fluids". In: *Physical Review A* 38.1 (1988), pp. 328–334. DOI: 10.1103/PhysRevA.38.328
- Max D Gunzburger and Janet S Peterson. "On finite element approximations of the streamfunction-vorticity and velocity-vorticity equations". In: *International journal for numerical methods in fluids* 8.10 (1988), pp. 1229–1240. DOI: 10.1002/fld.1650081010
- G.F. Davies. "Role of the lithosphere in mantle convection". In: *Journal of Geophysical Research: Solid Earth* 93.B9 (1988), pp. 10451–10466. DOI: 10.1029/JB093iB09p10451
- Ulrich Hansen and Adolf Ebel. "Time-dependent thermal convection—a possible explanation for a multiscale flow in the Earth's mantle". In: *Geophysical Journal International* 94.2 (1988), pp. 181–191. DOI: 10.1111/j.1365-246X.1988.tb05895.x
- 1989** ●U.R. Christensen and D.A. Yuen. "Time-dependent convection with non-Newtonian viscosity". In: *Journal of Geophysical Research* 94.B1 (1989), pp. 814–820. DOI: 10.1029/JB094iB01p00814
- G.F. Davies. "Mantle convection model with a dynamic plate: topography, heat flow and gravity anomalies". In: *Geophysical Journal International* 98.3 (1989), pp. 461–464. DOI: 10.1111/j.1365-246X.1989.tb02283.x
- S.A. Weinstein, P.L. Olson, and D.A. Yuen. "Time-dependent large aspect-ratio thermal convection in the earth's mantle". In: *Geophysical & Astrophysical Fluid Dynamics* 47.1-4 (1989), pp. 157–197. DOI: 10.1080/03091928908221820
- Ulrich Hansen and David A Yuen. "Dynamical influences from thermal-chemical instabilities at the core-mantle boundary". In: *Geophysical Research Letters* 16.7 (1989), pp. 629–632. DOI: 10.1029/GL016i007p00629
- Harro Schmeling. "Compressible convection with constant and variable viscosity: The effect on slab formation, geoid, and topography". In: *Journal of Geophysical Research: Solid Earth* 94.B9 (1989), pp. 12463–12481
- Volker Steinbach, Ulrich Hansen, and Adolf Ebel. "Compressible convection in the earth's mantle: a comparison of different approaches". In: *Geophysical Research Letters* 16.7 (1989), pp. 633–636. DOI: 10.1029/GL016i007p00633
- P. Machetel and D.A. Yuen. "Penetrative convective flows induced by internal heating and mantle compressibility". In: *Journal of Geophysical Research* 94.B8 (1989), pp. 10, 609–10, 626. DOI: 10.1029/JB094iB08p10609
- 1990** ●GA Houseman. "The thermal structure of mantle plumes: axisymmetric or triple-junction?" In: *Geophysical Journal International* 102.1 (1990), pp. 15–24. DOI: 10.1111/j.1365-246X.1990.tb00527.x
- GA Houseman. "Boundary conditions and efficient solution algorithms for the potential function formulation of the 3-D viscous flow equations". In: *Geophysical Journal International* 100.1 (1990), pp. 33–38. DOI: 10.1111/j.1365-246X.1990.tb04565.x
- LP Solheim and WR Peltier. "Heat transfer and the onset of chaos in a spherical, axisymmetric, anelastic model of whole mantle convection". In: *Geophysical & Astrophysical Fluid Dynamics* 53.4 (1990), pp. 205–255
- B.J. Travis et al. "A benchmark comparison of numerical methods for infinite Prandtl number thermal convection in two-dimensional Cartesian geometry". In: *Geophysical & Astrophysical Fluid Dynamics* 55.3-4 (1990), pp. 137–160
- Bryan Travis, Peter Olson, and Gerald Schubert. "The transition from two-dimensional to three-dimensional planforms in infinite-Prandtl-number thermal convection". In: *Journal of Fluid Mechanics* 216 (1990), pp. 71–91. DOI: 10.1017/S0022112090000349



- U. Hansen and D.A. Yuen. “Nonlinear physics of double-diffusive convection in geological systems”. In: *Earth Science Reviews* 29.1-4 (1990), pp. 385–399. DOI: 10.1016/0012-8252(90)90050-6
- 1991** •William SD Wilcock and JA Whitehead. “The Rayleigh-Taylor instability of an embedded layer of low-viscosity fluid”. In: *Journal of Geophysical Research: Solid Earth* 96.B7 (1991), pp. 12193–12200. DOI: 10.1029/91JB00339
- Philip M Gresho. “Some current CFD issues relevant to the incompressible Navier-Stokes equations”. In: *Computer Methods in Applied Mechanics and Engineering* 87.2-3 (1991), pp. 201–252. DOI: 10.1016/0045-7825(91)90006-R
- A.V. Malevsky and D.A. Yuen. “Characteristics-based methods applied to infinite Prandtl number thermal convection in the hard turbulent regime”. In: *Physics of Fluids A* 3.9 (1991), pp. 2105–2115. DOI: 10.1063/1.857893
- Ulrich Hansen, David A Yuen, and Sherri E Kroening. “Effects of depth-dependent thermal expansivity on mantle circulations and lateral thermal anomalies”. In: *Geophysical Research Letters* 18.7 (1991), pp. 1261–1264. DOI: 10.1029/91GL01288
- GT Jarvis. “Two-dimensional numerical models of mantle convection”. In: *Advances in geophysics*. Vol. 33. 1991, pp. 1–80. DOI: 10.1016/S0065-2687(08)60440-9
- 1992** •R.R. Christensen. “An Eulerian technique for thermomechanical modeling of lithospheric extension”. In: *J. Geophys. Res.* 97.B2 (1992), pp. 2015–2036. DOI: 10.1029/91JB02642
- U. Hansen, D.A. Yuen, and S.E. Kroening. “Mass and Heat Transport in Strongly Time-Dependent Thermal Convection at Infinite Prandtl Number”. In: *Geophysical & Astrophysical Fluid Dynamics* 63.1-4 (1992), pp. 67–89. DOI: 10.1080/03091929208228278
- U. Hansen, D.A. Yuen, and A.V. Malevsky. “Comparison of steady-state and strongly chaotic thermal convection at high Rayleigh number”. In: *Physical Review A* 46.8 (1992), pp. 4742–4754. DOI: 10.1103/PhysRevA.46.4742
- Andrei V Malevsky, David A Yuen, and LM Weyer. “Viscosity and thermal fields associated with strongly chaotic non-Newtonian thermal convection”. In: *Geophysical research letters* 19.2 (1992), pp. 127–130
- A. Poliakov and Y. Podlachikov. “Diapirism and topography”. In: *Geophys. J. Int.* 109 (1992), pp. 553–564. DOI: 10.1111/j.1365-246X.1992.tb00117.x
- P. van Keken, D.A. Yuen, and A. van den Berg. “Pulsating diapiric flows: Consequences of vertical variations in mantle creep laws”. In: *Earth Planet. Sci. Lett.* 112 (1992), pp. 179–194. DOI: 10.1016/0012-821X(92)90015-N
- A.V. Malevsky and D.A. Yuen. “Strongly chaotic non-newtonian mantle convection”. In: *Geophysical & Astrophysical Fluid Dynamics* 65.1-4 (1992), pp. 149–171. DOI: 10.1080/03091929208225244
- 1993** •Ulrich Hansen, David A Yuen, SE Kroening, and TB Larsen. “Dynamical consequences of depth-dependent thermal expansivity and viscosity on mantle circulations and thermal structure”. In: *Physics of the earth and planetary interiors* 77.3-4 (1993), pp. 205–223. DOI: 10.1016/0031-9201(93)90099-U
- U. Hansen and D.A. Yuen. “High Rayleigh number regime of temperature-dependent viscosity convection and the Earth’s early thermal history”. In: *Geophysical Research Letters* 20.20 (1993), pp. 2191–2194. DOI: 10.1029/93GL02416
- A.V. Malevsky and D.A. Yuen. “Plume structures in the hard-turbulent regime of three-dimensional infinite Prandtl number convection”. In: *Geophysical Research Letters* 20.5 (1993), pp. 383–386. DOI: 10.1029/93GL00293
- Arie P van den Berg, Peter E van Keken, and David A Yuen. “The effects of a composite non-Newtonian and Newtonian rheology on mantle convection”. In: *Geophysical Journal International* 115.1 (1993), pp. 62–78. DOI: 10.1111/j.1365-246X.1993.tb05588.x
- David Bercovici. “A simple model of plate generation from mantle flow”. In: *Geophysical Journal International* 114.3 (1993), pp. 635–650. DOI: 10.1111/j.1365-246X.1993.tb06993.x
- G Guj and F Stella. “A vorticity-velocity method for the numerical solution of 3D incompressible flows”. In: *Journal of Computational Physics* 106.2 (1993), pp. 286–298
- G.T. Jarvis. “Effects of curvature on two-dimensional models of mantle convection: cylindrical polar coordinates”. In: *J. Geophys. Res.* 98.B3 (1993), pp. 4477–4485. DOI: 10.1029/92JB02117
- Julian P Lowman and Gary T Jarvis. “Mantle convection flow reversals due to continental collisions”. In: *Geophysical Research Letters* 20.19 (1993), pp. 2087–2090. DOI: 10.1029/93GL02047
- Volker Steinbach, David A Yuen, and Wuling Zhao. “Instabilities from phase transitions and the timescales of mantle thermal evolution”. In: *Geophysical research letters* 20.12 (1993), pp. 1119–1122. DOI: 10.1029/93GL01243
- P.E. van Keken, C.J. Spiers, A.P. van den Berg, and E.J. Muzert. “The effective viscosity of rocksalt: implementation of steady-state creep laws in numerical models of salt diapirism”. In: *Tectonophysics* 225 (1993), pp. 457–476
- P.E. van Keken, D.A. Yuen, and A.P. van den Berg. “The effects of shallow rheological boundaries in the upper mantle on inducing shorter time scales of diapiric flows”. In: *Geophysical Research Letters* 20.18 (1993), pp. 1927–1930. DOI: 10.1029/93GL01768
- Yoshitsugu Furukawa. “Depth of the decoupling plate interface and thermal structure under arcs”. In: *Journal of Geophysical Research: Solid Earth* 98.B11 (1993), pp. 20005–20013
- 1994** •F.H. Busse et al. “3D convection at infinite Prandtl number in Cartesian geometry - a benchmark comparison”. In: *Geophys. Astrophys. Fluid Dynamics* 75.1 (1994), pp. 39–59. DOI: 10.1080/03091929408203646
- Larry P Solheim and WR Peltier. “Avalanche effects in phase transition modulated thermal convection: A model of Earth’s mantle”. In: *Journal of Geophysical Research: Solid Earth* 99.B4 (1994), pp. 6997–7018
- Tomoeaki Nakakuki, Hiroki Sato, and Hiromi Fujimoto. “Interaction of the upwelling plume with the phase and chemical boundary at the 670 km discontinuity: Effects of temperature-dependent viscosity”. In: *Earth Planet. Sci. Lett.* 121.3–4 (1994), pp. 369–384. DOI: 10.1016/0012-821X(94)90078-7
- B.A. Buffet, C.W. Gable, and R.J. R.J. O’Connell. “Linear stability of a layered fluid with mobile surface plates”. In: *J. Geophys. Res.* 99(B10) (1994), pp. 19, 885–19, 900
- J. Schmalzl and U. Hansen. “Mixing the Earth’s mantle by thermal convection: A scale dependent phenomenon”. In: *Geophysical Research Letters* 21.11 (1994), pp. 987–990. DOI: 10.1029/94GL00049
- N.J. Vlaar, P.E. van Keken, and A.P. van den Berg. “Cooling of the Earth in the Archean: Consequences of pressure-release melting in a hotter mantle”. In: *Earth Planet. Sci. Lett.* 121 (1994), pp. 1–18. DOI: 10.1016/0012-821X(94)90028-0
- Bryan Travis and Peter Olson. “Convection with internal heat sources and thermal turbulence in the Earth’s mantle”. In: *Geophysical Journal International* 118.1 (1994), pp. 1–19. DOI: 10.1111/j.1365-246X.1994.tb04671.x
- David A Yuen, DM Reuteler, S Balachandar, V Steinbach, AV Malevsky, and JJ Smedsmo. “Various influences on three-dimensional mantle convection with phase transitions”. In: *Physics of the Earth and Planetary interiors* 86.1-3 (1994), pp. 185–203. DOI: 10.1016/0031-9201(94)05068-6
- S. Zhong and M. Gurnis. “Role of plates and temperature-dependent viscosity in phase change dynamics”. In: *Journal of Geophysical Research* 99.B8 (1994), p. 15903. DOI: 10.1029/94JB00545
- U. Hansen and D.A. Yuen. “Effects of depth-dependent thermal expansivity on the interaction of thermal-chemical plumes with



a compositional boundary". In: *Physics of the Earth and Planetary Interiors* 86.1-3 (1994), pp. 205–221. DOI: 10.1016/0031-9201(94)05069-4

- 1995** •J. Schmalzl, G.A. Houseman, and U. Hansen. "Mixing properties of three-dimensional (3-D) stationary convection". In: *Physics of Fluids* 7.5 (1995), pp. 1027–1033. DOI: 10.1063/1.868614
- D. Bittner and H. Schmeling. "Numerical modelling of melting processes and induced diapirism in the lower crust". In: *Geophy. J. Int.* 123 (1995), pp. 59–70. DOI: 10.1111/j.1365-246X.1995.tb06661.x
- NM Ribe, UR Christensen, and J Theissing. "The dynamics of plume-ridge interaction, 1: Ridge-centered plumes". In: *Earth and Planetary Science Letters* 134.1-2 (1995), pp. 155–168
- Volker Steinbach and David A Yuen. "The non-adiabatic nature of mantle convection as revealed by passive tracers". In: *Earth and Planetary Science Letters* 136.3-4 (1995), pp. 241–250. DOI: 10.1016/0012-821X(95)00166-A
- SCR Dennis and JD Hudson. "Methods of solution of the velocity-vorticity formulation of the Navier-Stokes equations". In: *Journal of Computational Physics* 122.2 (1995), pp. 300–306
- Julian P Lowman and Gary T Jarvis. "Mantle convection models of continental collision and breakup incorporating finite thickness plates". In: *Physics of the Earth and Planetary Interiors* 88.1 (1995), pp. 53–68
- VS Solomatov. "Scaling of temperature-and stress-dependent viscosity convection". In: *Physics of Fluids* 7.2 (1995), pp. 266–274
- 1996** •T.B. Larsen, D.A. Yuen, A.V. Malevsky, and J.L. Smedsmo. "Dynamics of strongly time-dependent convection with non-Newtonian temperature-dependent viscosity". In: *Physics of the Earth and Planetary Interiors* 94.1-2 (1996), pp. 75–103. DOI: 10.1016/0031-9201(95)03082-4
- Terence D Barr and Gregory A Houseman. "Deformation fields around a fault embedded in a non-linear ductile medium". In: *Geophysical Journal International* 125.2 (1996), pp. 473–490. DOI: 10.1111/j.1365-246X.1996.tb00012.x
- David Bercovici. "Plate generation in a simple model of lithosphere-mantle flow with dynamic self-lubrication". In: *Earth and Planetary Science Letters* 144.1-2 (1996), pp. 41–51
- D Breuer, H Zhou, David A Yuen, and T Spohn. "Phase transitions in the Martian mantle: Implications for the planet's volcanic history". In: *Journal of Geophysical Research: Planets* 101.E3 (1996), pp. 7531–7542. DOI: 10.1029/96JE00117
- Ulrich R Christensen. "The influence of trench migration on slab penetration into the lower mantle". In: *Earth and Planetary Science Letters* 140.1-4 (1996), pp. 27–39. DOI: 10.1016/0012-821X(96)00023-4
- A. Lenardic and W. M. Kaula. "Near-surface thermal/chemical boundary layer convection at infinite Prandtl number: two-dimensional numerical experiments". In: *Geophysical Journal International* 126.3 (1996), pp. 689–711. DOI: 10.1111/j.1365-246X.1996.tb04698.x
- Julian P Lowman and Gary T Jarvis. "Continental collisions in wide aspect ratio and high Rayleigh number two-dimensional mantle convection models". In: *Journal of Geophysical Research: Solid Earth* 101.B11 (1996), pp. 25485–25497. DOI: 10.1029/96JB02568
- J.X. Mitrovica, R.N. Pysklywec, C. Beaumont, and A. Rutt. "The Devonian to Permian sedimentation of the Russian platform: An example of subduction-controlled long-wavelength tilting of continents". In: *Journal of Geodynamics* 22.1-2 (1996), pp. 79–96. DOI: 10.1016/0264-3707(96)00008-7
- NM Ribe. "The dynamics of plume-ridge interaction - II. Off-ridge plumes". In: *Journal of Geophysical Research: Solid Earth* 101.B7 (1996), pp. 16195–16204
- J. Schmalzl, G.A. Houseman, and U. Hansen. "Mixing in vigorous, time-dependent three-dimensional convection and application to Earth's mantle". In: *Journal of Geophysical Research B: Solid Earth* 101.B10 (1996), pp. 21847–21858
- Arie P van den Berg and David A Yuen. "Is the lower-mantle rheology Newtonian today?" In: *Geophysical research letters* 23.16 (1996), pp. 2033–2036. DOI: 10.1029/96GL02065
- S. Zhong, M. Gurnis, and L. Moresi. "Free-surface formulation of mantle convection-I. Basic theory and application to plumes". In: *Geophysical Journal International* 127.3 (1996), pp. 708–718. DOI: 10.1111/j.1365-246X.1996.tb04049.x
- Shijie Zhong. "Analytic solutions for Stokes' flow with lateral variations in viscosity". In: *Geophys. J. Int.* 124.1 (1996), pp. 18–28. DOI: 10.1111/j.1365-246X.1996.tb06349.x
- 1997** •Tine B Larsen, David A Yuen, Jiří Moser, and Bengt Fornberg. "A high-order finite-difference method applied to large Rayleigh number mantle convection". In: *Geophysical & Astrophysical Fluid Dynamics* 84.1-2 (1997), pp. 53–83. DOI: 10.1080/03091929708208973
- D. Olbertz, M.J.R. Wortel, and U. Hansen. "Trench migration and subduction zone geometry". In: *Geophysical Research Letters* 24.3 (1997), pp. 221–224. DOI: 10.1029/96GL03971
- Uwe Walzer and Roland Hendel. "Tectonic episodicity and convective feedback mechanisms". In: *Physics of the earth and planetary interiors* 100.1-4 (1997), pp. 167–188. DOI: 10.1016/S0031-9201(96)03238-4
- Doris Breuer, Dave A Yuen, and Tilman Spohn. "Phase transitions in the Martian mantle: Implications for partially layered convection". In: *Earth and Planetary Science Letters* 148.3-4 (1997), pp. 457–469. DOI: 10.1016/S0012-821X(97)00049-6
- W. DeLandro-Clarke and G.T. Jarvis. "Numerical models of mantle convection with secular cooling". In: *Geophysical Journal International* 129.1 (1997), pp. 183–193
- T.B. Larsen and D.A. Yuen. "Ultrafast upwelling bursting through the upper mantle". In: *Earth and Planetary Science Letters* 146.3-4 (1997), pp. 393–399. DOI: 10.1016/S0012-821X(96)00247-6
- T.B. Larsen and D.A. Yuen. "Fast plumeheads: Temperature-dependent versus non-Newtonian rheology". In: *Geophysical Research Letters* 24.16 (1997), pp. 1995–1998. DOI: 10.1029/97GL01886
- P.E. van Keken, S.D. King, H. Schmeling, U.R. Christensen, D. Neumeister, and M.-P. Doin. "A comparison of methods for the modeling of thermochemical convection". In: *J. Geophys. Res.* 102.B10 (1997), pp. 22, 477–22, 495
- A.P. Van Den Berg and D.A. Yuen. "The role of shear heating in lubricating mantle flow". In: *Earth and Planetary Science Letters* 151.1-2 (1997), pp. 33–42. DOI: 10.1016/S0012-821X(97)00110-6
- Uwe Walzer and Roland Hendel. "Tectonic episodicity and convective feedback mechanisms". In: *Physics of the earth and planetary interiors* 100.1-4 (1997), pp. 167–188. DOI: 10.1016/S0031-9201(96)03238-4
- 1998** •Helmut Harder. "Phase transitions and the three-dimensional planform of thermal convection in the Martian mantle". In: *Journal of Geophysical Research: Planets* 103.E7 (1998), pp. 16775–16797. DOI: 10.1029/98JE01543
- B Schott and H Schmeling. "Delamination and detachment of a lithospheric root". In: *Tectonophysics* 296.3-4 (1998), pp. 225–247. DOI: 10.1016/S0040-1951(98)00154-1
- JA Gil and Maria José Jurado. "Geological interpretation and numerical modelling of salt movement in the Barbastro–Balaguer anticline, southern Pyrenees". In: *Tectonophysics* 293.3-4 (1998), pp. 141–155

- 1999** •Uwe Walzer and Roland Hendel. “A new convection–fractionation model for the evolution of the principal geochemical reservoirs of the Earth’s mantle”. In: *Physics of the Earth and Planetary Interiors* 112.3-4 (1999), pp. 211–256. DOI: 10.1016/S0031-9201(99)00035-7
- C Sotin and S Labrosse. “Three-dimensional thermal convection in an iso-viscous, infinite Prandtl number fluid heated from within and from below: applications to the transfer of heat through planetary mantles”. In: *Physics of the Earth and Planetary Interiors* 112.3-4 (1999), pp. 171–190. DOI: 10.1016/S0031-9201(99)00004-7
- 2008** •Juan-Luis Valera, Ana-María Negredo, and Antonio Villaseñor. “Asymmetric delamination and convective removal numerical modeling: comparison with evolutionary models for the Alboran Sea region”. In: *Pure appl. geophys.* 165 (2008), pp. 1683–1706. DOI: 10.1007/s00024-008-0395-8
- 2011** •JL Valera, Ana M Negredo, and Ivone Jiménez-Munt. “Deep and near-surface consequences of root removal by asymmetric continental delamination”. In: *Tectonophysics* 502.1-2 (2011), pp. 257–265. DOI: 10.1016/j.tecto.2010.04.002

# Chapter 14

## Heat Transfer & convection in a porous medium

I am by no means an expert when it comes to porous media. I hope to revisit the topic regularly in the coming years and improve this section. Any help or comment welcome.

QUESTOIN: What is head?

### 14.0.1 Darcy's law for groundwater movement

[Taken from MODFLOW manual] The three-dimensional movement of groundwater of constant density through porous earth material is described by Darcy's Law:

$$\vec{q} = -\mathbf{K} \cdot \vec{\nabla} h$$

where  $\vec{q}$  is a vector of specific discharge (L/T), or fluid-flux vector,  $\mathbf{K}$  is the hydraulic-conductivity tensor (L/T),  $h$  is the potentiometric head (L).

When combined with a water balance on a small control volume, Darcy's Law leads to a partial-differential equation that describes the distribution of hydraulic head:

$$SS \frac{\partial h}{\partial t} = -\vec{\nabla} \cdot \vec{q} + Q'_s = \vec{\nabla} \cdot (\mathbf{K} \cdot \vec{\nabla} h) + Q'_s \quad (14.1)$$

where  $Q'_s$  is a volumetric flux per unit volume representing sources and sinks of water, with  $Q'_s$  being negative for flow out of the groundwater system, and  $Q'_s$  being positive for flow into the system ( $T^{-1}$ ).  $SS$  is the specific storage of the porous material ( $L^{-1}$ ); and  $t$  is time ( $T$ ).

Eq. (14.1) describes transient groundwater flow in a heterogeneous and anisotropic medium. This equation, together with specification of flow and head conditions at the boundaries of an aquifer system and specification of initial-head conditions, constitutes a mathematical representation of a groundwater flow system.

QUESTION: why no gravity in there ?

At this stage we have to acknowledge similarities with the heat equation with the heat flux  $\vec{q}$  being given by

$$\vec{q} = -k \vec{\nabla} T$$

where  $k$  is the heat conductivity (which can also be a tensor if the medium is anisotropic) and

$$\rho C_p \frac{\partial h}{\partial t} = -\vec{\nabla} \cdot \vec{q} + H = \vec{\nabla} \cdot (k \vec{\nabla} T) + H$$

This means that in the absence of other physics in the system, we know how to solve the groundwater equation, as explained in Chapter 6.

## 14.0.2 The equations of non-isothermal fluid flow in a porous medium

A porous medium is a material containing pores. These pores can be filled with a gas or a fluid. Often the pore space forms a network which allows fluids to pass through.

The equations under consideration are the following:

- Darcy's law is an equation that describes the flow of a fluid through a porous medium. The law was formulated by Henry Darcy<sup>1</sup> based on results of experiments on the flow of water through beds of sand:

$$\vec{v} = -\frac{\mathbf{K}}{\eta}(\vec{\nabla}p + \rho_f \vec{g}) \quad (14.2)$$

- mass conservation (incompressibility condition):

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (14.3)$$

- heat transport: Usually it is a good approximation to assume that the solid and fluid phases are in local thermal equilibrium (LTE) but there are situations, such as highly transient problems and some steady-state problems, where this is not so. Now this is commonly referred to as local thermal nonequilibrium (LTNE). If one wishes to allow for heat transfer between solid and fluid (that is, one no longer has local thermal equilibrium), then the equations are

$$(1 - \phi)(\rho C_p)_s \frac{\partial T_s}{\partial t} = (1 - \phi)\vec{\nabla} \cdot (k_s \vec{\nabla} T_s) + (1 - \phi)q_s + h(T_f - T_s) \quad (14.4)$$

$$\phi(\rho C_p)_f \frac{\partial T_f}{\partial t} + (\rho C_p)_f \vec{v} \cdot \vec{\nabla} T_f = \phi \vec{\nabla} \cdot (k_f \vec{\nabla} T_f) + \phi q_f + h(T_s - T_f) \quad (14.5)$$

When it is assumed that there is local thermal equilibrium then  $T_f = T_s = T$  (Section 2.1 of Nield & Bejan's book). Then one can add the equations together and obtain

$$(\rho C_p)_m \frac{\partial T}{\partial t} + (\rho C_p)_f \vec{v} \cdot \vec{\nabla} T = \vec{\nabla} \cdot (k_m \vec{\nabla} T) + q_m$$

with

$$(\rho C_p)_m = (1 - \phi)(\rho C_p)_s + \phi(\rho C_p)_f \quad (14.6)$$

$$k_m = (1 - \phi)k_s + \phi k_f \quad (14.7)$$

$$q_m = (1 - \phi)q_s + \phi q_f \quad (14.8)$$

- linearised equation of state in the form of the Oberbeck-Boussinesq approximation (Section 2.3 of Nield & Bejan's book):

$$\rho = \rho_0(1 - \alpha(T_f - T_0))$$

In the equations above the subscript  $f$  stands for “fluid”, and the subscript  $s$  for “solid”.  $\vec{v}$  is the velocity ( $\text{m s}^{-1}$ ),  $p$  is the pressure (Pa),  $\vec{g}$  is the gravitational acceleration ( $\text{m s}^{-2}$ ),  $\phi$  is the porosity,  $\mathbf{K}$  the permeability tensor ( $\text{m}^2$ ),  $\rho$  the mass density ( $\text{kg m}^{-3}$ ),  $T$  the temperature (K),  $h$  the coefficient of heat transfer between solid and fluid (unit?),  $\alpha$  the coefficient of thermal expansion,  $k$  the heat conductivity,  $C_p$  the heat capacity, and  $\eta$  the dynamic viscosity (Pa s).

<sup>1</sup>[https://en.wikipedia.org/wiki/Henry\\_Darcy](https://en.wikipedia.org/wiki/Henry_Darcy)

For the case of an isotropic medium the permeability is a scalar, i.e.  $\mathbf{K} = K\mathbf{1}$  so that

$$\vec{\mathbf{v}} = -\frac{K}{\eta}(\vec{\nabla}p + \rho_f \vec{g}) \quad (14.9)$$

Values of  $K$  for natural materials vary widely. Typical values for soils, in terms of the unit  $\text{m}^2$ , are: clean gravel  $10^{-7} - 10^{-9}$ , clean sand  $10^{-9} - 10^{-12}$ , peat  $10^{-11} - 10^{-13}$ , stratified clay  $10^{-13} - 10^{-16}$ , and unweathered clay  $10^{-16} - 10^{-20}$ . Workers concerned with geophysics often use as a unit of permeability the Darcy, which equals  $0.987 \cdot 10^{-12} \text{m}^2$ .

### 14.0.3 Weak form and discretisation

We wish to solve the equations of the previous section with the Finite Element method. There are four unknown fields:  $\vec{\mathbf{v}} = (u, v)$ ,  $T_s$ ,  $T_f$ ,  $p$  which we conveniently downsize to three assuming local thermal equilibrium (LTE), i.e.  $T = T_f = T_s$ .

The Cartesian domain is partitioned in non-overlapping elements. In each element the fields can be expressed as follows:

$$\mathbf{v}_x(x, y) = \sum_{i=1}^{m_v} \mathcal{N}_i^v(x, y) \mathbf{v}_{x,i} = \vec{\mathcal{N}}^v \cdot \vec{\mathbf{v}}_x \quad (14.10)$$

$$\mathbf{v}_y(x, y) = \sum_{i=1}^{m_v} \mathcal{N}_i^v(x, y) \mathbf{v}_{y,i} = \vec{\mathcal{N}}^v \cdot \vec{\mathbf{v}}_y \quad (14.11)$$

$$p(x, y) = \sum_{i=1}^{m_p} \mathcal{N}_i^p(x, y) p_i = \vec{\mathcal{N}}^p \cdot \vec{\mathbf{p}} \quad (14.12)$$

$$T(x, y) = \sum_{i=1}^{m_T} \mathcal{N}_i^t(x, y) T_i = \vec{\mathcal{N}}^t \cdot \vec{\mathbf{T}} \quad (14.13)$$

with

$$\vec{\mathcal{N}}^v = (\mathcal{N}_1^v, \mathcal{N}_2^v, \dots, \mathcal{N}_{m_v}^v) \quad (14.14)$$

$$\vec{\mathcal{N}}^p = (\mathcal{N}_1^p, \mathcal{N}_2^p, \dots, \mathcal{N}_{m_p}^p) \quad (14.15)$$

$$\vec{\mathcal{N}}^T = (\mathcal{N}_1^T, \mathcal{N}_2^T, \dots, \mathcal{N}_{m_T}^T). \quad (14.16)$$

The heat transport poses no real problem and the topic has been treated in Section 6.3 so we will not repeat it here. When solving this equation we assume that the velocity of the advection term is known.

There are actually two approaches to solve the mass and momentum conservation equations. We wish to find the velocity and pressure fields assuming the temperature known.

### Mixed variable approach

We have two coupled equations (14.2) and (14.3):

$$-\eta \mathbf{K}^{-1} \cdot \vec{\mathbf{v}} - \vec{\nabla}p = \rho \vec{g} \quad (14.17)$$

$$-\vec{\nabla} \cdot \vec{\mathbf{v}} = 0 \quad (14.18)$$

or, defining  $\mathbf{L} = \eta \mathbf{K}^{-1}$ , these become in 2D Cartesian coordinates:

$$-L_{xx}\mathbf{v}_x - L_{xy}\mathbf{v}_y - \partial_x p = \rho g_x \quad (14.19)$$

$$-L_{yx}\mathbf{v}_x - L_{yy}\mathbf{v}_y - \partial_y p = \rho g_y \quad (14.20)$$

$$-\partial_x \mathbf{v}_x - \partial_y \mathbf{v}_y = 0 \quad (14.21)$$

Let us go through each line separately and establish its weak form:

$$\begin{aligned}
& - \int \vec{\mathcal{N}}^\nu L_{xx} \mathbf{v}_x dV - \int \vec{\mathcal{N}}^\nu L_{xy} \mathbf{v}_y dV - \int \vec{\mathcal{N}}^\nu \partial_x p dV = \int \vec{\mathcal{N}}^\nu \rho g_x \\
& \underbrace{\left( - \int \vec{\mathcal{N}}^\nu L_{xx} \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{N}_{xx}} \cdot \vec{\mathcal{V}}_x + \underbrace{\left( - \int \vec{\mathcal{N}}^\nu L_{xy} \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{N}_{xy}} \cdot \vec{\mathcal{V}}_y + \underbrace{\left( - \int \vec{\mathcal{N}}^\nu \partial_x \vec{\mathcal{N}}^p dV \right)}_{\mathbb{G}_x} \cdot \vec{P} = \underbrace{\int \vec{\mathcal{N}}^\nu \rho g_x}_{\vec{f}_x}
\end{aligned}$$

The second line yields

$$\begin{aligned}
& \underbrace{\left( - \int \vec{\mathcal{N}}^\nu L_{yx} \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{N}_{yx}} \cdot \vec{\mathcal{V}}_x + \underbrace{\left( - \int \vec{\mathcal{N}}^\nu L_{yy} \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{N}_{yy}} \cdot \vec{\mathcal{V}}_y + \underbrace{\left( - \int \vec{\mathcal{N}}^\nu \partial_y \vec{\mathcal{N}}^p dV \right)}_{\mathbb{G}_y} \cdot \vec{P} = \underbrace{\int \vec{\mathcal{N}}^\nu \rho g_y}_{\vec{f}_y}
\end{aligned}$$

The third line yields

$$\begin{aligned}
& - \int \vec{\mathcal{N}}^p (\partial_x \mathbf{v}_x + \partial_y \mathbf{v}_y) dV = \vec{0} \\
& \underbrace{\left( - \int \vec{\mathcal{N}}^p \partial_x \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{H}_x} \cdot \vec{\mathcal{V}}_x + \underbrace{\left( - \int \vec{\mathcal{N}}^p \partial_y \vec{\mathcal{N}}^\nu dV \right)}_{\mathbb{H}_y} \cdot \vec{\mathcal{V}}_y = \vec{0}
\end{aligned} \tag{14.22}$$

In the end:

$$\begin{pmatrix} \mathbb{N}_{xx} & \mathbb{N}_{xy} & \mathbb{G}_x \\ \mathbb{N}_{yx} & \mathbb{N}_{yy} & \mathbb{G}_y \\ \mathbb{H}_x & \mathbb{H}_y & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}}_x \\ \vec{\mathcal{V}}_y \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_x \\ \vec{f}_y \\ \vec{h} \end{pmatrix}$$

In the case of an isotropic material and an isoviscous fluid, we have

$$\mathbb{N}_{xx} = \mathbb{N}_{yy} = -\frac{\eta}{K} \int \vec{\mathcal{N}}^\nu \vec{\mathcal{N}}^\nu dV = -\frac{\eta}{K} \mathbb{M}^\nu$$

where  $\mathbb{M}^\nu$  is the velocity mass matrix, while  $\mathbb{N}_{xy} = \mathbb{N}_{yx} = 0$  so that we then solve

$$\begin{pmatrix} -\frac{\eta}{K} \mathbb{M}^\nu & 0 & \mathbb{G}_x \\ 0 & -\frac{\eta}{K} \mathbb{M}^\nu & \mathbb{G}_y \\ \mathbb{H}_x & \mathbb{H}_y & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}}_x \\ \vec{\mathcal{V}}_y \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_x \\ \vec{f}_y \\ \vec{h} \end{pmatrix}$$

**About the  $\mathbb{G}$  blocks** What is above has one major disadvantage: the  $\mathbb{G}$  blocks contain the biquadratic basis velocity and the derivatives of the bilinear pressure basis functions. We could be tempted to integrate  $\int \vec{\mathcal{N}}^\nu \partial_x p dV$  by parts in order to bring the space derivative on the velocity basis functions and thereby recover the  $\mathbb{H}$  blocks. However there is a surface term  $[\vec{\mathcal{N}}^\nu p]_\Gamma$  which I am not too sure what to do about... I have implemented this in the code (while disregarding the surface term) and found that the convergence was much much worse.

**Block scaling** As explained in Section 7.5.4, we need to scale the blocks so as to insure an accurate solution. Eq. (14.17) can be written

$$-\eta L^2 \mathbf{K}^{-1} \cdot \frac{\vec{\mathbf{v}}}{L^2} - \vec{\nabla} p = \rho \vec{g}$$

where  $L$  is a characteristic length. The term  $\vec{\mathbf{v}}/L^2$  has the same dimensions as the Laplacian of the velocity in the Stokes equations and we obviously find that the dimension of the  $\eta' = \eta L^2 \mathbf{K}^{-1}$  term

is one of viscosity. Following the reasoning in Section 7.5.4 the scaling coefficient for the  $\mathbb{G}$  and  $\mathbb{H}$  blocks is

$$\frac{\eta'}{L} = \frac{\eta L^2}{\tilde{K} L} = \frac{\eta L}{\tilde{K}}$$

where  $\tilde{K}$  is a representative quantity of the  $\mathbf{K}$  tensor. In our case, we find that taking  $L = h_x$  yields blocks which coefficient magnitudes are very well matched. After each elemental  $\mathbb{G}$  or  $\mathbb{H}$  block is built it is multiplied by the factor above and assembled. After the solve, the obtained pressure must then be multiplied by this factor to recover the proper magnitude.

## Second approach

Inserting Eq. (14.2) in Eq. (14.3) we obtain

$$\vec{\nabla} \cdot \left( -\frac{\mathbf{K}}{\eta} (\vec{\nabla} p + \rho \vec{g}) \right) = 0 \quad (14.23)$$

If we assume that the permeability tensor, the viscosity and the gravity are constant, then

$$\frac{\mathbf{K}}{\eta} (\Delta p + \vec{\nabla} \rho \cdot \vec{g}) = 0$$

or simply

$$\Delta p + \vec{\nabla} \rho \cdot \vec{g} = 0.$$

and we end up with a simple Poisson equation.

Let us now establish the weak form of Eq. (14.23) (without the above assumption):

$$\int \mathcal{N}_i^p \vec{\nabla} \cdot \left( \frac{\mathbf{K}}{\eta} (\vec{\nabla} p + \rho \vec{g}) \right) + \int \mathcal{N}_i^p \vec{\nabla} \cdot (\rho \vec{g}) = 0$$

After integration by parts + neglecting surface term (for now) we obtain

$$- \left( \int (\vec{\nabla} \mathcal{N}^p)^T \cdot \frac{\mathbf{K}}{\eta} \cdot \vec{\nabla} \mathcal{N} dV \right) \cdot \vec{P} + \int \mathcal{N}^p (\vec{\nabla} \rho \cdot \vec{g}) dV = 0$$

We denote here  $\mathbf{B} = \vec{\nabla} \mathcal{N}^p$  so

$$\left( \int \mathbf{B}^T \cdot \frac{\mathbf{K}}{\eta} \cdot \mathbf{B} dV \right) \cdot \vec{P} = \int \mathcal{N}^p (\vec{\nabla} \rho \cdot \vec{g}) dV$$

Using the equation of state, we find that

$$\vec{\nabla} \rho = -\alpha \rho_0 \vec{\nabla} T = -\alpha \rho_0 \vec{\nabla} (\vec{N}^\theta \cdot \vec{T})$$

Note that the discarded surface term are not trivial to formulate and that this approach does not easily allow to prescribe velocities anywhere in the domain since the velocity field is not solved for. In fact prescribing flow in the boundary is akin to pressure Neumann boundary conditions.

## 14.0.4 The equations in dimensionless form

This follows Palm *et al.* (1972) [971]. The field variables may conveniently be made dimensionless by choosing

$$h, \quad \Delta T, \quad \frac{\eta \kappa}{K}, \quad \frac{\kappa}{h}$$

as units of length, temperature, pressure, and velocity respectively.

Let us start with Eq. (14.9). Dividing each side by the reference velocity  $\kappa/h$  yields:

$$\vec{v}' = \frac{\vec{v}}{\kappa/h} = -\frac{Kh}{\eta\kappa}(\vec{\nabla}p + \rho\vec{g}) \quad (14.24)$$

Defining the dynamic pressure  $\tilde{p}$  as

$$\tilde{p} = p - p_{hydr} = p - \rho_0 g(L_y - y)$$

then  $\vec{\nabla}p = \vec{\nabla}\tilde{p} - \rho_0 g\vec{e}_y$  and introducing the temperature-dependence of the density in the equation yields

$$\vec{v}' = -\frac{Kh}{\eta\kappa}(\vec{\nabla}\tilde{p} - \rho_0 g\vec{e}_y + \rho_0(1 - \alpha(T - T_0)g\vec{e}_y) \quad (14.25)$$

$$= -\frac{Kh}{\eta\kappa}(\vec{\nabla}\tilde{p} - \rho_0\alpha(T - T_0)g\vec{e}_y) \quad (14.26)$$

We now define the dimensionless temperature  $T'$  as

$$T' = \frac{T - T_0}{\Delta T}$$

then

$$\vec{v}' = -\frac{Kh}{\eta\kappa}(\vec{\nabla}\tilde{p} - \rho_0\alpha T'\Delta T g\vec{e}_y) \quad (14.27)$$

$$= -\vec{\nabla}'\tilde{p}' + \frac{\alpha\rho_0 g K \Delta T h}{\kappa\eta}T'\vec{e}_y \quad (14.28)$$

$$= -\vec{\nabla}'\tilde{p}' + \text{Ra}T'\vec{e}_y \quad (14.29)$$

The other two equations (mass and energy conservation) are trivial, so dropping the primes, the (steady state form of the) equations takes the form

$$-\vec{\nabla}p + \text{Ra}T\vec{e}_z - \vec{v} = \vec{0} \quad (14.30)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (14.31)$$

$$\vec{v} \cdot \vec{\nabla}T = \vec{\nabla}^2 T \quad (14.32)$$

where Ra is a Rayleigh number defined by

$$\boxed{\text{Ra} = \frac{K\rho_0 g \alpha \Delta T h}{\kappa\eta}}$$

Following Palm *et al.* (1972) [971] and Kuo (1961) [735], we introduce  $T = T_0 - z + \theta$  where  $T_0$  is a dimensionless temperature, eliminating the pressure by applying the curl operator and applying the equation of continuity.

It is know that the curl of a gradient is zero, so  $\vec{\nabla} \times \vec{\nabla}p = 0$ . We then have

$$\vec{\nabla} \times [\text{Ra}T\vec{e}_z - \vec{v}] = \vec{0}$$

or,

$$-\text{Ra}\frac{\partial T}{\partial x} - \left(\frac{\partial v_x}{\partial z} - \frac{\partial v_z}{\partial x}\right) = 0$$



We take the partial derivative with respect to  $x$  of the above equation to obtain

$$-\text{Ra} \frac{\partial^2 T}{\partial x^2} - \left( \frac{\partial^2 \mathbf{v}_x}{\partial x z} - \frac{\partial \mathbf{v}_z}{\partial x^2} \right) = 0$$

using the incompressibility condition we have  $\partial_x \mathbf{v}_x = -\partial_z \mathbf{v}_z$  so

$$\begin{aligned} -\text{Ra} \frac{\partial^2 T}{\partial x^2} - \left( -\frac{\partial^2 \mathbf{v}_z}{\partial z^2} - \frac{\partial^2 \mathbf{v}_z}{\partial x^2} \right) &= 0 \\ -\text{Ra} \frac{\partial^2 T}{\partial x^2} + \Delta \mathbf{v}_z &= 0 \end{aligned}$$

We take the Laplacian of this equation:

$$-\text{Ra} \frac{\partial^2 \Delta T}{\partial x^2} + \Delta^2 \mathbf{v}_z = 0$$

Since  $\Delta T = \vec{\nabla} \cdot \vec{\nabla} T$  then

$$-\text{Ra} \frac{\partial^2 (\vec{\nabla} \cdot \vec{\nabla} T)}{\partial x^2} + \Delta^2 \mathbf{v}_z = 0$$

We have  $\vec{\nabla} T = -\vec{e}_z + \vec{\nabla} \theta$  so  $\vec{\nabla} \cdot \vec{\nabla} T = -\mathbf{v}_z + \vec{\nabla} \cdot \vec{\nabla} \theta$  and finally

$$-\text{Ra} \frac{\partial^2 (-\mathbf{v}_z + \vec{\nabla} \cdot \vec{\nabla} \theta)}{\partial x^2} + \Delta^2 \mathbf{v}_z = 0$$

or,

$$\vec{\nabla}^4 \mathbf{v}_z + \text{Ra} \frac{\partial^2 \mathbf{v}_z}{\partial x^2} = \frac{\partial^2 (\vec{\nabla} \cdot \vec{\nabla} \theta)}{\partial x^2}$$

Finally we recover Eqs. (2.11-13) of Palm *et al.* (1972) [971]:

$$\vec{\nabla}^4 \mathbf{v}_y + \text{Ra} \frac{\partial^2 \mathbf{v}_z}{\partial x^2} = \frac{\partial^2 (\vec{\nabla} \cdot \vec{\nabla} \theta)}{\partial x^2} \quad (14.33)$$

$$\vec{\nabla}^2 \theta + \mathbf{v}_z = \vec{\nabla} \cdot \vec{\nabla} \theta \quad (14.34)$$

$$\vec{\nabla} \cdot \vec{\nabla} = 0 \quad (14.35)$$

The boundary conditions are then  $\mathbf{v}_y = \theta = 0$  for  $y = 0, 1$ .

This formulation of the equation forms the basis of the convection benchmark in the coming section.

# Chapter 15

## Adjoint methods

w.i.p.

adjoint methods in geodynamics [188, 458, 590, 625, 780, 1341, 1338, 1369]. Also see `johnson_Notes on Adjoint Methods.pdf` and `bradley-PDE-constrained optimization and the adjoint method.pdf`

Advection-diffusion: [https://en.wikipedia.org/wiki/Adjoint\\_equation](https://en.wikipedia.org/wiki/Adjoint_equation)

Derivation of the adjoint poisson equation: <https://math.stackexchange.com/questions/2269111/derivation-of-the-adjoint-poisson-equation>

Inverting PDEs with adjoints (esp Poisson) <https://joelcfd.com/inverting-pdes-with-adjoints/>

Video: Introduction to the adjoint method [https://youtu.be/EybH\\_Q-QTZ8](https://youtu.be/EybH_Q-QTZ8)

Video: adjoint-based optimization <https://youtu.be/Yiz92Ekn7vU>

# Chapter 16

## Elasticity: physics, formulations and FEM

chapter17.tex

Let us start by clarifying notations:

| variable name               | symbol                        | unit              |
|-----------------------------|-------------------------------|-------------------|
| full stress tensor          | $\boldsymbol{\sigma}$         | Pa                |
| deviatoric stress tensor    | $\boldsymbol{\tau}$           | Pa                |
| strain tensor               | $\boldsymbol{\epsilon}$       | -                 |
| elastic strain tensor       | $\boldsymbol{\epsilon}_e$     | -                 |
| visco-plastic strain tensor | $\boldsymbol{\epsilon}_{vp}$  | -                 |
| total strain tensor         | $\boldsymbol{\epsilon}_T$     | -                 |
| strain rate tensor          | $\dot{\boldsymbol{\epsilon}}$ | s <sup>-1</sup>   |
| Lamé parameter              | $\lambda$                     | Pa                |
| Shear modulus               | $\mu$                         | Pa                |
| Bulk modulus                | $K$                           | Pa                |
| Poisson ratio               | $\nu$                         | -                 |
| Young's modulus             | $E$                           | Pa                |
| viscosity                   | $\eta$                        | Pa s              |
| displacement                | $\vec{u}$                     | m                 |
| velocity                    | $\vec{v}$                     | m s <sup>-1</sup> |

What follows is a compilation of various sources, such as the Becker & Kaus lecture notes [66], the excellent paper by Beuchert and Podladchikov [86] (2010), the syllabus of R. Hassani [XX] and various books such as Sadd [1094].

One will find in the literature either 'elasto-viscosity' or 'visco-elasticity'. In what follows I have adopted the former notation with the acronym EV.

Once the equations have been laid out, one must then make a fundamental choice with regards to the type of code/calculations in the case of elasto-viscous rheologies: will the primary variable be displacement  $\vec{u}$  or velocity  $\vec{v}$ ? The latter is the common approach in the geodynamics community. The vast majority of codes are fluid flow solvers, formulated in velocity and pressure. Elasticity is usually added to such codes way after they were first used/written (eg ELEFANT, ASPECT, ...). For a purely elastic code displacement is the meaningful primary variable since the stress is formulated as a function of strain (not strain rate).

## 16.1 Basic equations

The strong form of the PDE that governs force balance in a medium is given by

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{f} = \vec{0}$$

where  $\boldsymbol{\sigma}$  is the full stress tensor and  $\vec{f}$  is a body force (typically  $\rho\vec{g}$ ).

The stress tensor is related to the strain tensor through the generalised Hooke's law<sup>1</sup>:

$$\sigma_{ij} = \sum_{kl} C_{ijkl} \varepsilon_{kl} \quad \text{or} \quad \boldsymbol{\sigma} = \mathbf{C} : \boldsymbol{\varepsilon} \quad (16.1)$$

where  $\mathbf{C}$  is the fourth-order elastic tensor (which contains  $3^4 = 81$  coefficients). The strain tensor is related to the displacement  $\vec{u}$  as follows:

$$\boldsymbol{\varepsilon}(\vec{u}) = \frac{1}{2}(\vec{\nabla}\vec{u} + (\vec{\nabla}\vec{u})^T) \quad (16.2)$$

Due to the inherent symmetries of  $\boldsymbol{\sigma}$ ,  $\boldsymbol{\varepsilon}$ , and  $\mathbf{C}$ , only 21 elastic coefficients of the latter are independent. For isotropic linear media (which have the same physical properties in any direction),  $\mathbf{C}$  can be reduced to only two independent numbers (for example the bulk modulus  $K$  and the shear modulus<sup>2</sup>  $\mu$  that quantify the material's resistance to changes in volume and to shearing deformations, respectively). We find that

$$C_{ijkl} = \lambda \delta_{ij} \delta_{kl} + \mu (\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk})$$

so that Eq. (16.1) becomes:

$$\sigma_{ij} = \lambda \varepsilon_{kk} \delta_{ij} + 2\mu \varepsilon_{ij}$$

or

$$\begin{aligned} \boldsymbol{\sigma} &= \lambda \text{Tr}[\boldsymbol{\varepsilon}(\vec{u})] \mathbf{1} + 2\mu \boldsymbol{\varepsilon}(\vec{u}) \\ &= \lambda (\vec{\nabla} \cdot \vec{u}) \mathbf{1} + 2\mu \boldsymbol{\varepsilon}(\vec{u}) \end{aligned} \quad (16.3)$$

where  $\lambda$  is the Lamé parameter and  $\mu$  is the shear modulus. The term  $\vec{\nabla} \cdot \vec{u} = \text{Tr}[\boldsymbol{\varepsilon}(\vec{u})]$  is the isotropic dilation.

Very explicitly, and since the stress and strain tensors are symmetric, we have

$$\begin{aligned} \sigma_{xx} &= (\lambda + 2\mu) \varepsilon_{xx} + \lambda \varepsilon_{yy} + \lambda \varepsilon_{zz} \\ \sigma_{yy} &= \lambda \varepsilon_{xx} + (\lambda + 2\mu) \varepsilon_{yy} + \lambda \varepsilon_{zz} \\ \sigma_{zz} &= \lambda \varepsilon_{xx} + \lambda \varepsilon_{yy} + (\lambda + 2\mu) \varepsilon_{zz} \\ \sigma_{xy} &= 2\mu \varepsilon_{xy} \\ \sigma_{xz} &= 2\mu \varepsilon_{xz} \\ \sigma_{yz} &= 2\mu \varepsilon_{yz} \end{aligned} \quad (16.4)$$

This can be re-written in the 6-dimensional stress/strain space as

$$\underbrace{\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix}}_{\vec{\sigma}} = \underbrace{\begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix}}_{\mathbf{D}} \cdot \underbrace{\begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{zz} \\ 2\varepsilon_{xy} \\ 2\varepsilon_{xz} \\ 2\varepsilon_{yz} \end{pmatrix}}_{\vec{\varepsilon}} \quad (16.5)$$

<sup>1</sup>[https://en.wikipedia.org/wiki/Hooke's\\_law](https://en.wikipedia.org/wiki/Hooke's_law)

<sup>2</sup>It is also sometimes written  $G$

Let us define the matrices

$$\mathbf{\Lambda} = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \mathbf{\Xi} = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

so that

$$\mathbf{D} = \lambda \mathbf{\Xi} + \mu \mathbf{\Lambda}$$

If we define the Young's modulus  $E$  and the Poisson's ratio as

$$E = \frac{\mu(3\lambda + 2\mu)}{\lambda + \mu} \quad \text{or} \quad \nu = \frac{\lambda}{2(\lambda + \mu)} \quad (16.6)$$

Then

$$\begin{aligned} 1 - \nu &= 1 - \frac{\lambda}{2(\lambda + \mu)} = \frac{2(\lambda + \mu)}{2(\lambda + \mu)} - \frac{\lambda}{2(\lambda + \mu)} = \frac{\lambda + 2\mu}{2(\lambda + \mu)} \\ 1 + \nu &= 1 + \frac{\lambda}{2(\lambda + \mu)} = \frac{2(\lambda + \mu)}{2(\lambda + \mu)} + \frac{\lambda}{2(\lambda + \mu)} = \frac{3\lambda + 2\mu}{2(\lambda + \mu)} = \frac{E}{2\mu} \\ 1 - 2\nu &= 1 - 2\frac{\lambda}{2(\lambda + \mu)} = \frac{2(\lambda + \mu)}{2(\lambda + \mu)} - \frac{2\lambda}{2(\lambda + \mu)} = \frac{\mu}{\lambda + \mu} \\ \frac{E}{(1 + \nu)(1 - 2\nu)}(1 - \nu) &= E \frac{\lambda + 2\mu}{2(\lambda + \mu)} \cdot \frac{2\mu}{E} \cdot \frac{\lambda + \mu}{\mu} = \lambda + 2\mu \\ \frac{E}{(1 + \nu)(1 - 2\nu)} \frac{1}{2}(1 - 2\nu) &= E \frac{\mu}{\lambda + \mu} \cdot \frac{2\mu}{E} \cdot \frac{1}{2} \frac{\lambda + \mu}{\mu} = \mu \\ \frac{E}{(1 + \nu)(1 - 2\nu)} \nu &= E \frac{\lambda}{2(\lambda + \mu)} \cdot \frac{2\mu}{E} \cdot \frac{\lambda + \mu}{\mu} = \lambda \end{aligned}$$

and in the end:

$$\mathbf{D}_{3D} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} = \frac{E}{(1 + \nu)(1 - 2\nu)} \begin{pmatrix} 1 - \nu & \nu & \nu & 0 & 0 & 0 \\ \nu & 1 - \nu & \nu & 0 & 0 & 0 \\ \nu & \nu & 1 - \nu & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1 - 2\nu}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1 - 2\nu}{2} & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1 - 2\nu}{2} \end{pmatrix} \quad (16.7)$$

This matrix is the same as Eq. (3.6) on page 294 of Braess [128]. It is SPD for  $0 \leq \nu < 1/2$ .

In terms of the compliance matrix  $\mathbf{D}^{-1}$ ,

$$\vec{\varepsilon} = \mathbf{D}^{-1} \cdot \vec{\sigma}$$

with [check these!](#)

$$\mathbf{D}^{-1} = \frac{1}{\mu(3\lambda + 2\mu)} \begin{pmatrix} \lambda + \mu & -\lambda/2 & -\lambda/2 & 0 & 0 & 0 \\ -\lambda/2 & \lambda + \mu & -\lambda/2 & 0 & 0 & 0 \\ -\lambda/2 & -\lambda/2 & \lambda + \mu & 0 & 0 & 0 \\ 0 & 0 & 0 & 3\lambda + 2\mu & 0 & 0 \\ 0 & 0 & 0 & 0 & 3\lambda + 2\mu & 0 \\ 0 & 0 & 0 & 0 & 0 & 3\lambda + 2\mu \end{pmatrix}$$

then

$$\mathbf{D}^{-1} = \frac{1}{E} \begin{pmatrix} 1 & -\nu & -\nu & 0 & 0 & 0 \\ -\nu & 1 & -\nu & 0 & 0 & 0 \\ -\nu & -\nu & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2(1+\nu) & 0 & 0 \\ 0 & 0 & 0 & 0 & 2(1+\nu) & 0 \\ 0 & 0 & 0 & 0 & 0 & 2(1+\nu) \end{pmatrix}$$

Note that the determinant of  $\mathbf{D}^{-1}$  is  $8(1+\nu)^5(1-2\nu)E^{-6}$ , so that when  $\nu \rightarrow 1/2$  (incompressible material), the compliance matrix is singular and the stress cannot be given as a function of strain [816].

The above equation also leads to:

$$\begin{aligned} E\varepsilon_{xx} &= \sigma_{xx} - \nu(\sigma_{yy} + \sigma_{zz}) \\ E\varepsilon_{yy} &= \sigma_{yy} - \nu(\sigma_{xx} + \sigma_{zz}) \\ E\varepsilon_{zz} &= \sigma_{zz} - \nu(\sigma_{xx} + \sigma_{yy}) \\ E\varepsilon_{xy} &= (1+\nu)\sigma_{xy} \\ E\varepsilon_{xz} &= (1+\nu)\sigma_{xz} \\ E\varepsilon_{yz} &= (1+\nu)\sigma_{yz} \end{aligned} \tag{16.8}$$

The incompressibility (or bulk modulus)  $K$  is defined as  $p = -K\vec{\nabla} \cdot \vec{\mathbf{v}}$  where  $p$  is the pressure with

$$\begin{aligned} p &= -\frac{1}{3}\text{tr}(\boldsymbol{\sigma}) \\ &= -\frac{1}{3}[\lambda(\vec{\nabla} \cdot \vec{\mathbf{v}}) \text{tr}(\mathbf{1}) + 2\mu \text{tr}[\boldsymbol{\varepsilon}(\vec{\mathbf{v}})]] \\ &= -\frac{1}{3}[\lambda(\vec{\nabla} \cdot \vec{\mathbf{v}})3 + 2\mu(\vec{\nabla} \cdot \vec{\mathbf{v}})] \\ &= -\left(\lambda + \frac{2}{3}\mu\right)(\vec{\nabla} \cdot \vec{\mathbf{v}}) \end{aligned} \tag{16.9}$$

so that

$$p = -K\vec{\nabla} \cdot \vec{\mathbf{v}} \quad \text{with} \quad K = \lambda + \frac{2}{3}\mu \quad \text{and} \quad \boldsymbol{\sigma} = -p\mathbf{1} + 2\mu\boldsymbol{\varepsilon}^d(\vec{\mathbf{v}})$$

**Remark.** Eq. (16.1) and (16.3) are analogous to the ones that one has to solve in the context of viscous flow using the penalty method. In this case  $\lambda$  is the penalty coefficient,  $\vec{\mathbf{v}}$  is the velocity, and  $\mu$  is the dynamic viscosity.

**Remark.** Note that sometimes authors define  $p = -\lambda\vec{\nabla} \cdot \vec{\mathbf{v}}$  instead so that then  $\boldsymbol{\sigma} = -p\mathbf{1} + 2\mu\boldsymbol{\varepsilon}(\vec{\mathbf{v}})$  (to be very clear, strain tensor is not deviatoric), see for instance Sanan, May, Bollhöfer, and Schenk [1106] (2020) or Hansbo, Larson, and Larson [528] (2001).

The Lamé parameter  $\lambda$  and the shear modulus  $\mu$  are also linked to the Poisson ratio  $\nu$ , and  $E$ , Young's modulus:

$$\lambda = \mu \frac{2\nu}{1-2\nu} = \frac{\nu E}{(1+\nu)(1-2\nu)} \quad \text{with} \quad E = 2\mu(1+\nu)$$

The shear modulus, expressed often in GPa, describes the material's response to shear stress. The Poisson ratio describes the response in the direction orthogonal to uniaxial stress. The Young's modulus<sup>3</sup>, expressed in GPa, describes the material's strain response to uniaxial stress in the direction of this stress.

In the future we will also need to express the deviatoric part of a tensor as a function of the tensor itself, all in vector format. Let us consider the stress tensor. Then we have  $\vec{\tau} = \text{dev}(\boldsymbol{\sigma}) = \boldsymbol{\sigma} - \frac{1}{3}\text{tr}[\boldsymbol{\sigma}]\mathbf{1}$  which becomes

$$\vec{\tau} = \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} - \frac{1}{3} \begin{pmatrix} \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\ \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\ \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{2}{3}\sigma_{xx} - \frac{1}{3}\sigma_{yy} - \frac{1}{3}\sigma_{zz} \\ -\frac{1}{3}\sigma_{xx} + \frac{2}{3}\sigma_{yy} - \frac{1}{3}\sigma_{zz} \\ -\frac{1}{3}\sigma_{xx} - \frac{1}{3}\sigma_{yy} + \frac{2}{3}\sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} & 0 & 0 & 0 \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} & 0 & 0 & 0 \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{\tilde{\Lambda}^d} \cdot \vec{\sigma}$$

or,

$$\vec{\tau} = \tilde{\Lambda}^d \cdot \vec{\sigma} \quad (16.10)$$

where  $\tilde{\Lambda}^d$  is a deviatoric projection matrix.

Using the definition of  $K$  above, we have  $\lambda + 2\mu = K + \frac{4}{3}\mu$  and  $\lambda = K - \frac{2}{3}\mu$  so that the  $\mathbf{D}$  matrix can also be written as a function of  $K, \mu$ :

$$\begin{aligned} \mathbf{D}_{3D} &= \begin{pmatrix} K + \frac{4}{3}\mu & K - \frac{2}{3}\mu & K - \frac{2}{3}\mu & 0 & 0 & 0 \\ K - \frac{2}{3}\mu & K + \frac{4}{3}\mu & K - \frac{2}{3}\mu & 0 & 0 & 0 \\ K - \frac{2}{3}\mu & K - \frac{2}{3}\mu & K + \frac{4}{3}\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \\ &= \begin{pmatrix} K & K & K & 0 & 0 & 0 \\ K & K & K & 0 & 0 & 0 \\ K & K & K & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} \frac{4}{3}\mu & -\frac{2}{3}\mu & -\frac{2}{3}\mu & 0 & 0 & 0 \\ -\frac{2}{3}\mu & \frac{4}{3}\mu & -\frac{2}{3}\mu & 0 & 0 & 0 \\ -\frac{2}{3}\mu & -\frac{2}{3}\mu & \frac{4}{3}\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \\ &= K\mathbf{\Xi} + \mu\mathbf{\Lambda}^d \end{aligned} \quad (16.11)$$

$$\mathbf{D}_{3D} = K\mathbf{\Xi} + \mu\mathbf{\Lambda}^d$$

The expression above is to be found at [https://en.wikipedia.org/wiki/Linear\\_elasticity](https://en.wikipedia.org/wiki/Linear_elasticity)

<sup>3</sup>[https://en.wikipedia.org/wiki/Young's\\_modulus](https://en.wikipedia.org/wiki/Young's_modulus)

## 16.2 Plane strain

Typically one of the spatial dimensions (e.g.  $z$ ) is very large compared to the other two. As a consequence displacements  $\mathbf{v}_z$  and displacement derivatives  $\partial_z$  in the  $z$ -direction are assumed to be negligible, i.e.  $\varepsilon_{zz} = \varepsilon_{xz} = \varepsilon_{yz} = 0$  and Eqs. (16.8) become:

$$\begin{aligned} E\varepsilon_{xx} &= \sigma_{xx} - \nu(\sigma_{yy} + \sigma_{zz}) \\ E\varepsilon_{yy} &= \sigma_{yy} - \nu(\sigma_{xx} + \sigma_{zz}) \\ E\varepsilon_{zz} &= \sigma_{zz} - \nu(\sigma_{xx} + \sigma_{yy}) \\ E\varepsilon_{xy} &= (1 + \nu)\sigma_{xy} \\ E\varepsilon_{xz} &= (1 + \nu)\sigma_{xz} \\ E\varepsilon_{yz} &= (1 + \nu)\sigma_{yz} \end{aligned}$$

leading to  $\sigma_{xz} = \sigma_{yz} = 0$ ,  $\sigma_{zz} = \nu(\sigma_{xx} + \sigma_{yy})$  and

$$\begin{aligned} E\varepsilon_{xx} &= \sigma_{xx} - \nu(\sigma_{yy} + \sigma_{zz}) \\ &= \sigma_{xx} - \nu(\sigma_{yy} + \nu(\sigma_{xx} + \sigma_{yy})) \\ &= (1 - \nu^2)\sigma_{xx} - \nu(1 + \nu)\sigma_{yy} \\ E\varepsilon_{yy} &= \sigma_{yy} - \nu(\sigma_{xx} + \sigma_{zz}) \\ &= \sigma_{yy} - \nu(\sigma_{xx} + \nu(\sigma_{xx} + \sigma_{yy})) \\ &= -\nu(1 + \nu)\sigma_{xx} + (1 - \nu^2)\sigma_{yy} \\ E\varepsilon_{xy} &= (1 + \nu)\sigma_{xy} \end{aligned}$$

or,

$$\begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{xy} \end{pmatrix} = \frac{1}{E} \begin{pmatrix} 1 - \nu^2 & -\nu(1 + \nu) & 0 \\ -\nu(1 + \nu) & 1 - \nu^2 & 0 \\ 0 & 0 & 1 + \nu \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} \quad (16.12)$$

or<sup>4</sup>

$$\begin{aligned} \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} &= \frac{E}{(1 + \nu)(1 - 2\nu)} \begin{pmatrix} 1 - \nu & \nu & 0 \\ \nu & 1 - \nu & 0 \\ 0 & 0 & 1 - 2\nu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{xy} \end{pmatrix} \\ &= \frac{E}{(1 + \nu)(1 - 2\nu)} \begin{pmatrix} 1 - \nu & \nu & 0 \\ \nu & 1 - \nu & 0 \\ 0 & 0 & \frac{1}{2}(1 - 2\nu) \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ 2\varepsilon_{xy} \end{pmatrix} \end{aligned} \quad (16.13)$$

We then have

$$\mathbf{D}_{\text{plane strain}} = \frac{E}{(1 + \nu)(1 - 2\nu)} \begin{pmatrix} 1 - \nu & \nu & 0 \\ \nu & 1 - \nu & 0 \\ 0 & 0 & \frac{1}{2}(1 - 2\nu) \end{pmatrix} \quad (16.14)$$

**Remark.** The compliance matrix for plane strain is not found by removing columns and rows from the general isotropic compliance matrix!

Let us also look at another notation used in Simpson's book [1172] (Eq. 12.3). We start from Eq. (16.3):

<sup>4</sup>[https://www.efunda.com/formulae/solid\\_mechanics/mat\\_mechanics/hooke\\_plane\\_strain.cfm](https://www.efunda.com/formulae/solid_mechanics/mat_mechanics/hooke_plane_strain.cfm)



$$\begin{aligned}
\sigma_{xx} &= \lambda(\varepsilon_{xx} + \varepsilon_{yy}) + 2\mu\varepsilon_{xx} \\
&= (\lambda + 2\mu)\varepsilon_{xx} + \lambda\varepsilon_{yy} \\
\sigma_{yy} &= \lambda(\varepsilon_{xx} + \varepsilon_{yy}) + 2\mu\varepsilon_{yy} \\
&= \lambda\varepsilon_{xx} + (\lambda + 2\mu)\varepsilon_{yy} \\
\sigma_{xy} &= 2\mu\varepsilon_{xy}
\end{aligned} \tag{16.15}$$

so that

$$\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ 2\varepsilon_{xy} \end{pmatrix}$$

Since  $K = \lambda + 2\mu/3$ , we also have

$$\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} = \begin{pmatrix} K + \frac{4}{3}\mu & K - \frac{2}{3}\mu & 0 \\ K - \frac{2}{3}\mu & K + \frac{4}{3}\mu & 0 \\ 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ 2\varepsilon_{xy} \end{pmatrix}$$

or,

$$\mathbf{D}_{\text{plane strain}} = \begin{pmatrix} K + \frac{4}{3}\mu & K - \frac{2}{3}\mu & 0 \\ K - \frac{2}{3}\mu & K + \frac{4}{3}\mu & 0 \\ 0 & 0 & \mu \end{pmatrix}$$

The stress tensor is then as follows:

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} & 0 \\ \sigma_{yx} & \sigma_{yy} & 0 \\ 0 & 0 & \sigma_{zz} \end{pmatrix}$$

and its second moment invariant is given by

$$\mathcal{I}_2(\boldsymbol{\sigma}) = \frac{1}{2}\boldsymbol{\sigma} : \boldsymbol{\sigma} = \frac{1}{2}(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \sigma_{xy}^2$$

The deviatoric stress is given by  $\boldsymbol{\tau} = \boldsymbol{\sigma} - \frac{1}{3}\text{tr}(\boldsymbol{\sigma})\mathbf{1}$  with in this case

$$\begin{aligned}
\text{tr}(\boldsymbol{\sigma}) &= \sigma_{xx} + \sigma_{yy} + \sigma_{zz} \\
&= \sigma_{xx} + \sigma_{yy} + \nu(\sigma_{xx} + \sigma_{yy}) \\
&= (1 + \nu)(\sigma_{xx} + \sigma_{yy})
\end{aligned} \tag{16.16}$$

so that

$$\begin{aligned}
\boldsymbol{\tau} &= \boldsymbol{\sigma} - \frac{1 + \nu}{3}(\sigma_{xx} + \sigma_{yy})\mathbf{1} \\
&= \frac{1}{3} \begin{pmatrix} 3\sigma_{xx} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 3\sigma_{xy} & 0 \\ 3\sigma_{yx} & 3\sigma_{yy} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 0 \\ 0 & 0 & 3\sigma_{zz} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) \end{pmatrix} \\
&= \frac{1}{3} \begin{pmatrix} 3\sigma_{xx} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 3\sigma_{xy} & 0 \\ 3\sigma_{yx} & 3\sigma_{yy} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 0 \\ 0 & 0 & 3\nu(\sigma_{xx} + \sigma_{yy}) - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) \end{pmatrix} \\
&= \frac{1}{3} \begin{pmatrix} 3\sigma_{xx} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 3\sigma_{xy} & 0 \\ 3\sigma_{yx} & 3\sigma_{yy} - (1 + \nu)(\sigma_{xx} + \sigma_{yy}) & 0 \\ 0 & 0 & (2\nu - 1)(\sigma_{xx} + \sigma_{yy}) \end{pmatrix}
\end{aligned} \tag{16.17}$$

and also

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{2}\boldsymbol{\tau} : \boldsymbol{\tau} = \frac{1}{2}(\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2) + \tau_{xy}^2$$

**Remark.** In the case of a (near-)incompressible material,  $\nu \rightarrow \frac{1}{2}$  then  $\tau_{zz} \rightarrow 0$ .

## 16.3 Plane stress

For thin geometries. Let  $z$  be the direction perpendicular to the plate. Traction on the  $z$ -surface are assumed to be negligible, e.g.  $\sigma_{zz} = \sigma_{yz} = \sigma_{xz} = 0$

$$\begin{aligned} E\varepsilon_{xx} &= \sigma_{xx} - \nu(\sigma_{yy} + \sigma_{zz}) \\ E\varepsilon_{yy} &= \sigma_{yy} - \nu(\sigma_{xx} + \sigma_{zz}) \\ E\varepsilon_{zz} &= \sigma_{zz} - \nu(\sigma_{xx} + \sigma_{yy}) \\ E\varepsilon_{xy} &= (1 + \nu)\sigma_{xy} \\ E\varepsilon_{xz} &= (1 + \nu)\sigma_{xz} \\ E\varepsilon_{yz} &= (1 + \nu)\sigma_{yz} \end{aligned}$$

Immediately we have  $\varepsilon_{xz} = \varepsilon_{yz} = 0$ . Furthermore,

$$E\varepsilon_{xx} + E\varepsilon_{yy} = \sigma_{xx} - \nu\sigma_{yy} + \sigma_{yy} - \nu\sigma_{xx} = (1 - \nu)(\sigma_{xx} + \sigma_{yy})$$

so that the third equation can be written

$$\varepsilon_{zz} = -\frac{\nu}{1 - \nu}(\varepsilon_{xx} + \varepsilon_{yy})$$

Then,

$$\begin{aligned} E\varepsilon_{xx} &= \sigma_{xx} - \nu\sigma_{yy} \\ E\varepsilon_{yy} &= \sigma_{yy} - \nu\sigma_{xx} \\ E\varepsilon_{xy} &= (1 + \nu)\sigma_{xy} \end{aligned} \tag{16.18}$$

or,

$$\begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{xy} \end{pmatrix} = \frac{1}{E} \begin{pmatrix} 1 & -\nu & 0 \\ -\nu & 1 & 0 \\ 0 & 0 & 1 + \nu \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} \tag{16.19}$$

or<sup>5</sup>,

$$\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \end{pmatrix} = \frac{E}{(1 - \nu^2)} \begin{pmatrix} 1 & \nu & 0 \\ \nu & 1 & 0 \\ 0 & 0 & 1 - \nu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{xx} \\ \varepsilon_{yy} \\ \varepsilon_{xy} \end{pmatrix} \tag{16.20}$$

We then have

$$\mathbf{D}_{\text{plane stress}} = \frac{E}{(1 - \nu^2)} \begin{pmatrix} 1 & \nu & 0 \\ \nu & 1 & 0 \\ 0 & 0 & \frac{1}{2}(1 - \nu) \end{pmatrix} \tag{16.21}$$

**Remark.** The stiffness matrix for plane stress is not found by removing columns and rows from the general isotropic stiffness matrix.

<sup>5</sup>[https://www.efunda.com/formulae/solid\\_mechanics/mat\\_mechanics/hooke\\_plane\\_stress.cfm](https://www.efunda.com/formulae/solid_mechanics/mat_mechanics/hooke_plane_stress.cfm)

## 16.4 The axisymmetric case

We start from

$$\boldsymbol{\sigma} = \lambda(\vec{\nabla} \cdot \vec{v}) \mathbf{1} + 2\mu\boldsymbol{\varepsilon}(\vec{v}) \quad (16.22)$$

In cylindrical coordinates the velocity gradient is given by

$$\vec{\nabla}\vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & \frac{1}{r} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r} & \frac{\partial v_r}{\partial z} \\ \frac{\partial v_\theta}{\partial r} & \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} + \frac{v_r}{r} & \frac{\partial v_\theta}{\partial z} \\ \frac{\partial v_z}{\partial r} & \frac{1}{r} \frac{\partial v_z}{\partial \theta} & \frac{\partial v_z}{\partial z} \end{pmatrix}$$

In the case of axisymmetry, and in this case symmetry about the  $z$  axis, there is invariance with respect to the rotation around the axis so stresses and other quantities are independent of the  $\theta$  coordinate, or simply put  $\partial_\theta \rightarrow 0$ . The velocity gradient simplifies to:

$$\vec{\nabla}\vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & -\frac{v_\theta}{r} & \frac{\partial v_r}{\partial z} \\ \frac{\partial v_\theta}{\partial r} & \frac{v_r}{r} & \frac{\partial v_\theta}{\partial z} \\ \frac{\partial v_z}{\partial r} & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix}$$

Also, it follows logically that  $v_\theta = 0$  so that ultimately:

$$\vec{\nabla}\vec{v} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & 0 & \frac{\partial v_r}{\partial z} \\ 0 & \frac{v_r}{r} & 0 \\ \frac{\partial v_z}{\partial r} & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix}$$

and the strain tensor is then given by

$$\boldsymbol{\varepsilon}(\vec{v}) = \frac{1}{2} (\vec{\nabla}\vec{v} + \vec{\nabla}\vec{v}^T) = \begin{pmatrix} \frac{\partial v_r}{\partial r} & 0 & \frac{1}{2}(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) \\ 0 & \frac{v_r}{r} & 0 \\ \frac{1}{2}(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) & 0 & \frac{\partial v_z}{\partial z} \end{pmatrix} \quad (16.23)$$

The term  $\vec{\nabla} \cdot \vec{v}$  is simply the trace of  $\boldsymbol{\varepsilon}(\vec{v})$  so

$$\vec{\nabla} \cdot \vec{v} = \frac{\partial v_r}{\partial r} + \frac{v_r}{r} + \frac{\partial v_z}{\partial z}$$

Finally the full stress tensor is then

$$\begin{aligned} \boldsymbol{\sigma} &= \begin{pmatrix} \lambda(\frac{\partial v_r}{\partial r} + \frac{v_r}{r} + \frac{\partial v_z}{\partial z}) + 2\mu\frac{\partial v_r}{\partial r} & 0 & \mu(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) \\ 0 & \lambda(\frac{\partial v_r}{\partial r} + \frac{v_r}{r} + \frac{\partial v_z}{\partial z}) + 2\mu\frac{v_r}{r} & 0 \\ \mu(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) & 0 & \lambda(\frac{\partial v_r}{\partial r} + \frac{v_r}{r} + \frac{\partial v_z}{\partial z}) + 2\mu\frac{\partial v_z}{\partial z} \end{pmatrix} \\ &= \begin{pmatrix} (\lambda + 2\mu)\frac{\partial v_r}{\partial r} + \lambda(\frac{v_r}{r} + \frac{\partial v_z}{\partial z}) & 0 & \mu(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) \\ 0 & (\lambda + 2\mu)\frac{v_r}{r} + \lambda(\frac{\partial v_r}{\partial r} + \frac{\partial v_z}{\partial z}) & 0 \\ \mu(\frac{\partial v_z}{\partial r} + \frac{\partial v_r}{\partial z}) & 0 & (\lambda + 2\mu)\frac{\partial v_z}{\partial z} + \lambda(\frac{\partial v_r}{\partial r} + \frac{v_r}{r}) \end{pmatrix} \end{aligned}$$

As we did in the 2D case, we rewrite the six independent stress terms in to a vector  $\vec{\sigma}$  and we use Eq. (16.22) to arrive at:

$$\vec{\sigma} = \begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{r\theta} \\ \sigma_{rz} \\ \sigma_{\theta z} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 & 0 & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 & 0 & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{r\theta} \\ 2\varepsilon_{rz} \\ 2\varepsilon_{\theta z} \end{pmatrix} = \vec{\varepsilon}(\vec{\mathbf{v}})$$

or  $\vec{\sigma} = \mathbf{D} \cdot \vec{\varepsilon}(\vec{\mathbf{v}})$ . The components of the  $\vec{\varepsilon}(\vec{\mathbf{v}})$  vector are

$$\vec{\varepsilon}(\vec{\mathbf{v}}) = \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{r\theta} \\ 2\varepsilon_{rz} \\ 2\varepsilon_{\theta z} \end{pmatrix} = \begin{pmatrix} \frac{\partial \mathbf{v}_r}{\partial r} \\ \frac{\mathbf{v}_r}{r} \\ \frac{\partial \mathbf{v}_z}{\partial z} \\ 0 \\ \frac{\partial \mathbf{v}_z}{\partial r} + \frac{\partial \mathbf{v}_r}{\partial z} \\ 0 \end{pmatrix}$$

We see that there are two zeroes and consequently we'll find that  $\sigma_{r\theta}$  and  $\sigma_{\theta z}$  are also identically zero, so we discard these and end up with only four stress components :

$$\vec{\sigma} = \begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{rz} \end{pmatrix}$$

Note that in the literature the above relationship is often written

$$\begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix} = \frac{E}{(1+\nu)(1-2\nu)} \begin{pmatrix} 1-\nu & \lambda & \nu & 0 \\ \nu & 1-\nu & \nu & 0 \\ \nu & \nu & 1-\nu & 0 \\ 0 & 0 & 0 & (1-2\nu)/2 \end{pmatrix} \cdot \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ 2\varepsilon_{rz} \end{pmatrix}$$

which is equivalent since  $E = 2\mu(1+\nu)$  and  $\lambda = \frac{\nu E}{(1+\nu)(1-2\nu)}$  (see for instance Section 5.2.4 in [1430]).

#### about the implementation:

Only displacements in the  $r$  and  $z$  directions remain (note that  $\varepsilon_{\theta\theta}$  is in fact equal to  $\mathbf{v}_r/r$ ). In what follows I rename  $u = \mathbf{v}_r$  and  $\mathbf{v}_z = w$  to simplify notations. Then, inside an element we have

$$\begin{aligned} u^h(r, z) &= \sum_{i=1}^m \mathcal{N}_i(r, z) u_i \\ w^h(r, z) &= \sum_{i=1}^m \mathcal{N}_i(r, z) w_i \end{aligned} \tag{16.24}$$

where  $\mathcal{N}_i$  are the basis functions attached to the  $m$  nodes of the element. We compute the elements

of the  $\boldsymbol{\varepsilon}$  tensor of Eq. (16.23) as follows:

$$\varepsilon_{rr} = \frac{\partial u^h}{\partial r} = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, z) u_i \quad (16.25)$$

$$\varepsilon_{\theta\theta} = \frac{u_r^h}{r} = \frac{1}{r} \sum_{i=1}^m \mathcal{N}_i(r, z) u_i \quad (16.26)$$

$$\varepsilon_{zz} = \frac{\partial w^h}{\partial z} = \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial z}(r, z) w_i \quad (16.27)$$

$$\varepsilon_{rz} = \frac{1}{2} \frac{\partial u^h}{\partial z} + \frac{1}{2} \frac{\partial w^h}{\partial r} = \frac{1}{2} \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial z}(r, z) u_i + \frac{1}{2} \sum_{i=1}^m \frac{\partial \mathcal{N}_i}{\partial r}(r, z) w_i \quad (16.28)$$

Let us take  $m = 3$ , i.e. linear triangles, for simplicity. Then the strain vector  $\bar{\boldsymbol{\varepsilon}}^h$  is given by

$$\bar{\boldsymbol{\varepsilon}}^h = \begin{pmatrix} \varepsilon_{rr} \\ \varepsilon_{\theta\theta} \\ \varepsilon_{zz} \\ \textcolor{teal}{2}\varepsilon_{rz} \end{pmatrix} = \begin{pmatrix} \frac{\partial u^h}{\partial r} \\ \frac{u_r^h}{r} \\ \frac{\partial w^h}{\partial z} \\ \frac{\partial u^h}{\partial z} + \frac{\partial w^h}{\partial r} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial \mathcal{N}_1}{\partial r} & 0 & \frac{\partial \mathcal{N}_2}{\partial r} & 0 & \frac{\partial \mathcal{N}_3}{\partial r} & 0 \\ \frac{\mathcal{N}_1}{r} & 0 & \frac{\mathcal{N}_2}{r} & 0 & \frac{\mathcal{N}_3}{r} & 0 \\ 0 & \frac{\partial \mathcal{N}_1}{\partial z} & 0 & \frac{\partial \mathcal{N}_2}{\partial z} & 0 & \frac{\partial \mathcal{N}_3}{\partial z} \\ \frac{\partial \mathcal{N}_1}{\partial z} & \frac{\partial \mathcal{N}_1}{\partial r} & \frac{\partial \mathcal{N}_2}{\partial z} & \frac{\partial \mathcal{N}_2}{\partial r} & \frac{\partial \mathcal{N}_3}{\partial z} & \frac{\partial \mathcal{N}_3}{\partial r} \end{pmatrix}}_{\mathbf{B}(4 \times 6)} \cdot \underbrace{\begin{pmatrix} u1 \\ w1 \\ u2 \\ w2 \\ u3 \\ w3 \end{pmatrix}}_{\vec{\mathcal{U}}(6 \times 1)}$$

or  $\bar{\boldsymbol{\varepsilon}}^h = \mathbf{B} \cdot \vec{\mathcal{U}}$  and finally

$$\underbrace{\begin{pmatrix} \sigma_{rr} \\ \sigma_{\theta\theta} \\ \sigma_{zz} \\ \sigma_{rz} \end{pmatrix}}_{\vec{\boldsymbol{\sigma}}} = \underbrace{\begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix}}_{\mathbf{D}} \cdot \underbrace{\begin{pmatrix} \frac{\partial \mathcal{N}_1}{\partial r} & 0 & \frac{\partial \mathcal{N}_2}{\partial r} & 0 & \frac{\partial \mathcal{N}_3}{\partial r} & 0 \\ \frac{\mathcal{N}_1}{r} & 0 & \frac{\mathcal{N}_2}{r} & 0 & \frac{\mathcal{N}_3}{r} & 0 \\ 0 & \frac{\partial \mathcal{N}_1}{\partial z} & 0 & \frac{\partial \mathcal{N}_2}{\partial z} & 0 & \frac{\partial \mathcal{N}_3}{\partial z} \\ \frac{\partial \mathcal{N}_1}{\partial z} & \frac{\partial \mathcal{N}_1}{\partial r} & \frac{\partial \mathcal{N}_2}{\partial z} & \frac{\partial \mathcal{N}_2}{\partial r} & \frac{\partial \mathcal{N}_3}{\partial z} & \frac{\partial \mathcal{N}_3}{\partial r} \end{pmatrix}}_{\mathbf{B}(4 \times 6)} \cdot \underbrace{\begin{pmatrix} u1 \\ w1 \\ u2 \\ w2 \\ u3 \\ w3 \end{pmatrix}}_{\vec{\mathcal{U}}(6 \times 1)}$$

or,

$$\boxed{\vec{\boldsymbol{\sigma}} = \mathbf{D} \cdot \mathbf{B} \cdot \vec{\mathcal{U}}}$$

Note that in 2D, the matrix  $\mathbf{D}$  is  $3 \times 3$  and  $\mathbf{B}$  is  $3 \times 6$ .

I do not know yet how to arrive at what follows

The  $6 \times 6$  stiffness matrix is then

$$\mathbb{K} = \iiint \mathbf{B}^T \cdot \mathbf{D} \cdot \mathbf{B} \, dV$$

with  $dV = r dr d\theta dz$  in cylindrical coordinates. The integral over the  $\theta$  coordinate yields a factor  $2\pi$  so

$$\mathbb{K} = 2\pi \iint \mathbf{B}^T \cdot \mathbf{D} \cdot \mathbf{B} \, \textcolor{teal}{r} dr dz$$

The integration can now be performed as simply as was the case in the plane stress problem.

write the derivation for the rhs

Note that in practice the matrix  $\mathbf{D}$  is computed as follows (see for example Stone 63):

$$\mathbf{D} = \begin{pmatrix} \lambda + 2\mu & \lambda & \lambda & 0 \\ \lambda & \lambda + 2\mu & \lambda & 0 \\ \lambda & \lambda & \lambda + 2\mu & 0 \\ 0 & 0 & 0 & \mu \end{pmatrix} = \lambda \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} + \mu \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

## 16.5 FEM: Incompressible formulation from Zienkiewicz & Taylor book

This is from Volume 1-The Basis, page 307. Note that the authors use a different sign convention for pressure so that what follows is adapted from the book.

The authors start by defining the vector  $\vec{m}^T = (1, 1, 1, 0, 0, 0)$ . Using the vector notation of stress, the mean stress or pressure is given by

$$p = -\frac{1}{3}\text{tr}(\boldsymbol{\sigma}) = -\frac{1}{3}(\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) = -\frac{1}{3}(1, 1, 1, 0, 0, 0) \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} = -\frac{1}{3}\vec{m}^T \cdot \vec{\sigma}$$

For isotropic behaviour the 'pressure' is related to the volumetric strain by the bulk modulus of the material,  $K$ . Thus,

$$p = -K\text{tr}(\boldsymbol{\varepsilon}(\vec{v})) = -K\vec{m}^T \cdot \vec{\varepsilon}(\vec{v})$$

For an incompressible material  $K \rightarrow \infty$  and the volumetric strain is simply zero. Since  $\text{tr}(\boldsymbol{\varepsilon}(\vec{v})) = \vec{m}^T \cdot \vec{\varepsilon}(\vec{v})$  the deviatoric strain is defined by

$$\vec{\varepsilon}^d(\vec{v}) = \vec{\varepsilon}(\vec{v}) - \frac{1}{3}\vec{m} \text{tr}(\boldsymbol{\varepsilon}(\vec{v})) = \vec{\varepsilon}(\vec{v}) - \frac{1}{3}\vec{m}\vec{m}^T \cdot \vec{\varepsilon}(\vec{v}) = \left(\mathbf{1} - \frac{1}{3}\vec{m}\vec{m}^T\right) \cdot \vec{\varepsilon}(\vec{v}) = \mathbf{I}_d \cdot \vec{\varepsilon}(\vec{v})$$

where  $\mathbf{I}_d$  is the already defined deviatoric projection matrix (see Eq. (16.10)).

change notation to big lambda matrix?

In isotropic elasticity the deviatoric strain is related to the deviatoric stress by the shear modulus

$\mu$  as  $\boldsymbol{\tau} = 2\mu\boldsymbol{\varepsilon}^d(\vec{\mathbf{v}})$ , or

$$\begin{aligned}
\vec{\boldsymbol{\tau}} &= \underbrace{\begin{pmatrix} 2\mu & 0 & 0 & 0 & 0 & 0 \\ 0 & 2\mu & 0 & 0 & 0 & 0 \\ 0 & 0 & 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix}}_{\mathbf{C}_\mu} \cdot \vec{\boldsymbol{\varepsilon}}^d(\vec{\mathbf{v}}) \\
&= \mathbf{C}_\mu \cdot \mathbf{I}_d \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) \\
&= \mathbf{C}_\mu \cdot \left( \mathbf{1} - \frac{1}{3} \vec{m} \vec{m}^T \right) \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) \\
&= \mathbf{C}_\mu \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) - \frac{1}{3} \begin{pmatrix} 2\mu & 0 & 0 & 0 & 0 & 0 \\ 0 & 2\mu & 0 & 0 & 0 & 0 \\ 0 & 0 & 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & \mu & 0 & 0 \\ 0 & 0 & 0 & 0 & \mu & 0 \\ 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix} \cdot \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) \\
&= \mathbf{C}_\mu \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) - \frac{1}{3} \begin{pmatrix} 2\mu & 2\mu & 2\mu & 0 & 0 & 0 \\ 2\mu & 2\mu & 2\mu & 0 & 0 & 0 \\ 2\mu & 2\mu & 2\mu & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) \\
&= \underbrace{\left( \mathbf{C}_\mu - \frac{2\mu}{3} \vec{m} \vec{m}^T \right)}_{\mathbf{D}_d} \cdot \vec{\boldsymbol{\varepsilon}}(\vec{\mathbf{v}}) \tag{16.29}
\end{aligned}$$

We can also write  $\mathbf{D}_d$  in a more explicit manner more amenable to implementation:

$$\mathbf{D}_d = \frac{\mu}{3} \begin{pmatrix} 4 & -2 & -2 & 0 & 0 & 0 \\ -2 & 4 & -2 & 0 & 0 & 0 \\ -2 & -2 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{pmatrix} \tag{16.30}$$

CHANGE matrix names/notations!

In the mixed form considered next we shall use as variables the displacement  $\vec{\mathbf{v}}$  and the pressure  $p$ . Following the usual approach we arrive at a linear system  $\mathcal{A} \cdot \vec{X} = \vec{b}$  with (may be be a bit careful



about signs again)

$$\vec{X} = \begin{pmatrix} \vec{\mathcal{U}} \\ \vec{\mathcal{P}} \end{pmatrix}$$

$$\vec{b} = \begin{pmatrix} \vec{f} \\ \vec{0} \end{pmatrix}$$

$$\mathcal{A} = \begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{M}_K \end{pmatrix}$$

$$\mathbb{K} = \int_{\Omega} \boldsymbol{B}^T \cdot \boldsymbol{D}_d \cdot \boldsymbol{B} \, dV$$

$$\mathbb{G} = \int_{\Omega} \boldsymbol{B}^T \boldsymbol{N}_p \, dV$$

$$\mathbb{M}_K = \int_{\Omega} \frac{1}{K} \boldsymbol{N}_p^T \cdot \boldsymbol{N}_p \, dV$$

**Remark.** A similar approach is taken in Sanan, May, Bollhöfer, and Schenk [1106] (2020). However the authors define a new auxiliary pressure  $p = -\lambda \vec{\nabla} \cdot \vec{\mathbf{v}}$  (as opposed to  $p = -K \vec{\nabla} \cdot \vec{\mathbf{v}}$  above).

# 16.6 Elastic parameter values for Earth materials

what are  $E$ ,  $\nu$ ,  $\mu$ , etc ... for Earth ? bounds ? etc ...

| material  | Young's modulus (in $10^6$ bars) | shear modulus (in $10^6$ bars) |
|-----------|----------------------------------|--------------------------------|
| ice       | 0.1                              | 0.03                           |
| shale     | 0.2–0.3                          | 0.15                           |
| limestone | 0.4–0.7                          | 0.22–0.26                      |
| granite   | 0.3–0.6                          | 0.2                            |
| basalt    | 0.7–0.9                          | 0.3                            |
| steel     | 2.1                              | 0.83                           |

Taken from <https://www.britannica.com/science/rock-geology/Stress-strain-relationships>

Note that 1bar = 0.1MPa, so  $10^6\text{bar}=10^5\text{MPa}$ .  
In Farrington, Moresi, and Capitanio [387] (2014), we find “Elastic properties within the lithosphere and mantle are relatively well constrained with a shear modulus between  $10^{10}$  and  $10^{11}$  Pa.”

# 16.7 Benchmarks and analytical solutions

- Cook’s membrane problem (see Lamichhane [743] (2009) and refs therein; see Lamichhane [742] (2014)).
- Rectangular beam problem (see Lamichhane [743] (2009) and refs therein; see Lamichhane [742] (2014)).
- Thick-walled sphere under internal pressure (see Lamichhane [743] (2009) and refs therein).

# Chapter 17

## Visco-elasticity: physics, formulations and FEM

Work in progress!!!

### 17.1 A remark

In the absence of gravity, density does not enter the equations. As such compressible elasticity is fine.

If gravitational forces are taken into account, compressible elasticity means that density should change in space/time. However, a lot of the experiments hereafter are written with a constant density, meaning that they are reproducible or valid only for incompressible elastic (and viscou) flows!

### 17.2 Analytical Benchmarks

#### 17.2.1 the 1D solution

We wish to find the general solution  $\tau(t)$  of from the first order ODE

$$\frac{1}{2\mu} \frac{d\tau}{dt} + \frac{1}{2\eta} \tau = \dot{\varepsilon}_0 \quad \text{or,} \quad \frac{d\tau}{dt} + \frac{\mu}{\eta} \tau = 2\mu \dot{\varepsilon}_0$$

There is a standard technique to solve such equations, and we start with the following equation instead:

$$\frac{1}{2\mu} \frac{d\tau}{dt} + \frac{1}{2\eta} \tau = 0$$

It can be rewritten

$$\frac{d\tau}{\tau} = -\frac{\mu}{\eta} dt$$

so

$$\int \frac{d\tau}{\tau} = -\frac{\mu}{\eta} \int dt$$

$$\ln \tau - \ln \tau_0 = -\frac{\mu}{\eta} t$$

$$\ln \tau = -\frac{\mu}{\eta} t + \ln \tau_0$$

$$\tau(t) = \exp \left( -\frac{\mu}{\eta} t + \ln \tau_0 \right) = \exp \left( -\frac{\mu}{\eta} t \right) \exp (\ln \tau_0) = \tau_0 \exp \left( -\frac{\mu}{\eta} t \right)$$

or, using the Maxwell time  $t_M = \eta/\mu$ ,

$$\tau(t) = \tau_0 \exp\left(-\frac{t}{t_M}\right)$$

Based on this solution we now consider the following equation

$$\begin{aligned} \frac{d}{dt} \left[ \tau \exp\left(\frac{t}{t_M}\right) \right] &= \frac{d\tau}{dt} \exp\left(\frac{t}{t_M}\right) + \tau \frac{1}{t_M} \exp\left(\frac{t}{t_M}\right) \\ &= \underbrace{\left( \frac{d\tau}{dt} + \frac{\mu}{\tau} \tau \right)}_{2\mu\dot{\epsilon}_0} \exp\left(\frac{t}{t_M}\right) \end{aligned}$$

Then we proceed to integrate both sides:

$$\begin{aligned} \int \frac{d}{dt} \left[ \tau \exp\left(\frac{t}{t_M}\right) \right] dt &= 2\mu\dot{\epsilon}_0 \int \exp\left(\frac{t}{t_M}\right) dt \\ \tau(t) \exp\left(\frac{t}{t_M}\right) - \tau(0) \exp\left(\frac{0}{t_M}\right) &= 2\mu\dot{\epsilon}_0 t_M [\exp\left(\frac{t}{t_M}\right) - \exp\left(\frac{0}{t_M}\right)] \\ \tau(t) \exp\left(\frac{t}{t_M}\right) - \tau(0) &= 2\mu\dot{\epsilon}_0 \frac{\eta}{\mu} [\exp\left(\frac{t}{t_M}\right) - 1] \\ \tau(t) \exp\left(\frac{t}{t_M}\right) &= 2\mu\dot{\epsilon}_0 \frac{\eta}{\mu} [\exp\left(\frac{t}{t_M}\right) - 1] + \tau(0) \\ \tau(t) &= 2\eta\dot{\epsilon}_0 [1 - \exp\left(-\frac{t}{t_M}\right)] + \tau(0) \exp\left(-\frac{t}{t_M}\right) \end{aligned}$$

which is the solution of Eq. (4) of Kaus and Becker [680] (2007).

Alternative:

The general solution can be arrived at by means of the Laplace transform (!) and is given by:

$$\tau(t) = \tau(t_0) \exp\left(-\frac{t-t_0}{t_M}\right) + \exp\left(-\frac{t}{t_M}\right) \int_{t_0}^t 2\mu\dot{\epsilon}_T \exp\left(\frac{t'}{t_M}\right) dt'$$

If  $t_0 = 0$  and  $\tau(t_0) = 0$  then

$$\tau(t) = \exp\left(-\frac{t}{t_M}\right) \int_0^t 2\mu\dot{\epsilon}_T \exp\left(\frac{t'}{t_M}\right) dt'$$

If the strain rate and shear modulus are constant in time, then

$$\begin{aligned} \tau(t) &= \exp\left(-\frac{t}{t_M}\right) 2\mu\dot{\epsilon}_T \int_0^t \exp\left(\frac{t'}{t_M}\right) dt' \\ &= \exp\left(-\frac{t}{t_M}\right) 2\mu\dot{\epsilon}_T t_M \left[ \exp\left(\frac{t}{t_M}\right) - 1 \right] \\ &= 2\eta\dot{\epsilon}_T \left[ 1 - \exp\left(-\frac{t}{t_M}\right) \right] \end{aligned}$$

since  $t_M = \eta/\mu$ .

## 17.2.2 Pure shear

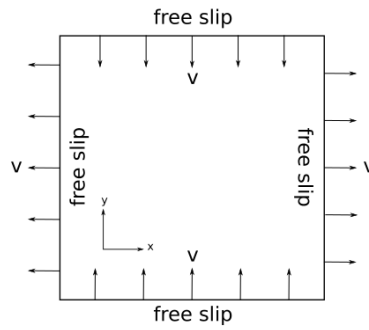
### fully incompressible

The first benchmark performed to test the viscoelastic implementation considers the stress build-up present in a viscoelastic Maxwell body. Contrary to stressed viscous materials, viscoelastic materials gradually build-up stress when sheared after which a transition to viscous deformation occurs.

An unstressed, incompressible viscoelastic Maxwell medium is subjected to a velocity field resulting in pure shear. The increase of the accumulated stress with time is given by an analytical solution:

$$\tau = 2\eta \dot{\epsilon} \left(1 - e^{-\frac{\mu}{\eta}t}\right) \quad (17.1)$$

with  $t$  time,  $\eta$  the prescribed material viscosity and  $\mu$  the prescribed material shear modulus. The domain size is  $100 \times 100 \text{ km}$ . The velocity prescribed at all boundaries equals  $v = 1 \text{ cm/yr}$  in magnitude yielding a constant background strain rate of  $\dot{\epsilon} = 2 \text{ cm/yr} / 100 \text{ km} \simeq 6.342 \times 10^{-15} \text{ s}^{-1}$ . The viscosity is  $\eta = 10^{21} \text{ Pa.s}$ , the shear modulus is  $\mu = 10^{10} \text{ Pa}$  and the gravity is set to zero. We set  $\delta t = 100 \text{ yr}$ .



Set up of the stress build-up benchmark. All domain sides have a free slip boundary condition, and pure shear velocity conditions are prescribed. Adapted from Gerya (2010) [455].

We have

$$\eta_{eff} = \frac{\eta \delta t}{\delta t + \eta/\mu} = \frac{10^{21} \cdot 3.154 \times 10^9}{3.154 \times 10^9 + 10^{21}/10^{10}} \simeq 3.0592 \times 10^{19} \text{ Pa.s} \quad \text{and} \quad Z = \frac{\eta_{eff}}{\mu \delta t} \simeq 0.9694$$

The Maxwell time is  $t_M = \frac{\eta}{\mu} = 10^{11} \text{ s} \simeq 3171 \text{ yr}$ . In the absence of elasticity (purely viscous behaviour), we have  $\dot{\epsilon}_{xx} = 6.342 \times 10^{-15}$  and  $\eta = 10^{21}$  so the deviatoric stress  $\tau_{xx}$  is equal to

$$\tau_{xx} = 2 \cdot 10^{21} \cdot 6.342 \times 10^{-15} \simeq 12.68 \times 10^6 \text{ Pa}$$

The first time that the Stokes system is solved, there is no stored stress, i.e. the elastic rhs is identically zero, so that the system is solved with a viscosity equal to  $\eta_{eff}$ . We can easily compute the analytical solution, and we see that  $\dot{\epsilon}_{xy} = 0$  and  $\dot{\omega}_{xy} = 0$ , which we recover:

Results in Stone 64!

The expected stress value for  $\tau_{xx}$  after the first Stokes solve is

$$\tau_{xx} = 2\eta_{eff}\dot{\epsilon}_{xx} = 2 \cdot 3.057 \times 10^{19} \cdot 6.342 \times 10^{-15} \simeq 38.775 \times 10^4 \text{ Pa}$$

Also check Gerya's book page 358 2nd edition ?

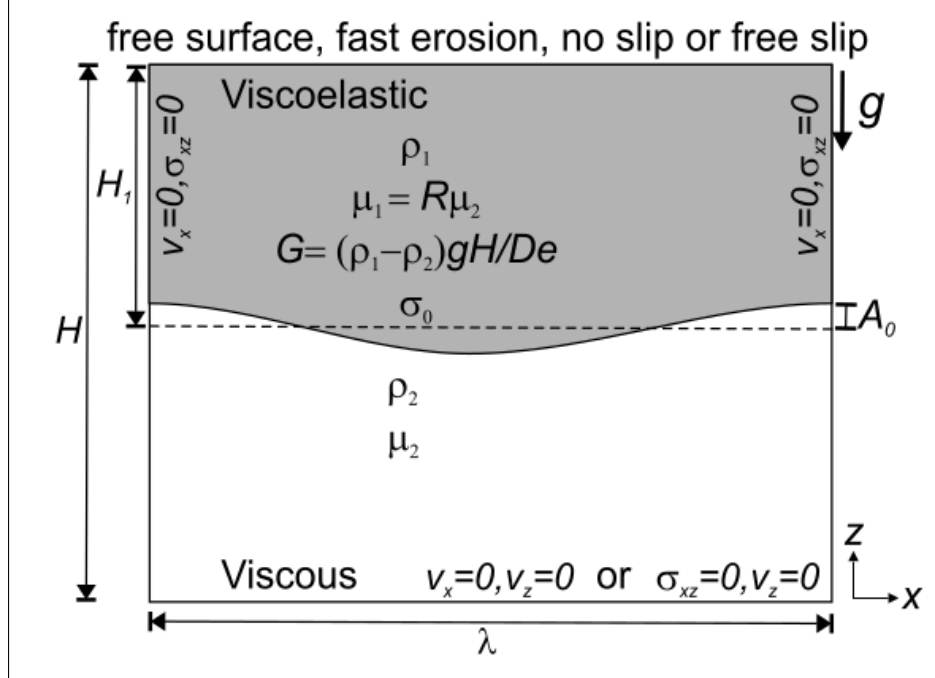
## 17.2.3 simple shear

## 17.2.4 Rayleigh-Taylor instability

### fully incompressible

This experiment is presented in Kaus and Becker [680] (2007).

The model consists of a viscoelastic layer of thickness  $H_1$ , with density  $\rho_1$ , viscosity  $\eta_1$  and elastic shear module  $\mu_1$  that overlies a viscous layer of thickness  $H - H_1$ , with density  $\rho_2$  and viscosity  $\mu_2$ . The interface between the two fluids is perturbed sinusoidally according to  $h(x) = (H - H_1) + A_0 \cos(2\pi x/\lambda)$ , where  $H_1$  is the thickness of the upper layer,  $H$  the height of the model,  $A_0$  the initial amplitude and  $\lambda$  the wavelength of the perturbation. If  $\rho_1 > \rho_2$ , the system is gravitationally unstable.



Taken from Kaus and Becker [680].

The Deborah number is given by

$$De = \frac{(\rho_1 - \rho_2)gH}{\mu}$$

and is here defined as the ratio between the viscous (Stokes) timescale  $((\rho_1 - \rho_2)gH/\eta_1)$  and the viscoelastic timescale  $\eta_1/\mu$  of the upper layer. The authors state: “ Interestingly, the Deborah number, which is a measure of the importance of elasticity [...], is independent on the viscosity of the system. This is due to the fact that the magnitude of stress is solely dependent on the density difference for purely buoyancy-driven flow.” Also: “ In the present definition of the Deborah number realistic values for lithospheric-scale deformation are  $10^{-4} \leq De \leq 1$  (with  $\rho = 10 - 330 \text{ kg m}^{-3}$ ,  $g = 10 \text{ m s}^{-2}$ ,  $H = 100 - 3000 \text{ km}$ ,  $\mu = 10^{10} - 10^{11} \text{ Pa}$ ).”

Another parameter controlling the dynamics of the system is the viscosity contrast between the upper and the lower layer, expressed by  $R = \eta_1/\eta_2$ .

We here set  $H = 500 \text{ km}$ ,  $H_1 = 100 \text{ km}$ ,  $\rho_2 = 3300 \text{ kg m}^{-3}$ ,  $\eta_2 = 10^{21} \text{ Pa s}$ ,  $\mu_1 = \mu_2 = 10^{10} \text{ Pa}$ ,  $\eta_1 = R\eta_2$ ,  $A_0 = 1 \text{ km}$ ,  $\lambda = L_x$ .

### 17.2.5 stress build-up inside an elastic inclusion in a viscous matrix (Beuchert & Podlachikov

#### fully incompressible

This is presented in Section 5.4 of Beuchert and Podlachikov [86] (2010). The authors test their viscoelastic FEM model against the analytical solution for stress built-up inside an elastic inclusion in a viscous matrix under pure shear.

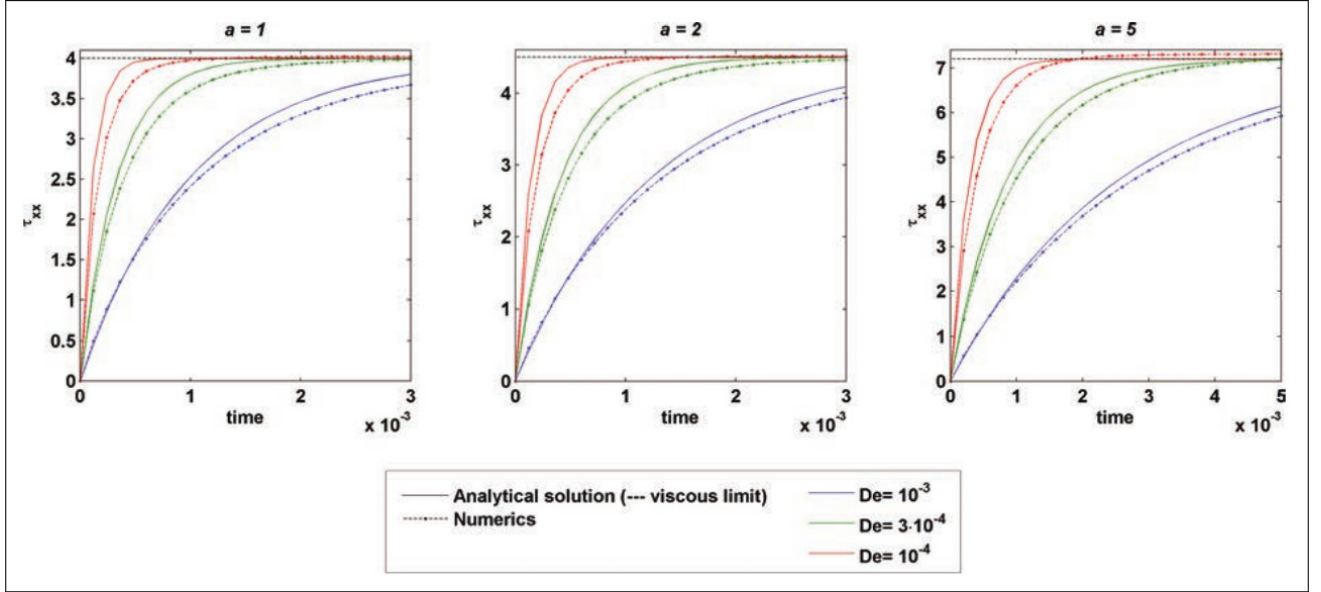
Note that the fluid in question is incompressible as shown by their Eq. (3).

An analytical solution applicable for the time-dependent stress evolution inside a viscoelastic inclusion embedded in a viscous matrix can be derived:

$$\tau_{xx} = -\frac{(a+1)^2 \left[-1 + \exp -\frac{2\mu at}{a^2+1}\right]}{a} \dot{\epsilon}_{xx}^d$$

where  $a$  is the aspect ratio of the inclusion and the matrix viscosity is assumed to be unity.

For this equation to be applicable to our viscoelastic model, we prescribe a high viscosity inside the inclusion ( $\eta_{inclusion} = 10^5$ ) and thus obtain an elastic response inside the inclusion. For the matrix, we apply a viscous rheology with viscosity  $\eta_{matrix} = 1$ .

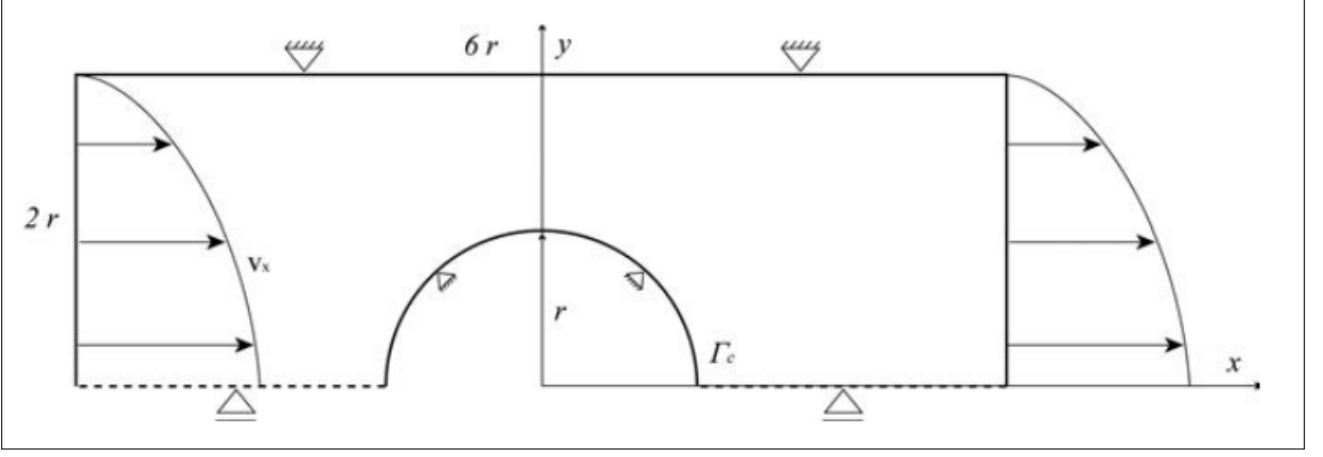


Taken from Beuchert and Podladchikov [86]. Stress built-up inside an elliptical elastic inclusion with aspect ratios  $a = 1, 2, 5$  embedded in a viscous matrix under constant pure shear loading. Comparison of the analytical solution with the numerical solution obtained in their FEM model for different values of Deborah number  $De = 10^{-4}, 3 \cdot 10^{-4}, 10^{-3}$  at non-dimensional deviatoric strain rate  $\dot{\epsilon}_{xx}^d = 1$ . Time and deviatoric stress  $\tau_{xx}$  are non-dimensional. The test results show a good agreement of the numerical solution with the analytical solution. For time  $t \rightarrow \infty$ , the solution for stress inside the elastic inclusion converges towards the solution for stress inside a viscous inclusion (dashed horizontal line 'viscous limit').

### 17.2.6 Viscoelastic flow past a cylinder in a channel (Beuchert & Podladchikov)

This is presented in Section 5.4 of Beuchert and Podladchikov [86] (2010).

We tested the flow code against a numerical benchmark for iso-viscous, viscoelastic flow past a circular cylinder in a channel. Figure below shows the domain setup and boundary conditions for this benchmark. The radius of the circular cylinder  $r = 1$  is half the domain height and the domain aspect ratio is 3 : 1.



Taken from Beuchert and Podladchikov [86]. Setup for the benchmark of viscous flow past circular cylinder in a channel. The domain is symmetric about the horizontal axis and thus only the upper half of the channel is modelled.

At the inflow and outflow boundaries, an established Poiseuille flow  $v_x = 3(R^2 - y^2)/(2R^2)$  is imposed;  $v_y = 0$  at the sides. Both upper and lower boundaries are fixed in  $y$ -direction. We apply no-slip boundary conditions at the top and along the cylinder wall and free-slip (zero traction) conditions at the bottom (symmetry axis). For the inflow conditions on  $\sigma$ , we use the analytical solution for simple shear of a Jaumann fluid, which is valid for Poiseuille flow. In that case, the Jaumann derivative equations are

$$\tau_{xx} + 2De\omega\tau_{xy} = 0 \quad (17.2)$$

$$\tau_{yy} + 2De\omega\tau_{xy} = 0 \quad (17.3)$$

$$\tau_{xy} + 2De\omega(\tau_{yy} - \tau_{xx}) = 2\mu\dot{\epsilon}_{xy} \quad (17.4)$$

$$(17.5)$$

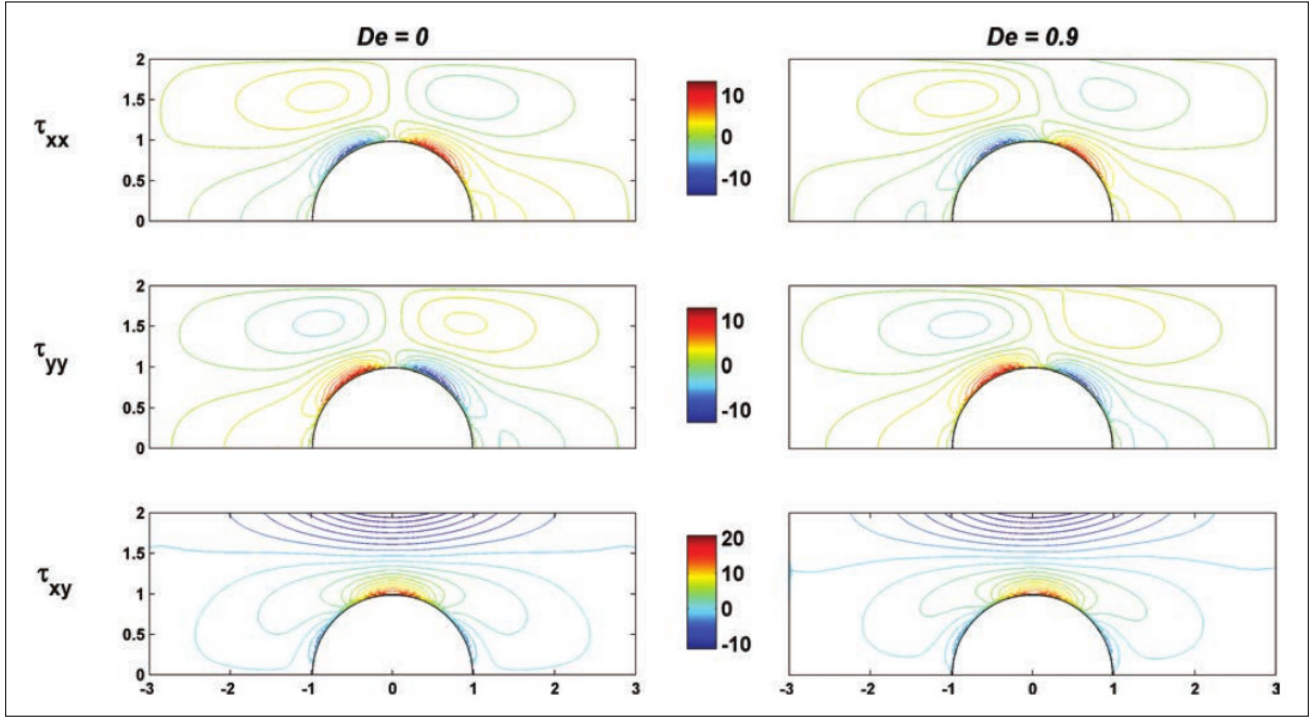
Given that  $\dot{\gamma} = 2\dot{\epsilon}_{xy}$  and  $\omega = -1/2\dot{\gamma}$  for simple shear, we obtain from the equation above

$$\tau_{xy} = \frac{2\mu\dot{\epsilon}_{xy}}{1 + 4De^2\omega^2} = \frac{\mu\dot{\gamma}}{1 + 4De^2\omega^2} = \frac{\mu\dot{\gamma}}{1 + 2De^2\dot{\gamma}^2}$$

and can then solve for  $\tau_{xx}$  and  $\tau_{yy}$  by substitution. The strain rate  $\dot{\gamma}$  is given by  $\partial v_x / \partial y$ , that is, by differentiating the inflow condition  $v_x = 3(R^2 - y^2)/(2R^2)$  with respect to  $y$ , resulting in  $\dot{\gamma} = -3y/R^2$  at the inflow boundary.

The authors further measure the non-dimensional drag force  $C_d$  exerted by the passing fluid on the cylinder wall and compare it with published values.

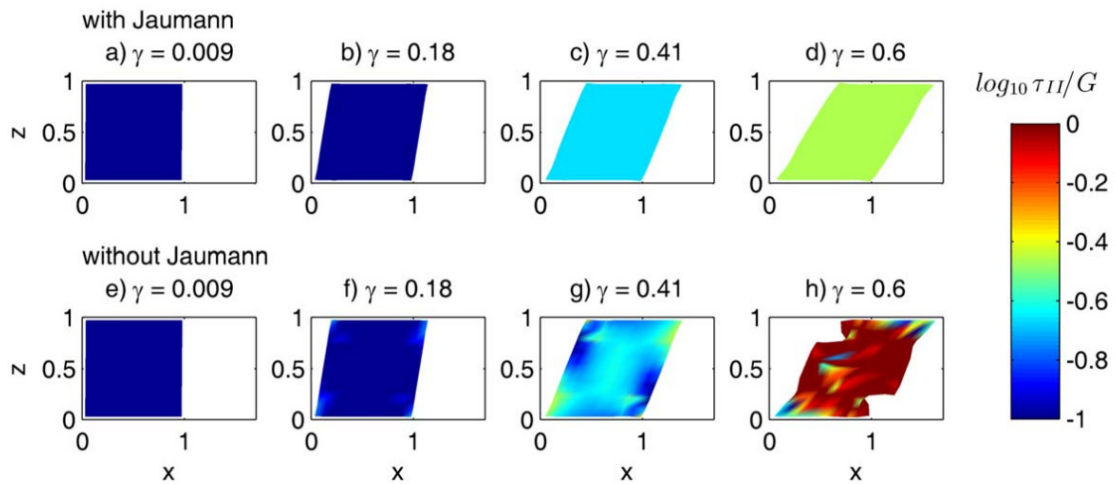




Taken from Beuchert and Podladchikov [86]. Non-dimensional deviatoric stresses  $\tau$  at  $De = 0$  and  $De = 0.9$  resulting from fluid flow past a cylinder in a channel. The benchmark is conducted on a regular, structured grid at resolution  $600 \times 200$  nodes.

### 17.2.7 Elastic Simple Shear, von Tschärner and Schmalholz [1328] (2015)

This experiment uses a homogeneous cube with the dimensions  $1 \times 1 \times 1$  which is deformed by simple shear. The bottom boundary is fixed (i.e., the boundary condition is no slip), the velocities on the top surface are prescribed in x-direction to generate simple shear and the top boundary is fixed in z-direction. The boundary conditions are free slip for two vertical boundaries ( $y=0$  and  $y=1$ ) and the two remaining vertical boundaries ( $x=0$  and  $x=1$ ) are free. The viscosity and the elastic shear modulus are  $\eta = 10^{10}$  and  $G = 1$ , respectively. All model dimensions and material parameters are given in dimensionless numbers using the model length, the elastic shear modulus and the background strain rate as characteristic parameters. The results for two simulations, one with Jaumann correction and one without Jaumann correction, are given in Figure B5 where the colors indicate the second invariant of the stress tensor. The results with the Jaumann correction show a homogeneous distribution of stress (Figures a–d), whereas the distribution of the second invariant of the stress tensor is inhomogeneous without the Jaumann correction and the simulation “crashes” for high strain (Figures e–h).



Deformation of a homogeneous cube under simple shear (a–d) with Jaumann correction and (e–h) without Jaumann correction for different amounts of bulk shear strain  $c$  (i.e., ratio of maximal horizontal displacement to model thickness). The colors indicate the second invariant of the stress tensor  $\tau_{ii}$ . Without the Jaumann corrections, the stress distribution becomes inhomogeneous and the simulation crashes for high strain.

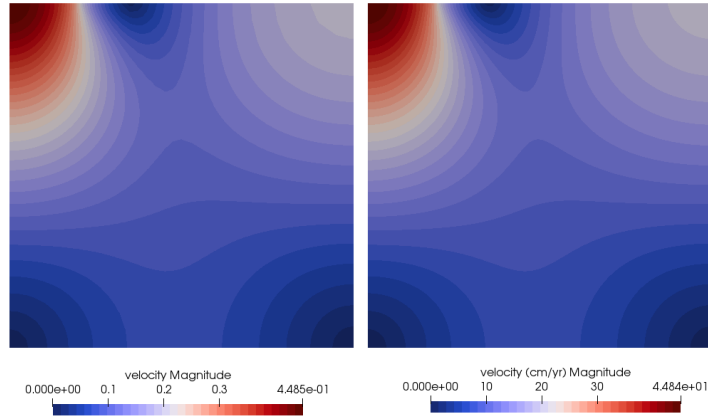
## 17.2.8 Response to load from ice sheet - Nakiboglu and Lambeck (1982)

The domain is  $500 \times 500$  km. Vertical gravity is 9.8, density is 3300, viscosity is  $3 \cdot 10^{20}$  Pa s, shear modulus is  $10^{10}$  Pa, free slip on left, right and bottom. A normal stress is imposed on the top for  $0 < x < 100$  km. It corresponds to an ice sheet of density  $\rho_i = 900$  of 1000 m height. The timestep is set to 100 yr. Resolution is set to  $50 \times 50$  elements. Stress/traction b.c. are explained in Section 7.13.

Analytical solution is provided in Nakiboglu and Lambeck (1982) [925]. Note however that this is a 2D setup while the original solution is for a cylindrical load and also for a semi-infinite domain.

Effective viscosity is given by

$$\eta_{eff} = \frac{\eta \delta t}{\delta t + \eta/\mu} = 2.85362675546547e + 19$$



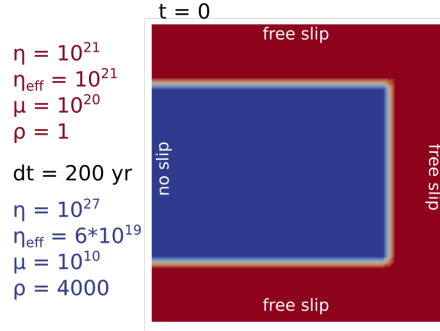
Left: ASPECT; Right: stone 64

It was run with ASPECT: prm file in stone 64 results folder.

## 17.3 Numerical Benchmarks

### 17.3.1 Bending of elastic slab (Gerya's book)

The sinking slab benchmark consists of a beam of elastic material which is placed in a weak and viscous surrounding medium. The initially unstressed beam is attached to the left domain boundary through boundary conditions. A stress is then applied to the beam in the form of gravity. The applied gravity force results in the deformation of the beam through bending. After 20 kyr, the gravity field is turned off and the elastic properties of the beam will then force itself to its original position. The set-up of the benchmark is given in the following figure:



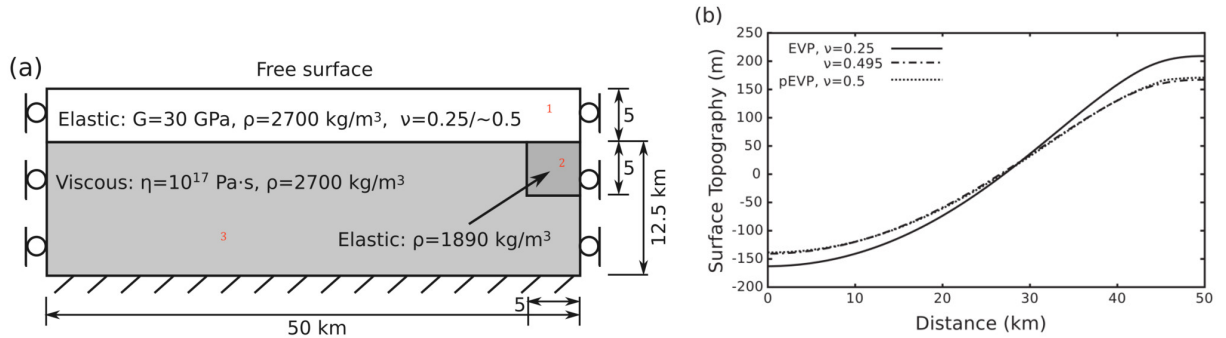
Set-up of the benchmark from [455]. The properties of the two materials are given on the left, together with the initial configuration of the benchmark.

The beam is surrounded by a low-density, low-viscosity and high shear modulus medium of which the specifications are given in the following table. The boundary conditions of the domain consist of a no slip condition at the left boundary where the slab is attached and free slip boundary conditions along all other sides. The results are calculated on a grid with a resolution of 50x50 elements containing 64 randomly distributed markers at startup. The time step is set to  $\delta t = 200yr$  (i.e. gravity is switched off after 100 time steps).

| Material properties                 | Elastic slab (fluid 1) | Surrounding medium (fluid 2) |
|-------------------------------------|------------------------|------------------------------|
| Density $\rho$ [kg/m <sup>3</sup> ] | 4000                   | 1                            |
| Viscosity $\eta$ [Pa·s]             | $10^{27}$              | $10^{21}$                    |
| Shear modulus $\mu$ [Pa]            | $10^{10}$              | $10^{20}$                    |
| Maxwell time $t_M$ [yr]             | 3.17Gyr                | $3.17 \times 10^{-7}yr$      |
| eff. visc. $\eta_{eff}$ [Pa·s]      | 6.307199602192306e+19  | 9.999999984145105e+20        |
| visco-elasticity factor $Z$ [-]     | 0.9999999369280039     | 1.5854895966744522e-09       |

### 17.3.2 Flexure of elastic plate (Choi et al)

This benchmark is presented in Choi, Tan, Lavier, and Calo [238] (2013). The setup is as follows:



Taken from Choi, Tan, Lavier, and Calo [238]. “Effect of elastic compressibility on the prediction of an elastic thin plate subject to an uplifting load. (a) Model setup for a finite length elastic layer subjected to a finite length buoyant load applied on the bottom. (b) Profiles of mean-subtracted surface topography.”

| Material properties               | elastic plate (1) | elastic block (2) | viscous mantle (3)     |
|-----------------------------------|-------------------|-------------------|------------------------|
| Density $\rho$ kg m <sup>-3</sup> | 2700              | 1890              | 2700                   |
| Viscosity $\eta$ Pa s             | $10^{35}$         | $10^{35}$         | $10^{17}$              |
| Shear modulus $\mu$ Pa            | $30 \cdot 10^9$   | $30 \cdot 10^9$   | $10^{50}$              |
| Maxwell time $t_M$ yr             | 10569930          | 10569930          | 3.1709791983764584e-41 |
| eff. visc. $\eta_{eff}$ Pa s      | 4.73039776e+18    | 4.73039776e+18    | 1e+17                  |

The value of  $\eta_1 = \eta_2 = 10^{35}$  for the elastic materials was obtained through personal communication. The value of  $\mu_3 = 10^{50}$  for the viscous material ensures that  $\eta_{eff} = \eta_3$ . Note that in

the publication the authors test both compressible and incompressible formulations, but we restrict ourselves to incompressible results since our code cannot handle compressible behavior. I also use  $dt=5\text{year}$ . Gravity is not specified in the paper.

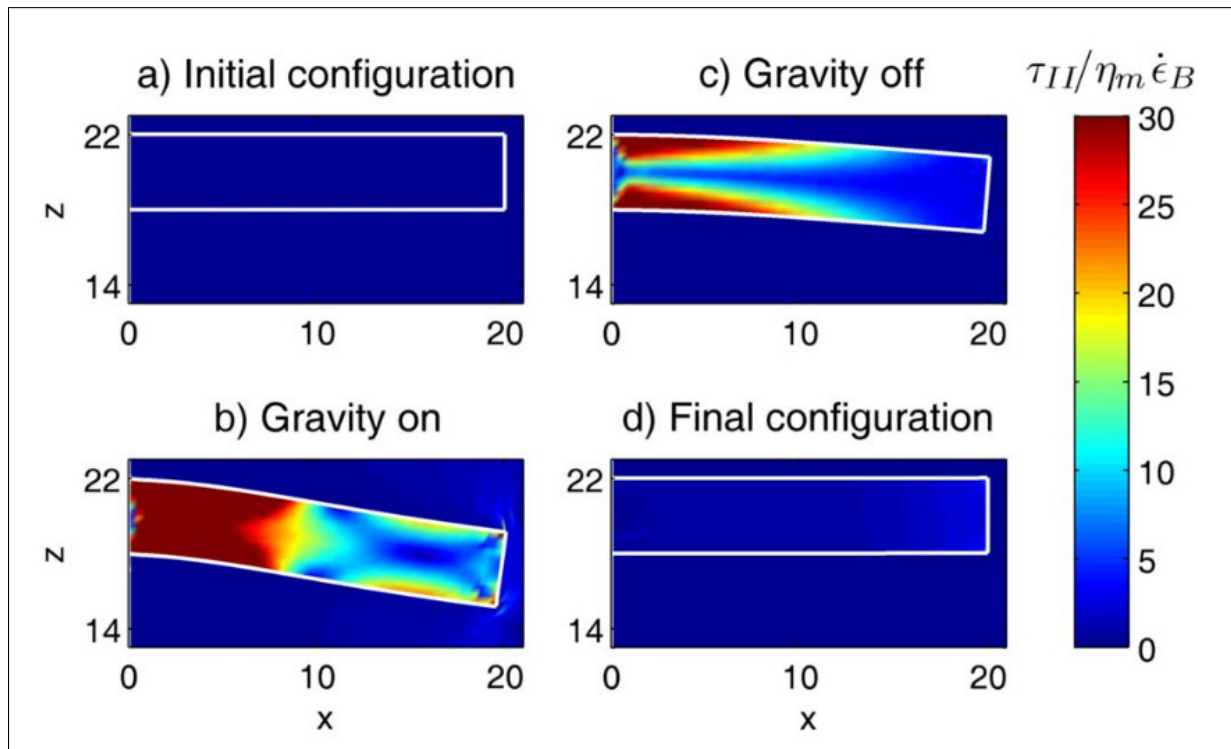
The authors report a converged total relief of 306-308m.

This benchmark requires either a sticky air layer (see Section 9.1) on top of the plate or a deformable mesh (ALE formulation, see Section 9.1).

### 17.3.3 Elastic beam in viscous matrix - von Tscharner and Schmalholz

What follows is taken from von Tscharner and Schmalholz [1328] (2015).

This experiment first is a cylindrical elastic beam in a viscous matrix under gravity. The model box has the dimensions of  $10 \times 0.25 \times 10$ . The elastic beam which is vertically located in the middle of the model box has a thickness of  $H = 1$  and a length of  $L = 5$ . The viscosity is  $\eta_m = 1$  and  $\eta_b = 10^{13}$  for the matrix and the beam, respectively. The elastic shear modulus is  $G_m = 5 \cdot 10^{10}$  and  $G_b = 5 \cdot 10^3$  for the matrix and beam, respectively, and the density difference is  $\rho_b - \rho_m = 300$ . These parameters provide a beam that is effectively elastic and a matrix that is effectively viscous. All model dimensions and material parameters are given in dimensionless numbers using the thickness of the beam  $H$ , the matrix viscosity and the background strain rate as characteristic parameters. The boundary conditions are free slip for all boundaries. The results of this simulation are shown in the figure below where the colors indicate the second invariant of the stress tensor. The elastic beam is deflected downward under vertical gravity (Figure b). When the gravity is turned off, the beam deflects upward due to the stored elastic energy and recovers the original rectangular shape which is stress free (Figure c). A similar test is given in Gerya's book [2010]. The test shows the reversible elastic deformation. The elastic beam recovers its original rectangular shape and stress state when the applied load is removed.



Reversible deformation of an elastic beam in a viscous matrix under gravity. (a) Unstressed initial configuration (gravity off). (b) Deformation of the elastic beam under gravity. (c) Gravity is turned off. (d) The elastic beam recovers the original rectangular shape with zero stress. Colors indicate the second invariant of the stress tensor  $\tau_{II}$ .

### 17.3.4 Elastic beam in viscous matrix - Keller, May, and Kaus

This benchmark comes from Appendix B of Keller et al (2013) [690].

The domain is  $7.5 \times 5$  km. A dominantly elastic beam is fixed to, and protrudes horizontally from the left wall of the model box. Surrounding the elastic beam is a viscous, but inelastic fluid. All boundaries are free slip, except for the left wall, which is set to no slip in order to keep the bending beam fixed to the wall. The beam has a higher density than the surrounding fluid and thus will bend down elastically driven by gravity. After the beam has accumulated some elastic strain through bending down, gravity is switched off. If the stress evolution is implemented accurately, the elastic beam should now, free from the pull of gravity, move upwards again and restore its initial position.

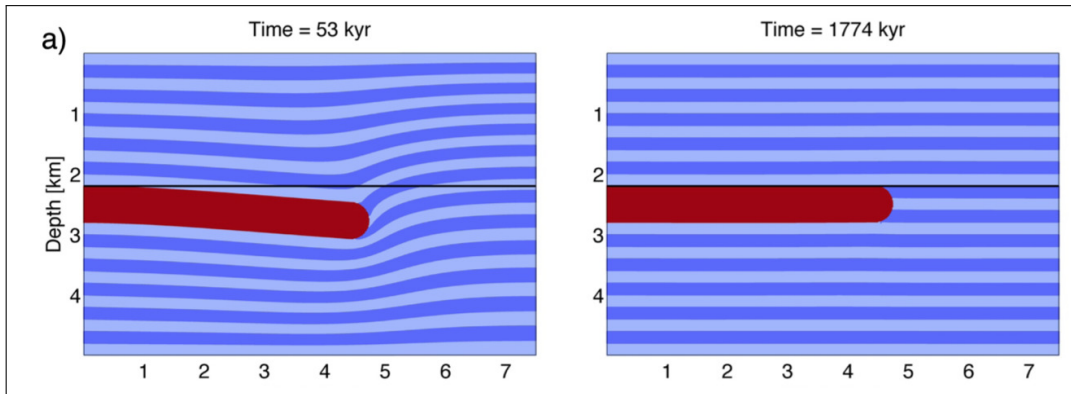
Material properties are as follows:

- beam:  $\rho = 1500$ ,  $\eta = 10^{24}$ ,  $\mu = 10^{10}$
- fluid:  $\rho = 1000$ ,  $\eta = 10^{18}$ ,  $\mu = 10^{11}$

This choice of parameters leads to a Maxwell time  $t_m = 0.32$  yr for the background fluid and Maxwell times of  $t_m = 3.2$  Myr for the beam, meaning that the deformation in this benchmark problem, which occurs on a timescale of thousands to a million years, will lead to dominantly viscous deformation in the fluid, and dominantly elastic behaviour of the beam.

Keller et al set the numerical resolution to  $300 \times 200$  elements, with 16 markers per elements for stress advection. Such a resolution is not feasible with our simple python implementation so the resolution is then set to  $96 \times 64$ .

The following plot comes from [690]:

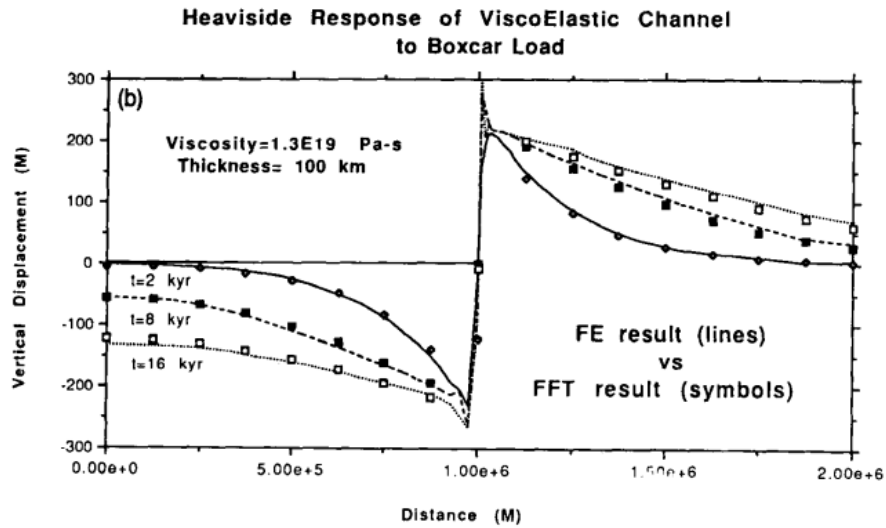


The elastic timestep is set to  $\delta t_e = 100$  yr and the tectonic timestep is set to the same value. This yields  $\eta_{eff} = 10^{18}$  in the fluid, and  $\eta_{eff} \simeq 3.15 \times 10^{19}$ . After 50 kyr, the gravity ( $|\vec{g}| = 10$ ) is switched off and the model is ran for another 500 kyr.

### 17.3.5 Boxcar load on an incompressible viscoelastic lithosphere - Wu (1992)

fully incompressible

This originates in Wu [1372] (1992). The domain is 2500 km long and it is either a halfspace in the vertical direction or a channel of 100 km width. The material is characterised by  $\rho = 3400 \text{ kg m}^{-3}$ , Young's modulus  $E = 1.13 \times 10^{11} \text{ Pa}$  and a Poisson ratio  $\nu = 0.5$  with  $|\vec{g}| = 9.82 \text{ m s}^{-2}$ . The author explores the effect of a linear viscous mantle vs dislocation creep. Time evolution/relaxation figures are available.



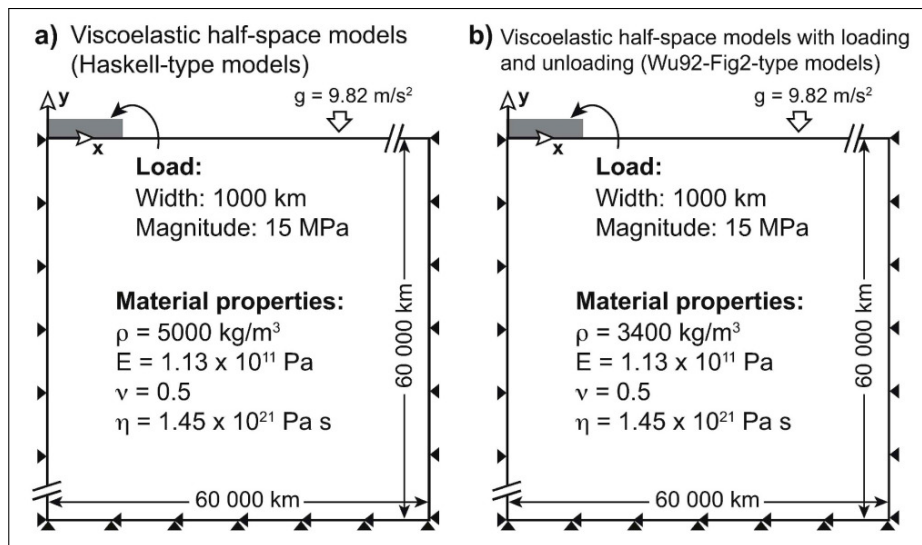
The Heaviside responses of linear viscoelastic earth models calculated with the finite element method (lines) are compared to those computed using the conventional transform method (symbols). For a 100 km thick channel.

also explore Wu & Peltier [1373] (1982).

### 17.3.6 Boxcar load on an incompressible viscoelastic lithosphere - Hampel et al. (2019)

Hampel, Lüke, Krause, and Hetzel [523] (2019)

fully incompressible



a) Viscoelastic Haskell-type half-space models. b) Viscoelastic half-space models used for modelling loading and subsequent unloading (after Wu, 1992; his Fig. 2). Both model types are meshed with 25x25 km large linear, rectangular plane strain elements suitable for incompressible materials. The same mesh and element type are used in all model runs. All viscoelastic half-space models have the same boundary conditions (indicated by black triangles): the model bottom is fixed in both the vertical and horizontal direction while the model sides are fixed in the horizontal direction. In models with an elastic foundation (Table 1), it is applied to the model surface. Abbreviations for model parameters are  $\rho$  density,  $E$  Young's modulus,  $\nu$  Poisson's ratio,  $\eta$  viscosity and  $g$  acceleration due to gravity. Following Wu (1992), we use a value of  $g = 9.82 \text{ m/s}^2$  in the viscoelastic half-space models. In all models, the load is applied (and removed, if applicable) instantaneously. The magnitude of the applied load is 15 MPa, which is equivalent to about 1.5 km of ice. The left end of the load coincides with the origin of the coordinate system at the beginning of the model run. See text for details.

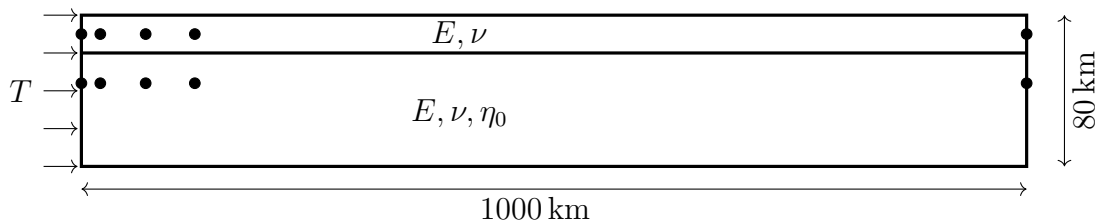


### 17.3.7 Kusznir and Bott (1977) experiment

The domain is a  $2000 \text{ km} \times 80 \text{ km}$  Cartesian box. However since there is a vertical axis of symmetry it is then reduced to  $1000 \text{ km} \times 80 \text{ km}$ . It consists of two layers: the top layer is  $20 \text{ km}$  thick and is purely elastic characterised by a Young's modulus  $E = 1.7 \times 10^{11} \text{ N m}^{-2}$  and a Poisson ratio  $\nu = 0.25$ . The lithosphere below is then  $60 \text{ km}$  thick and is characterised by an elasto-viscous rheology. The elastic parameters are identical to the upper layer while the viscosity is set to  $\eta_0 = 1 \times 10^{23} \text{ Pa s}$ .

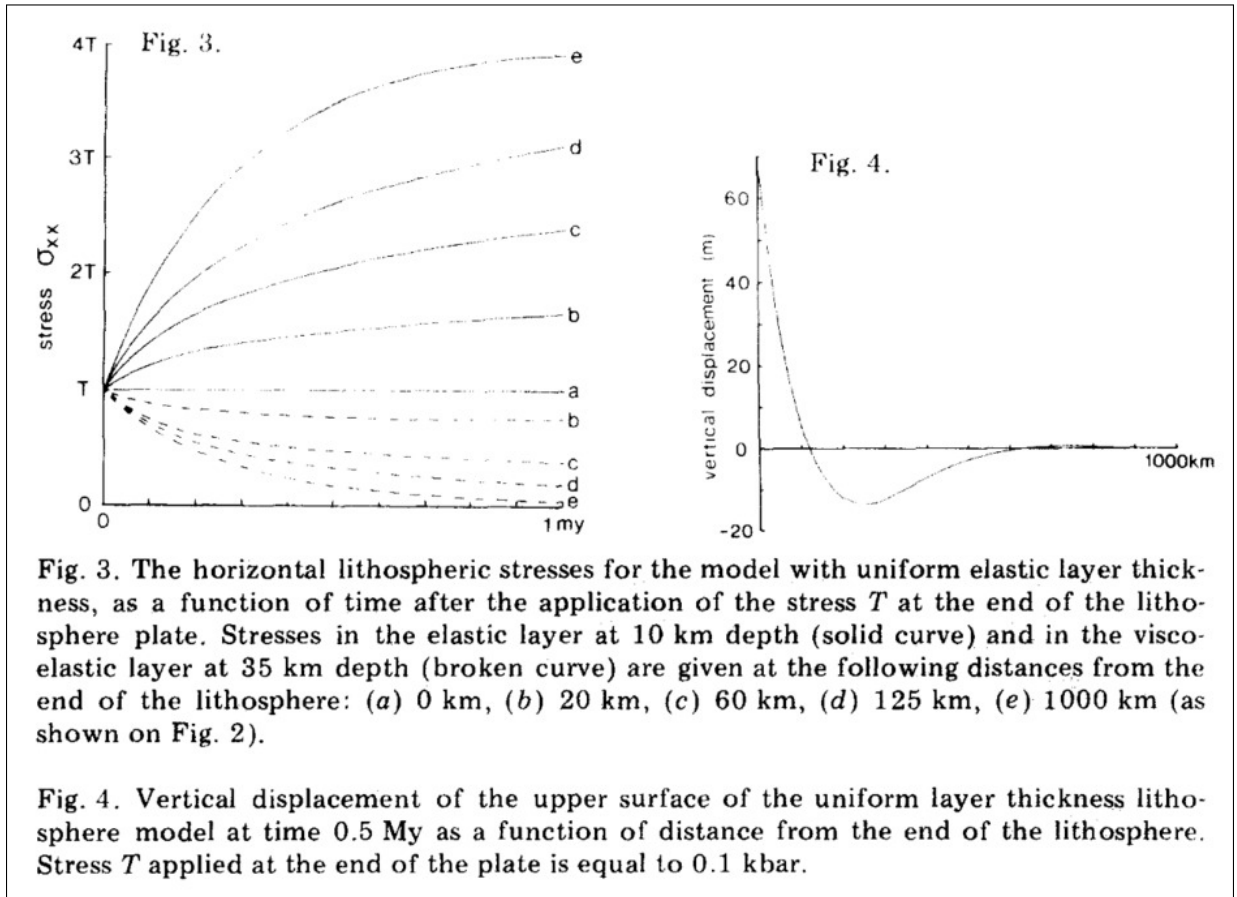
We arbitrarily design the right side as being the axis of symmetry and the boundary condition on that side are then free slip. On the left side a uniform horizontal stress  $T$  is applied at time  $t=0$ . The viscous drag of the asthenosphere below is neglected and the authors do not mention the bottom of the domain deforming so we'll also impose free slip boundary conditions at the bottom.

Fig. 4 of the paper shows the vertical displacement of the upper layer so we will use a free surface boundary condition at the top. The model is ran for  $1 \text{ Myr}$ .



After a thorough read of the paper, I have noticed quite a few problems:

- the paper is old and has been digitized but the figures are missing a lot of lines/shades/points ... this could be remedied by finding the article in a library.
- in the intro it is stated: "Here we investigate the response of a lithosphere divided into upper elastic and lower uniform visco-elastic layers to simple boundary force and body force systems." The authors later talk about 'isostatic forces' opposing flexure. This means that buoyancy forces should be taken into account but there is no information about densities or gravity values!
- little uncertainty about boundary conditions. is it free slip or no-slip on the right side ? Looking at fig 4, it looks like the vertical displacement is zero at  $x=1000\text{km}$ ?
- a Maxwell elasto-viscous rheology is used but this is only mentioned in the Appendix
- the dimensions of the thinned or thickened areas is simply not mentioned.
- Fig 1 shows triangular elements. No mention is made of resolution, type of element, ndofs, resolution tests, any numerical detail whatsoever. Given the age of the paper, I would guess  $P_1$  elements.
- In the appendix they equate the viscous strain rate to  $\sigma/4\eta$ . Why 4 ?
- it is also not clear whether the domain is ALE or fully Lagrangian: does it shorten?
- the value of  $T$  is never specified!
- the paper was published 45 years ago, it is extremely unlikely any of the two authors is still available



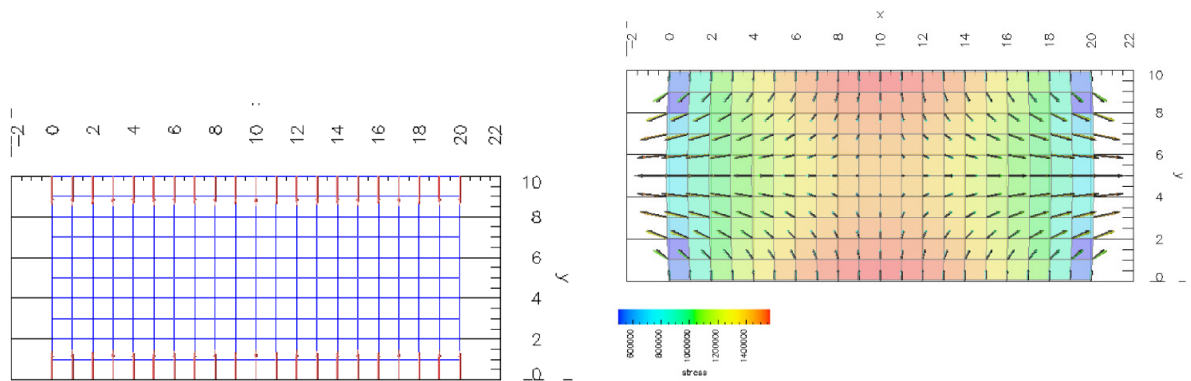
Taken from Kusznir and Bott [736]. Measurement locations are indicated on the setup figure above. probably should be reversed

### 17.3.8 Parallel-Plate Viscometer Problem - SNAC manual

A parallel-plate viscometer problem is simulated, in which viscoelastic material is squeezed between two parallel plates. The plates are moving at a constant velocity,  $v_0$ . Each plate has the length of  $2L$  and is at a distance  $2h$  from the other. No slip is assumed between the material and the plates. The approximate analytical solution is given in the book by Jaeger [631] (1969).

Model Setup:  $L = 10$  m,  $h = 5$  m, viscosity  $\eta = 10^9$  Pa s, bulk modulus  $K = 1.5$  GPa, shear modulus  $\mu = 500$  MPa,  $v_0 = 10^{-4}$  m s $^{-1}$ ,  $dt = 1$  s (results compared after 500 time steps), mesh size:  $20 \times 10$  m, each element is a 1 m cube.

Due to the assumption of the original problem setup, artificial forces should be added to left and right surfaces.



Left: The initial mesh (blue) with the velocity boundary condition (red arrows); Right: The second invariant of stress and velocities plotted on the deformed mesh. Colored arrows are for SNAC's solution, black ones for the analytic solution.



$$K = \lambda + \frac{2}{3}\mu \text{ so } \lambda = \frac{3500}{3} \text{ MPa} ?$$

From wikipedia<sup>1</sup>

$$\nu = \frac{3K - 2\mu}{2(3K + \mu)} = \frac{4500 - 2 * 500}{2(4500 + 500)} = \frac{3500}{10000} = 0.35$$

so we find that the material is **compressible**  $\nu = 0.35$

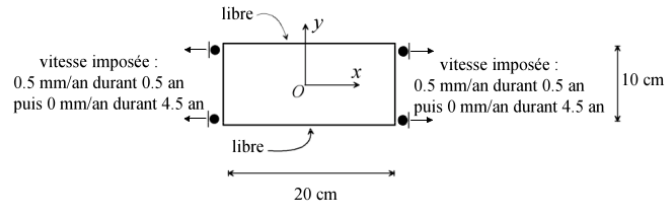
$$E = \frac{9K\mu}{3K + \mu} = \frac{9 * 1500 * 500 \text{ MPa}^2}{4500 + 500 \text{ MPa}} = \frac{9 * 1500 * 500}{5000} \text{ MPa} = 1350 \text{ MPa}$$

### 17.3.9 Relaxation after extention - Hassani syllabus

**compressible**  $\nu = 0.25$

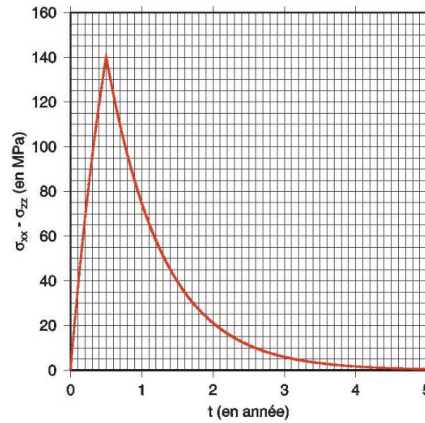
Un essai de relaxation consiste à imposer à un instant donné une déformation que l'on maintient constante par la suite. On observe alors comment évolue la contrainte au cours du temps.

The setup of the experiment is shown in the following figure:



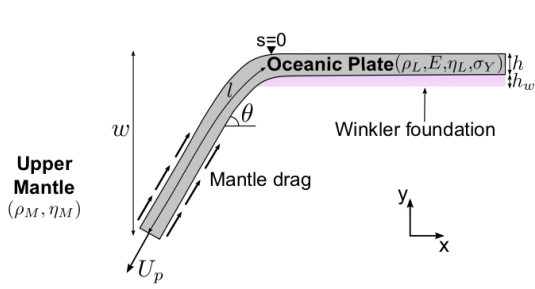
The total duration of the experiment is  $T = 5 \text{ yr} \simeq 15.75 \times 10^7 \text{ s}$ . The duration of the loading is  $T/10 = 6 \text{ months}$  while the duration of the subsequent relaxation is then  $9T/10 \simeq 4.5 \text{ yr}$ .

The loading velocity is  $v = 1 \text{ mm yr}^{-1} \simeq 3.17 \times 10^{-11} \text{ m s}^{-1}$ . The sample has size  $L \times L/2 = 20 \times 10 \text{ cm}$  and the strain rate is then  $v/L \simeq 1.5 \times 10^{-10} \text{ s}^{-1}$ . Young's modulus is set to  $1 \times 10^{11} \text{ Pa}$  and the Poisson ratio is 0.25, i.e.  $\mu = 40 \text{ GPa}$ . The viscosity is set to  $\eta_0 = 1 \times 10^{18} \text{ Pa s}$ . The Maxwell time is then  $t_M = \eta/\mu = 0.8 \text{ yr}$ , which is also the time it takes to reduce the maximum stress by a factor  $e$ .

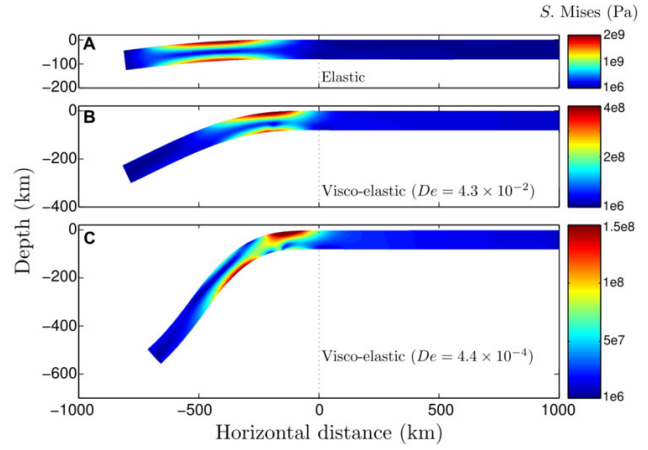


<sup>1</sup>[https://en.wikipedia.org/wiki/Lame\\_parameters](https://en.wikipedia.org/wiki/Lame_parameters)

### 17.3.10 Role of elasticity in slab bending - Fourel, Goes, and Morra [406] (2014)



**Figure 1.** Schematic of the model setup and definition of parameters. The slab of thickness  $h$ , density  $\rho_L$ , Young's modulus  $E$ , viscosity  $\eta_L$ , and in some cases yield stress  $\sigma_Y$  is freely sinking into an infinite mantle of density  $\rho_M$  and viscosity  $\eta_M$  achieving a subduction velocity  $U_p$  deflection  $w$  and slab dip  $\theta$ . Slab deformation is modeled in 2-D, but the analytically calculated mantle drag that is applied along the slab's sides and tip is for 3-D flow around a finite-width plate. The plate is supported in the top few kilometers by a Winkler foundation of thickness  $h_w$ . The top surface and the trailing edge are free.



### 17.3.11 Shear test in 2D - Farrington, Moresi, and Capitanio [387] (2014)

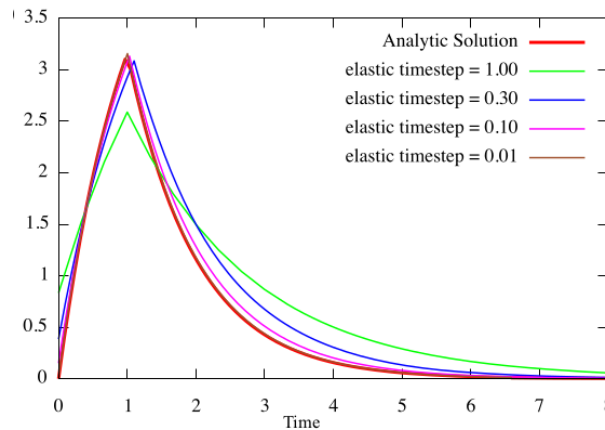
This experiment consists of a viscoelastic material undergoing simple shear at a constant rate from a time  $t_0$  though to  $t_{max}$ . At  $t_{max}$ , the shearing velocity is taken to zero using a no-slip velocity boundary condition. The viscoelastic stresses then decay with time while the material deformation rate remains zero. The  $xy$  component of stored stress, i.e. the nonviscous portion of the total stress, is given by

$$\begin{aligned} \tau_{xy}^{stor} &= \exp\left(-\frac{\mu}{\eta}t\right) \left(C_2 \cos\left(\frac{Vt}{h}\right) - C_1 \sin\left(\frac{Vt}{h}\right)\right) - C_2 & \text{if } t < t_{max} \\ &= \left[\exp\left(-\frac{\mu}{\eta}t_{max}\right) \left(C_2 \cos\left(\frac{Vt}{h}\right) - C_1 \sin\left(\frac{Vt}{h}\right)\right) - C_2\right] \exp\left(-\frac{\mu}{\eta}(t - t_{max})\right) & \text{if } t > t_{max} \end{aligned} \quad (17.6)$$

where  $V$  is the shear velocity along the top wall boundary,  $h$  is the height of the box,

$$C_1 = -\frac{V^2 \eta^2 \mu}{\mu^2 h^2 + V^2 \eta^2}$$

$$C_2 = -\frac{V h \eta \mu^2}{\mu^2 h^2 + V^2 \eta^2}$$



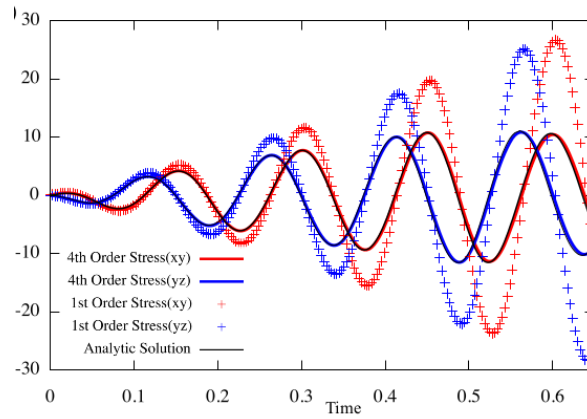
The  $xy$  component of dimensionless stored stress with dimensionless time for the 2-D viscoelastic material undergoing simple shear and relaxation. Nondimensional material parameters of  $\eta = 10^2$ ,  $\mu = 10^2$ ,  $t_M = 1$ ,  $V = 0.05$ ,  $h = 1$  with  $\Delta t_e \in [0.01, 1]$  and  $\Delta t_c = \frac{1}{3} \Delta t_e$ .

Figure above shows the resulting xy component of the deviatoric stored stress term for a 2-D viscoelastic material across a range of  $\Delta t_e$ . At longer  $\Delta t_e$ , the stored stress appears under resolved, indicating these larger elastic time step values capture dynamics that occur on time scales approximately equal to or greater than the Maxwell relaxation time, that is with a portion of viscous deformation. For shorter  $\Delta t_e$ , the numerically calculated stress approaches the analytic solution, indicating that the elastic time step is sufficiently small to fully capture the elastic stored stresses produced within the material under the applied strain rate.

### 17.3.12 Torsion test in 3D - Farrington, Moresi, and Capitanio [387] (2014)

insert here eq 9 of paper

The analytic solution outlined by equation (9) can be extended from this essentially 1-D test into 3-D by applying the shear velocity in the x-z plane. Testing of the full viscoelastic implementation including the rotation terms is possible by placing this 3-D shear test in a coordinate system under going solid body rotation.



The  $xy$  and  $yz$  stress components of a material undergoing simple shear within a 3D rotating reference frame. Nondimensional material parameters of  $\eta = 100$ ,  $\mu = 100$ ,  $\alpha = 1$ ,  $V = 0.3$ ,  $h = 1$ ,  $t_{max} = 0.5$  and  $\omega = 42$ . Note the coordinate system for this test has the  $y$  axis in the vertical with the  $z$  axis in plane. Results for a first (crosses) and fourth- (lines) order accurate Jaumann stress rate integration scheme are shown in comparison to the analytical solution (black) given by equation (9) within the rotating frame.

Figure above shows the evolution of the stress in comparison to the analytical solution of equation (9) placed within the rotating frame. The rotating frame is achieved by imposing a velocity boundary condition of a constant solid body rotation about the  $y$  axis in addition to the shearing rate. The stress within this rotating frame is given by equation (9), with the stress in the nonrotating frame found by applying a rotation matrix,  $R$ , to the nonrotating stress solution. That is,  $\tau = R^{-1}\tau'R$ , where  $'$  denotes the rotated frame,  $R$  is the rotation matrix about the  $y$  axis given by

$$R = \begin{pmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{pmatrix}$$

with  $\theta = \omega t$ ,  $\omega$  being the dimensionless rotation rate. The stress components shown in Figure 1b from within the nonrotating frame can then be given by  $\tau_{xy}^{stor} = \tau_{xy}^{stor'} \cos \theta$  and  $\tau_{yz}^{stor} = \tau_{yz}^{stor'} \sin \theta$ . It can be seen that, using a higher-order Jaumann stress rate advection scheme results in accurate stress advection and rotation within the full 3-D space plus time domain. It should be noted here that the velocity field used in this test was chosen to rigorously test the rotational terms in equation (6). Whether these rotational terms are required for individual models is dependent upon the model setup. For subducting slabs at a constant curvature, the individual parcels of material experience purely rotational effects, accounting for this within the stress history term would then be required for consistency.

finish!

### 17.3.13 Cylindrical tunnel - Segall book (?)

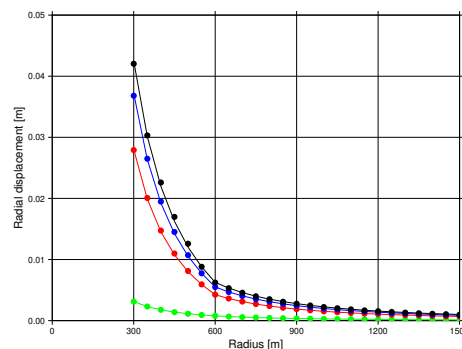
compressible  $\nu = 0.25$

communicated to me by L. van de Wiel.

This is a cylindrical tunnel (2D magma chamber) in an infinite space, with certain radius. The material close around the hole is warm, and and such viscous. The material further from the hole is cold, and purely elastic. (accomplished by setting and absurdly high viscosity) There is a clear transition radius between the two properties The hole contains a pressure, causing the space to expand.

After the initial elastic deformation, viscous deformation continues in the viscous region of the domain. I have the implementation of the analytical solution attached to save you time.

See plot with radial deformation for  $t=0$  (green),  $t=300s$  (red),  $t=600s$  (blue), and  $t=3000s$  (black). thin line is analytic, points are numeric. (with parameters and size parameters as in analytic.f)



### 17.3.14 Relevant literature & various notes

- convection of viscoelastic fluids: Harder [545], Moser, Matyska, Yuen, Malevsky, and Harder [909], Zhong, Gurnis, and Moresi [1410], Moresi, Dufour, and Mühlhaus [899], Mühlhaus and Regenauer-Lieb [915], Li and Khayat [784], Li and Khayat [785], Furuichi, Kameyama, and Kageyama [429]
- stress buildup associated with viscoelastic rheology Kusznir and Bott [736], Kusznir and Park [737], Poliakov, Cundall, Podlachikov, and Lyakhovsky [1007], Marques and Podladchikov [836]
- viscoelastic effects on geodynamical pbs involving gravitational instability Poliakov, Cundall, Podlachikov, and Lyakhovsky [1007], Kaus and Becker [680], Burov and Molnar [185], and Schmeling et al. [1124], Hanyk, Moser, Yuen, and Matyska [544]
- large strain eulerian viscoelasticity Schmalholz, Podladchikov, and Schmid [1117], Vasilyev, Podladchikov, and Yuen [1313], Cooper, Lenardic, Levander, Moresi, and Benn [278], Moresi, Quenette, Lemiale, Mériaux, Appelbe, and Mühlhaus [901], Furuichi, Kameyama, and Kageyama [429], and Popov and Sobolev [1011]

Farrington, Moresi, and Capitanio [387] states “Funiciello et al. [2003] implemented a viscoelastic rheology in numerical models of subduction, performing a range of numerical simulations to investigate its effect on subducting slab dynamics. Similar methodologies have addressed the details of viscoelastic stress within the bending zone during subduction, although a comparison between viscous and viscoelastic rheology was lacking [Capitanio et al., 2009; Capitanio and Morra, 2012]. Mühlhaus and Regenauer-Lieb [2005] and Moresi et al. [2002] have studied the role of elasticity in

mantle convection, comparing the viscous case to that of viscoelastic, and Kaus and Becker [2007] discussed the effect of elasticity on layered Rayleigh-Taylor instabilities. However, a systematic study into the effects of elastic stresses on models of free subduction has yet to be completed. Morra and Regenauer-Lieb [2006], Funicello et al. [2003], Capitanio et al. [2007], Yamato et al. [2007], and Royden and Husson [2006] have included a viscoelastic slab in subduction models, without explicitly studying the effects of the elastic component across a range of parameters.”

Check early paper by Braun & Beaumont (1987) [142]

Asgari and Moresi [31] (2012) Herwegh, Poulet, Karrech, and Regenauer-Lieb [566] (2014) Dansereau, Weiss, Saramito, and Lattes [305] (2016) Thielmann, Kaus, and Popov [1253] (2015) Beuchert, Podladchikov, Simon, and Rüpke [87] (2010) Sanan, May, Bollhöfer, and Schenk [1106] (2020) von Tscharnier and Schmalholz [1328] (2015) Naliboff, Lithgow-Bertelloni, Ruff, and Koker [928] (2012) Peltier [987] (1974)

# Chapter 18

## Geophysical data

## 18.1 The PREM model

prem.tex

**The density profile** Let us define  $x = r/R$ . Following table I of Dziewonski & Anderson (1981) [357] we have

- for the inner core  $0 < r < 1221.5\text{km}$  (or  $0 < x < 0.19172814314$ ):

$$\rho(x) = 13.0885 - 8.8381x^2$$

- for the outer core  $1221.5 < r < 3480\text{km}$  (or  $0.19172814314 < x < 0.5462250824$ ):

$$\rho(x) = 12.5815 - 1.2638x - 3.6426x^2 - 5.5281x^3$$

- for the Lower mantle  $3480 < r < 5701\text{km}$ :

$$\rho(x) = 7.9565 - 6.4761x + 5.5283x^2 - 3.0807x^3$$

- for the transition zone 1  $5701 < r < 5771\text{km}$ :

$$\rho(x) = 5.3197 - 1.4836x$$

- for the transition zone 2  $5771 < r < 5971\text{km}$ :

$$\rho(x) = 11.2494 - 8.0298x$$

- for the transition zone 3  $5971 < r < 6151\text{km}$ :

$$\rho(x) = 7.1089 - 3.8045x$$

- Low velocity zone  $6151 < r < 6291\text{km}$ :

$$\rho(x) = 2.6910 + 0.6924x$$

- LID  $6291 < r < 6346.6\text{km}$ :

$$\rho(x) = 2.6910 + 0.6924x$$

- Lower Crust  $6346.6 < r < 6356\text{km}$ :

$$\rho(x) = 2.9$$

- Upper Crust  $6356 < r < 6368\text{km}$ :

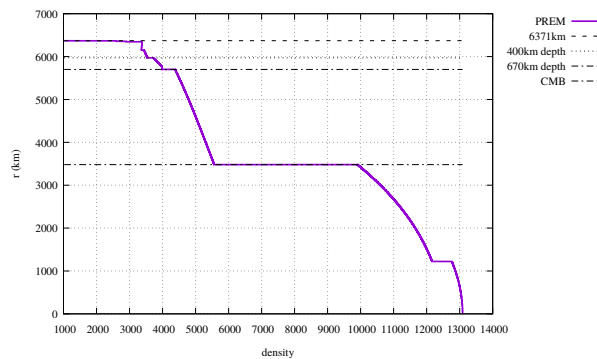
$$\rho(x) = 2.6$$

- Ocean  $6368 < r < 6371\text{km}$

$$\rho(x) = 1.020$$

Note that the returned densities should be multiplied by 1000 to obtain units of  $\text{kg/m}^3$ .

One can verify that the functions above yield the familiar PREM density profile:



**Corresponding gravity field** Following Eq. (10.39), the radial component of the gravitational acceleration at a position  $r$  outside of the Earth is given by:

$$g_r(r) = -\frac{1}{r^2} \int_0^r 4\pi\mathcal{G}\rho(r')r'^2 dr' = -\frac{4\pi\mathcal{G}}{r^2} \int_0^r \rho(r')r'^2 dr' \quad (18.1)$$

This integral can be broken up into layer integrals and we can compute the contribution of each layer to the gravity value  $g_r(r)$ . A layer is characterised by  $r_{in}$  and  $r_{out}$ . Inside  $r_{in}$  the gravity is zero. Between  $r_{in}$  and  $r_{out}$  the integral runs from  $r_{in}$  to  $r \leq r_{out}$ . Anout outside the layer the integral runs from  $r_{in}$  to  $r_{out}$ .

For simplicity the following integrals are computed for  $x \in [0, 1]$ :

- inner core

$$\begin{aligned} \text{inside} \quad g_{ic}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0^x (13.0885 - 8.8381x'^2)x'^2 dx' \\ &= \frac{4\pi\mathcal{G}}{x^2} (4.36283x^3 - 1.76762x^5) \\ &= 4\pi\mathcal{G}(4.36283x - 1.76762x^3) \end{aligned} \quad (18.2)$$

$$\begin{aligned} \text{outside} \quad g_{ic}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0^{0.19172814314} (13.0885 - 8.8381x'^2)x'^2 dx' \\ &\simeq 0.0302907 \frac{4\pi\mathcal{G}}{x^2} \end{aligned} \quad (18.3)$$

- outer core

$$\begin{aligned} \text{inside} \quad g_{oc}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_{0.19172814314}^x (12.5815 - 1.2638x' - 3.6426x'^2 - 5.5281x'^3)x'^2 dx' \\ &= 4\pi\mathcal{G}(-0.92135x^4 - 0.72852x^3 - 0.31595x^2 + 4.19383x - 0.028896(11/8.4)) \\ \text{outside} \quad g_{oc}(x) &= \frac{4\pi\mathcal{G}R^3}{x^2} \int_{0.19172814314}^{0.5462250824} (12.5815 - 1.2638x' - 3.6426x'^2 - 5.5281x'^3)x'^2 dx' \\ &\simeq 0.56663 \frac{4\pi\mathcal{G}}{x^2} \end{aligned} \quad (18.5)$$

- lower mantle

$$\begin{aligned} \text{inside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0^x ()x'^2 dx' \\ \text{outside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0{} ( )x'^2 dx' \end{aligned} \quad (18.6)$$

- transition zone 1

$$\begin{aligned} \text{inside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0^x ()x'^2 dx' \\ \text{outside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int_0{} ( )x'^2 dx' \end{aligned} \quad (18.7)$$



- transition zone 2

$$\begin{aligned}
\text{inside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int^x ()x'^2 dx' \\
\text{outside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int ()x'^2 dx'
\end{aligned} \tag{18.8}$$

- transition zone 3

$$\begin{aligned}
\text{inside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int^x ()x'^2 dx' \\
\text{outside} \quad g_{lm}(x) &= \frac{4\pi\mathcal{G}}{x^2} \int ()x'^2 dx'
\end{aligned} \tag{18.9}$$

$$\begin{aligned}
g_{lm}(r) &= \int_{3480/6371}^{5701/6371} (7.9565 - 6.4761x + 5.5283x^2 - 3.0807x^3)x^2 dx \simeq 0.904793 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{tz1}(r) &= \int_{5701/6371}^{5771/6371} (5.3197 - 1.4836x)x^2 dx \simeq 0.0354823 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{tz2}(r) &= \int_{5771/6371}^{5971/6371} (11.2494 - 8.0298x)x^2 dx \simeq 0.1026 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{tz3}(r) &= \int_{5971/6371}^{6151/6371} (7.1089 - 3.8045x)x^2 dx \simeq 0.0892215 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{lvz} &= \int_{6151/6371}^{6291/6371} (2.6910 + 0.6924x)x^2 dx \simeq 0.0705516 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{lid} &= \int_{6291/6371}^{6346.6/6371} (2.6910 + 0.6924x)x^2 dx \simeq 0.0289968 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{lc} &= \int_{6346.6/6371}^{6356/6371} 2.9x^2 dx \simeq 0.00425234 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_{uc} &= \int_{6356/6371}^{6368/6371} 2.6x^2 dx \simeq 0.00488337 \frac{4\pi\mathcal{G}R^3}{r^2} \\
g_o &= \int_{6368/6371}^1 1.020x^2 dx \simeq 0.000480075 \frac{4\pi\mathcal{G}R^3}{r^2}
\end{aligned} \tag{18.10}$$

Finally

$$\begin{aligned}
|g_r(r)| &= g_{ic}(r) + g_{oc}(r) + g_{lm}(r) + g_{tz1}(r) + g_{tz2}(r) + g_{tz3}(r) + g_{lvz}(r) + g_{lid}(r) + g_{lc}(r) + g_{uc}(r) + g_o(r) \\
&= \frac{4\pi\mathcal{G}R^3}{r^2} (0.0302907 + 0.56663 + 0.904793 + 0.0354823 + 0.1026 + 0.0892215 + \\
&\quad 0.0705516 + 0.0289968 + 0.00425234 + 0.00488337 + 0.000480075) \\
&\simeq \frac{4\pi\mathcal{G}R^3}{r^2} 1.838181685
\end{aligned}$$

At the surface of the Earth,  $r = R$  so we arrive at (after multiplying by 1000, see comment above):

$$g_{PREM}(R) \simeq 4\pi \cdot 6.67408 \times 10^{-11} \cdot 6371 \times 10^3 \cdot 1.838181685 \simeq 9.82194$$

All these calculations should be rechecked, although obviously the obtained value makes much sense.

## 18.2 From 1D tomography to density/temperature

The mantle is heterogeneous but it is also inaccessible. This means that one must rely on indirect methods to probe its structure. Seismic tomography is a technique for imaging the subsurface of the Earth with seismic waves produced by earthquakes or explosions. P-, S-, and surface waves can be used for tomographic models of different resolutions.

Seismic velocity is a meaningful parameter for the interior dynamics of the Earth because there exists a direct relation between seismic velocity and density. Such a relation was analysed experimentally by (for instance) Barton (1986) who used laboratory measurements of P-wave seismic velocity and density of rocks [51].

Fourty years later or so, a crucial question remains: what is the exact form of the relation between density and seismic velocity for the entire Earth's mantle?

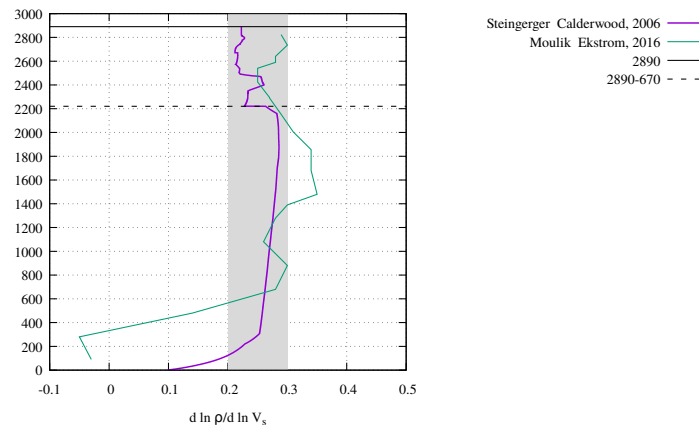
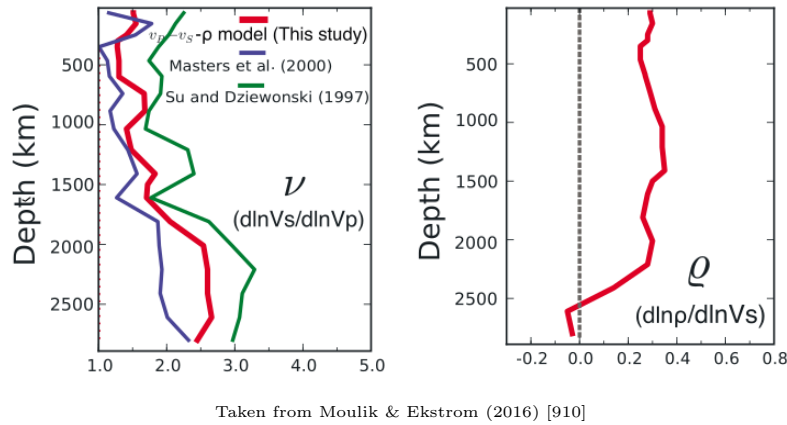
I will here not go into the details of the underlying theories and their approximations but will show a few useful results.

From tomography to density, the workflow is usually as follows:

$$d \ln V_p \rightarrow d \ln V_s \rightarrow d \ln \rho \rightarrow d \rho$$

Note that if the method is based on shear wave tomography the conversion  $d \ln V_p \rightarrow d \ln V_s$  is not necessary. Also the last step  $d \ln \rho \rightarrow d \rho$  requires a background density field, often taken to be either the PREM model of AK135 (see Section 18.1).

On the following plots are shown radial averages of the ratio  $d \ln V_s / d \ln V_p$  and  $d \ln \rho / d \ln V_s$ :



Profiles of scaling factor  $\xi = d \ln \rho / d \ln V_s$ . Data from Steinberger & Calderwood (2006) [1203] and Moulik & Ekstrom (2016) [910]. Data available in [images/dlnrhodlnvs/](#)

Then


$$\delta \ln(\rho(r, \theta, \phi)) = \xi(r) \cdot \delta \ln(V_s(r, \theta, \phi))$$

with

$$\delta \ln(\rho(r, \theta, \phi)) = \frac{\delta \rho(r, \theta, \phi)}{\rho_{ref}(r)}$$

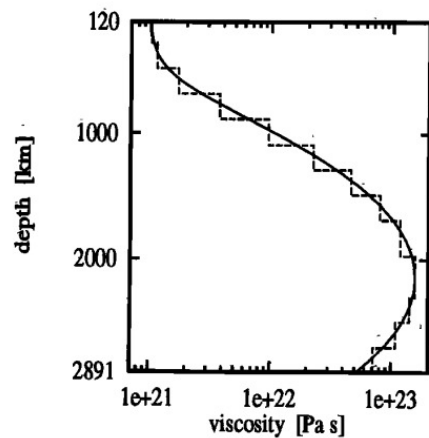
where  $\rho_{ref}$  is a radial profile, PREM for instance in the case of S40RTS so finally

$$\delta \rho(r, \theta, \phi) = \rho_{ref}(r) \cdot \delta \ln(\rho(r, \theta, \phi)) = \rho_{ref}(r) \cdot \xi(r) \cdot \delta \ln(V_s(r, \theta, \phi))$$

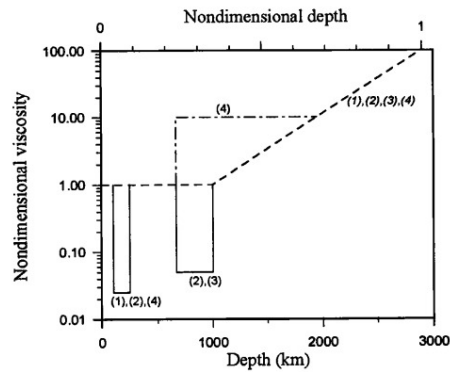
 Relevant Literature: [1083]

# 18.3 Earth radial viscosity profile

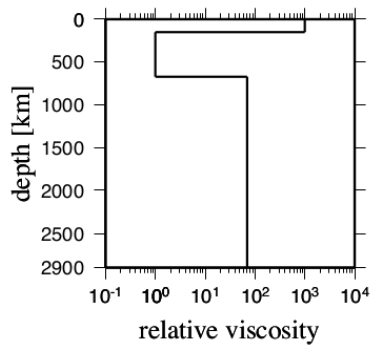
viscosity\_profile.tex



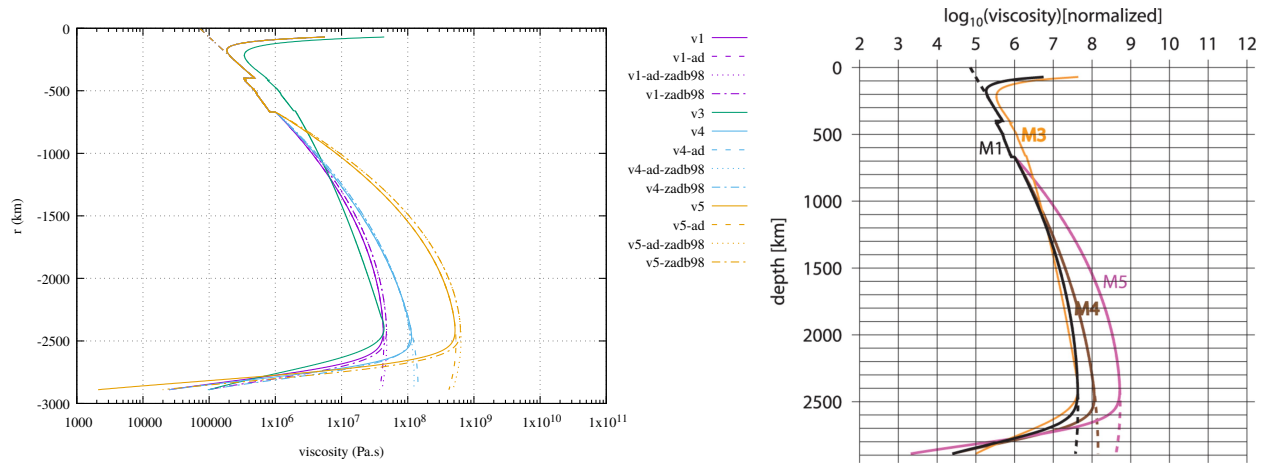
1995: Taken from Hanyk *et al.* (1995) [544].  $\eta(z) = (1 + 214.3z \exp -16.7(0.7 - z)^2) \times 10^{21}$  pascals



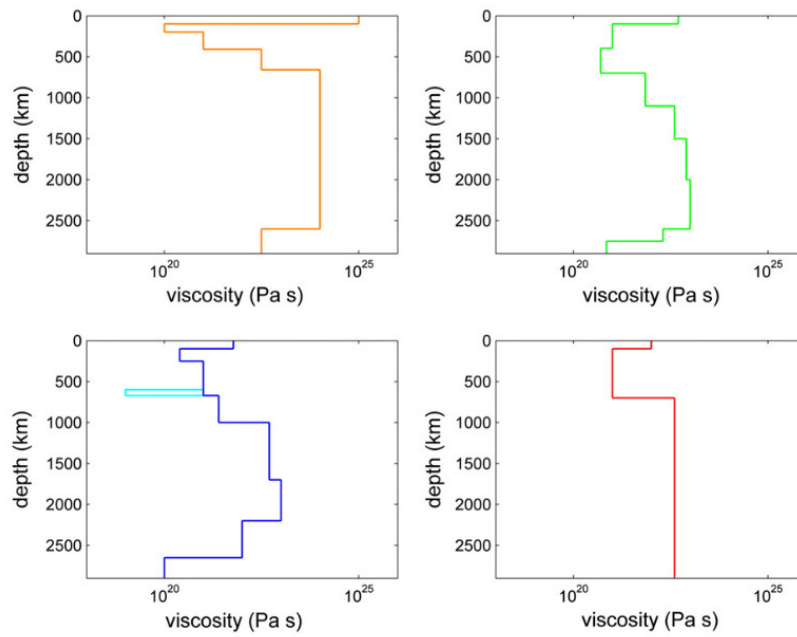
1997: Taken from Cserepes & Yuen (1997) [292]



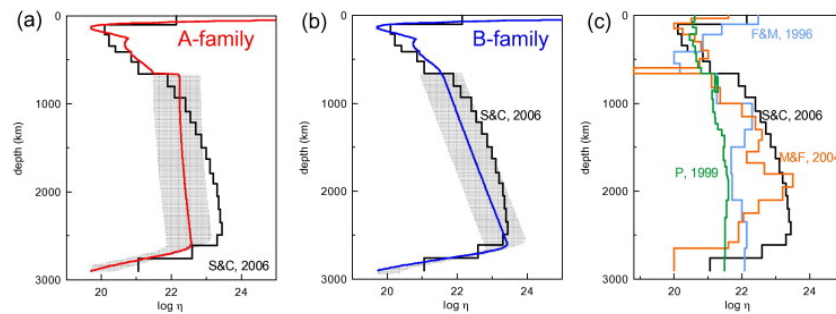
2001:Radial viscosity profile of the reference model. 3-layered model is adopted: the lithosphere (0 km to 150 km), the upper mantle (150 km to 670km) and the lower mantle (670 km to 2900 km). Taken from [1388]



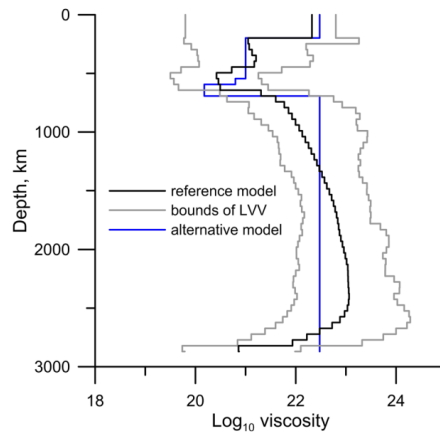
**2006:** Non-optimized, normalized viscosity profiles, as sent by B. Steinberger, corresponding to fig. 4 in Steinberger & Calderwood (2006) [1203].



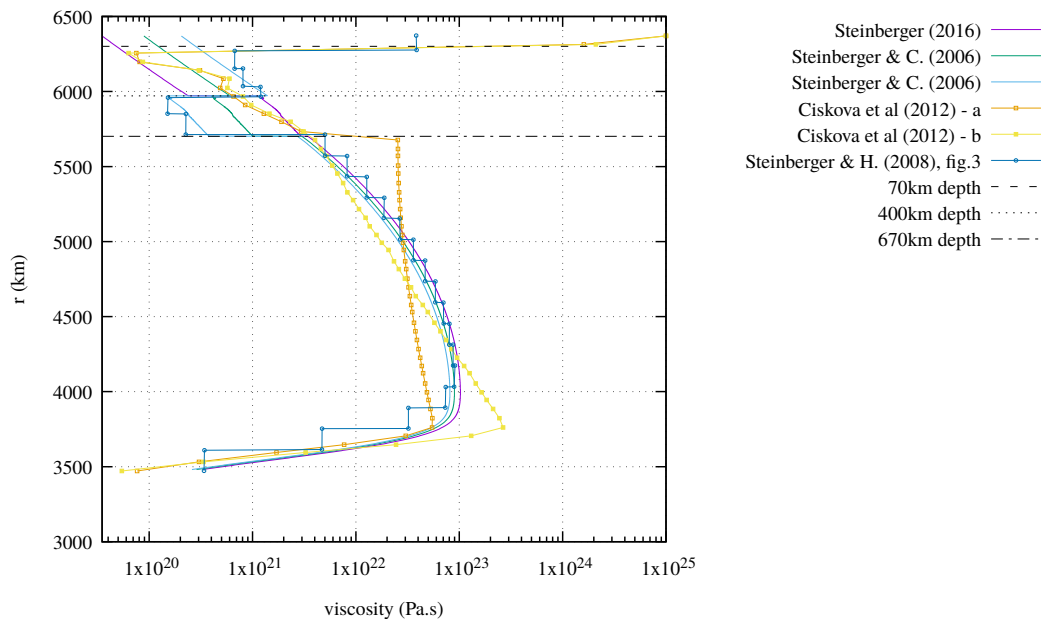
**2011:** Taken from Cadio *et al.* [200]. Mantle viscosity structure employed in calculating synthetic geoid anomalies. Red: VR (Ricard *et al.* , 1993); Blue: VMF (Mitrovica and Forte, 2004); Cyan: VMF-LVZ (Mitrovica and Forte, 2004); Green: VSC (Steinberger and Calderwood, 2006); Orange: VYN (Yoshida and Nakakuki, 2009).



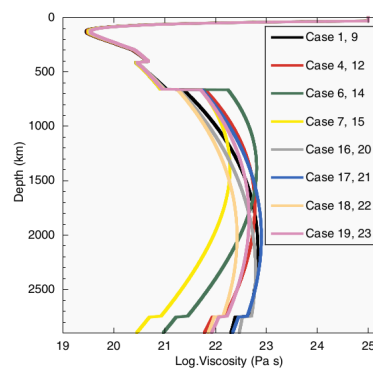
**2012:** Taken from Ciskova *et al.* [259].



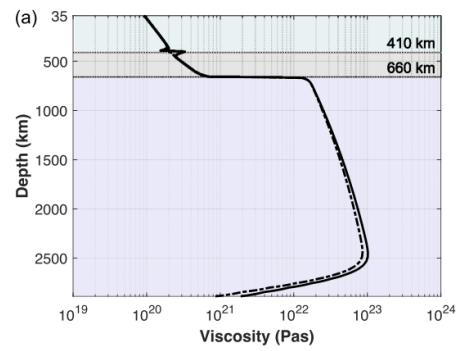
2014 Taken from Kaban *et al.* [660]. (Black) reference radial viscosity from (Steinberger and Calderwood 2006); (Blue) alternative viscosity model of Kaban and Trubitsyn (2012); (Light Grey) limits of viscosity variations in the model with LVV (Petrinin *et al.* 2013).



Steinberger (2016) [1204]; Ciskova *et al.* (2012) [259].



2019: Taken from Kaneko *et al.* [669].



2023: Taken from Neuharth and Mittelstaedt [934].

#### Relevant Literature:

- Viscosity profile of the lower mantle [370]
- Matyska *et al.* (2011) [844]
- Flament (2019) [396]
- Mitrovica & Forte [880]
- King & Masters [706]
- Rudolph *et al.* [1089]
- steinberger & holme [1202]
- steinberger [1204]
- Cadek & Yuen [198]
- supp of [49]



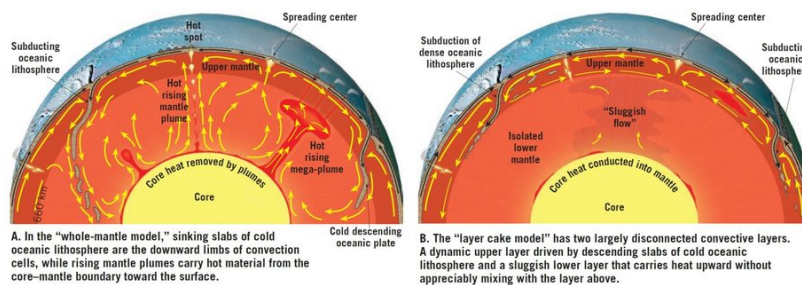
## 18.4 Earth radial temperature profile

adiabatic.tex

The method of manufactured solutions is a relatively simple way of carrying out

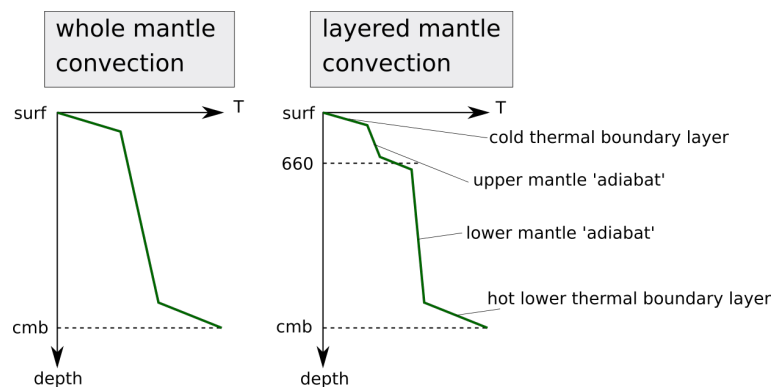
We start by the first sentence of Bunge (2005) [170]: "The average temperature increase through Earth's crust and mantle is called the geotherm. Its basic form is assumed to consist of adiabatic regions where temperatures rise only slightly with depth, and of narrow thermal boundary layers where temperatures increase rapidly over a depth of a few hundred kilometers (Jeanloz and Morris, 1986) [642]."

Before we look further at the equations behind this temperature profile, we must look at the basic assumption we make about the type of convection taking place in the mantle: layered convection vs. whole mantle convection, as depicted here:

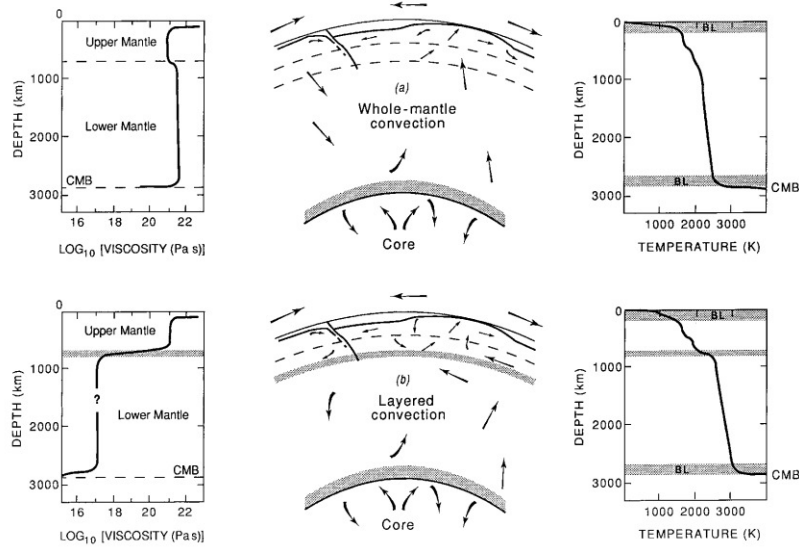


Taken from <https://geologyengineering.com/2020/05/mantle-convection/>

Indeed, the type of convection is then expected to have an strong influence on the radial temperature profile:

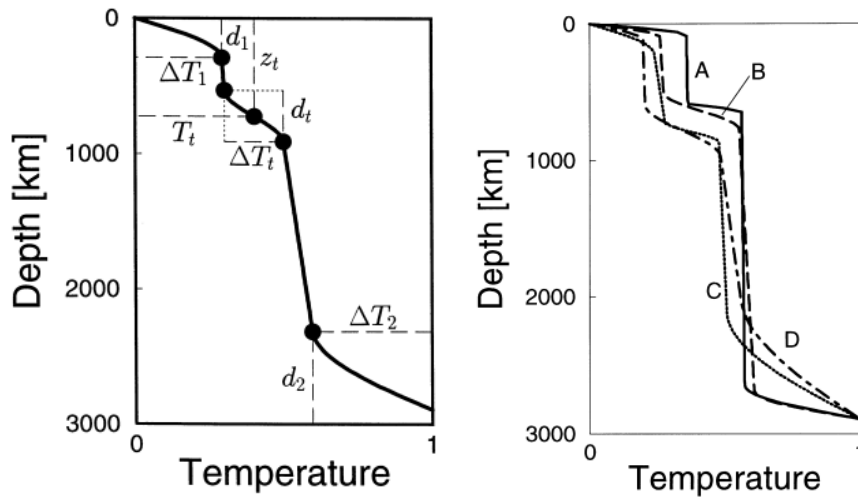


This is also to be found in Poirier's book [1006]:



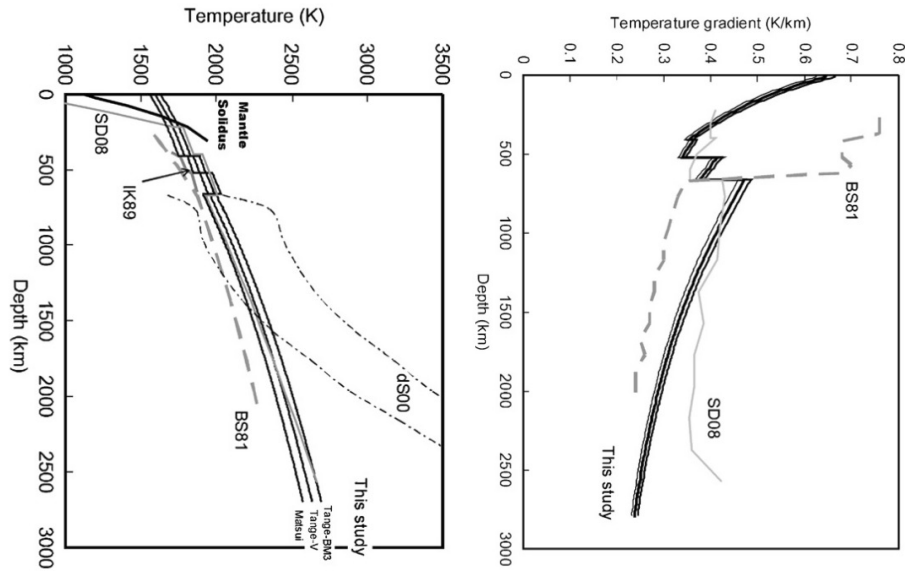
Schematic diagrams of (a) whole-mantle and (b) layered convection models, with corresponding temperature and viscosity profiles (after Peltier & Jarvis (1982) [989].)

The two-layered convection hypothesis relies essentially on the viscosity jump at the 660 discontinuity and seems to be generally accepted. It also forms the basis of Čadek & van den Berg (1998) [199] in which the authors carry out an inversion to obtain radial profiles of temperature and viscosity in the Earth's mantle inferred from the geoid and lateral seismic structure:

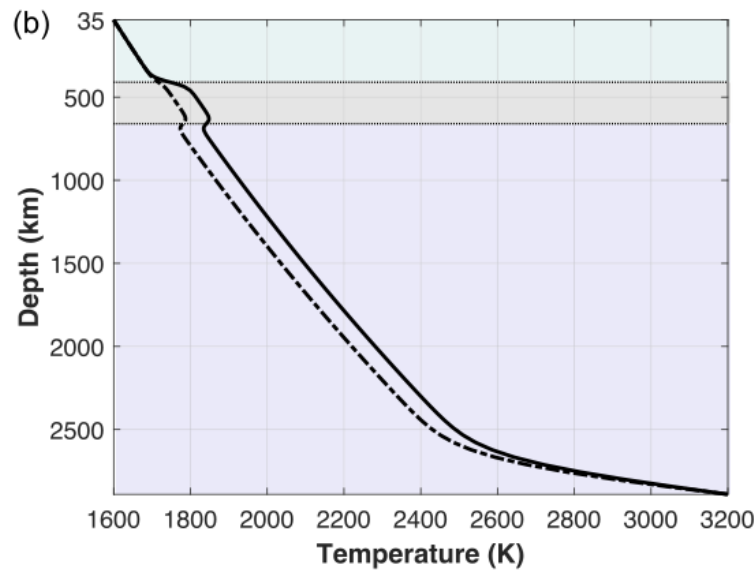


Taken from [199]. Left: Parameterization of the geotherm used in the paper; Right: Four model geotherms reducing the misfit by 70%. The differences between the models illustrate uncertainties of the solution.

More recently, Katsura *et al.* (2010) [676] have constructed the following mantle temperature and temperature gradient profiles:

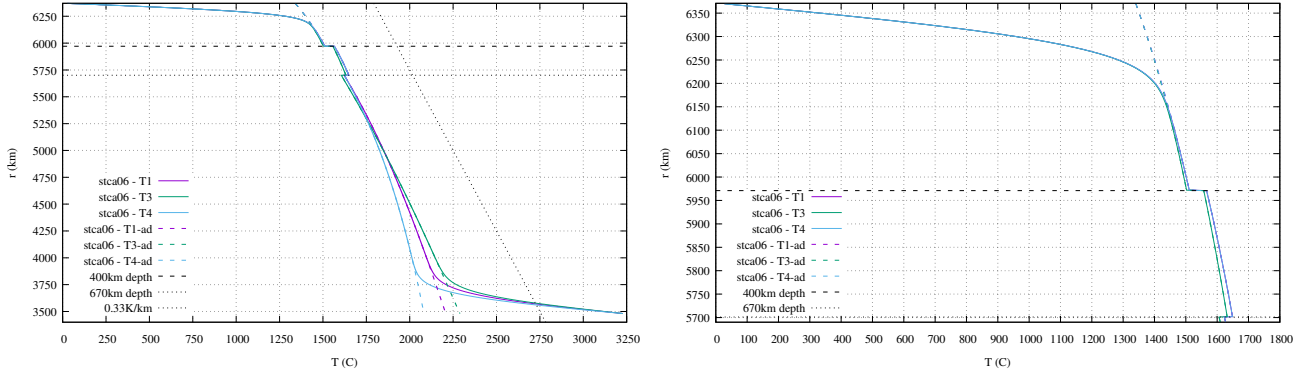


Taken from Katsura *et al.* (2010) [676]. Left: The adiabatic temperature distributions in the mantle. The three solid lines denote the temperature distributions proposed in this study using three different pressure scales. Those proposed by the previous studies are shown for comparison (BS81: Brown and Shankland, 1981; IK89: Ito and Katsura, 1989; dS00: da Silva *et al.*, 2000; SD08: Stacey and Davis, 2008). The mantle solidus proposed by Hirschmann (2000) is also shown. Right: Adiabatic temperature gradient in the mantle. The adiabatic temperature gradient abruptly increases in association with the olivine-wadsleyite, wadsleyite-ringwoodite, ringwoodite-perovskite + periclase transitions, as is the case for the thermal expansion. The adiabatic temperature gradients given in the previous studies are also shown for comparison (BS81: Brown and Shankland, 1981; SD08: Stacey and Davis, 2008 [1191]).



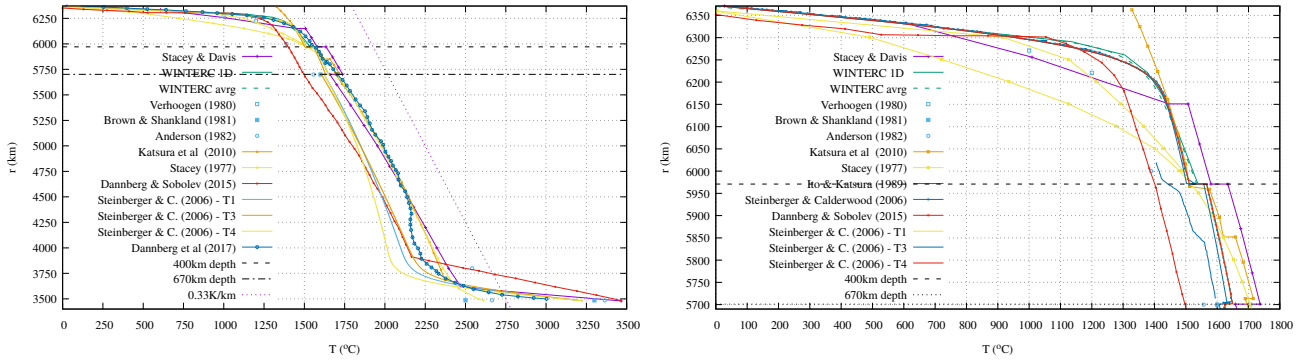
Taken from Neuharth and Mittelstaedt [934].

Prof. Steinberger was gracious to communicate to me the data of Steinberger & Calderwood (2006) [1203]:



Temperature profiles from Steinberger & Calderwood [1203].

One can gather data from various papers and books and this yields the following figure<sup>1</sup>:

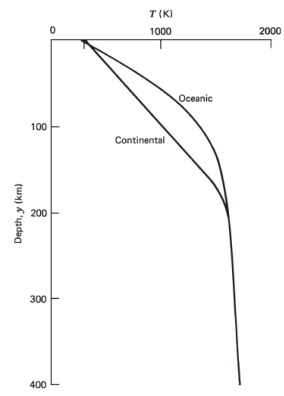


Note that the T values of Stacey (1977) [1190] inexplicably go to 540-550K at the surface. The value 545 has therefore been subtracted from the values in the paper. These then align remarkably well with those of Katsura *et al.*

We see that there is a remarkable agreement between many studies for the temperature down to 400km depth. Most show a step in the profile between 400 and 670km depth but the temperature at 670km depth seems to be between 1650°C and 1750°C. Further down, the discrepancy between studies increases (effectively boiling down to the value of the adiabatic gradient). 100km above the CMB all studies seem to indicate a temperature of 2300-2400°C. Temperature values at the CMB differ a lot and as noted in Bull *et al.* (2009) [167]: "Estimates of CMB temperature vary greatly. [...] We impose a surface temperature of 273 K and investigate three CMB temperatures: 3000 K (consistent with previous studies of this nature, e.g., Kellogg *et al.* , 1999 [693]), 3950 K (consistent with recent results on the double-crossing of Post- Perovskite, e.g., Hernlund *et al.* , 2005 [565]; van der Hilst *et al.* , 2007 [572]; Alfè *et al.* , 2002), and 4800 K (Knittle and Jeanloz, 1991). The 3950 K temperature is the 'reference' temperature which we use for the majority of cases." The CMB temperature is set to 2900 – 3200°C in Fowler [408].

Also, oceanic plates being thinner than continental plates (alongside with different radiogenic decay and thermal properties) one can draw two representative mantle geotherms [1288]:

<sup>1</sup>It must be said that finding actual data -not figures- is remarkably difficult.



**Figure 4.56** Representative oceanic and continental shallow upper mantle geotherms.

Taken from Turcotte & Schubert [1288]

**Isentropic gradient** The material properties also define the slope of the adiabat (the change in temperature with pressure at constant entropy) at all pressures and temperatures. Using the cyclic relation, we can define this slope in terms of partial differentials of the entropy with respect to pressure and temperature:

$$\left(\frac{\partial T}{\partial p}\right)_S = -\left(\frac{\partial T}{\partial S}\right)_p \left(\frac{\partial S}{\partial p}\right)_T \quad (18.11)$$

$$= -\left(\frac{T}{C_p}\right) \left(-\frac{\alpha}{\rho}\right) \quad (18.12)$$

$$= \frac{\alpha T}{\rho C_p} \quad (18.13)$$

This expression does not pose a constraint on the material properties, but in order to be self-consistent, the adiabat must be computed following this relation.

For complex material models, obtaining analytical functions which obey all these relations may be a non-trivial exercise. Furthermore, it is often not immediately clear when a given formulation is thermodynamically inconsistent. Indeed, both the thermodynamic and the geodynamic literature contain many equations of state and material parameterizations which do not obey these relations! This may not invalidate the results obtained with these models, but it is a point worth keeping in mind as the geodynamics community moves to more complicated and more realistic parameterizations.

*A final note of warning: Some compressible formulations in ASPECT (Section ??) use the isothermal compressibility, while others use the isentropic compressibility. Fully self-consistent material models must either specify what approximation of the compressible equations they are consistent with (see Section ??), or have a switch so that they use the correct compressibility for each of the different approximations. The conversion between isothermal and isentropic compressibilities is given in (??).*

---

**Initial conditions and the adiabatic pressure/temperature** The thermo-mechanically coupled (Navier-)Stokes equations require us to pose initial conditions for the temperature. Note that the equations themselves do not require that initial conditions are specified for the velocity and pressure variables (since there are no time derivatives on these variables in the model).

Nevertheless, a nonlinear solver will have difficulty converging to the correct solution if we start with a completely unphysical pressure for models in which coefficients such as density  $\rho$  and viscosity  $\eta$  depend on the pressure and temperature. To this end, ASPECT uses pressure and temperature fields  $p_{\text{ad}}(z), T_{\text{ad}}(z)$  computed in the adiabatic conditions model (see Section ??). By default, these fields satisfy adiabatic conditions:

$$\rho C_p \frac{d}{dz} T_{\text{ad}}(z) = \frac{\partial \rho}{\partial T} T_{\text{ad}}(z) g_z, \quad (18.14)$$

$$\frac{d}{dz} p_{\text{ad}}(z) = \rho g_z, \quad (18.15)$$

where strictly speaking  $g_z$  is the magnitude of the vertical component of the gravity vector field, but in practice we take the magnitude of the entire gravity vector.

These equations can be integrated numerically starting at  $z = 0$ , using the depth dependent gravity field and values of the coefficients  $\rho = \rho(p, T, z), C_p = C_p(p, T, z)$ . As starting conditions at  $z = 0$  we choose a pressure  $p_{\text{ad}}(0)$  equal to the average surface pressure (often chosen to be zero, see

Section ??), and an adiabatic surface temperature  $T_{\text{ad}}(0)$  that is also selected in the input parameter file.

**Note:** The adiabatic surface temperature is often chosen significantly higher than the actual surface temperature. For example, on earth, the actual surface temperature is on the order of 290 K, whereas a reasonable adiabatic surface temperature is maybe 1600 K. The reason is that the bulk of the mantle is more or less in thermal equilibrium with a thermal profile that corresponds to the latter temperature, whereas the very low actual surface temperature and the very high bottom temperature at the core-mantle boundary simply induce a thermal boundary layer. Since the temperature and pressure profile we compute using the equations above are simply meant to be good starting points for nonlinear solvers, it is important to choose this profile in such a way that it covers most of the mantle well; choosing an adiabatic surface temperature of 290 K would yield a temperature and pressure profile that is wrong almost throughout the entire mantle.

For instance, let us consider  $\alpha = 3 \cdot 10^{-5}$ ,  $g_z = 10$ ,  $C_p = 1250$ ,  $\rho = \rho_0(1 - \alpha(T - T_0))$  so that  $\frac{\partial \rho}{\partial T} = -\alpha \rho_0$  with  $\rho_0 = 3300$ .

Then we must solve the following equation

$$\rho_0(1 - \alpha(T - T_0))C_p \frac{dT}{dz} = -\alpha T^2 g_z$$

---

In Verhoogen (1951) [1319]: As is well known, the adiabatic gradient may be written as

$$\frac{dT}{dP} = \alpha T / \rho C_p$$

If hydrostatic equilibrium is assumed, the pressure varies with depth  $h$  as  $dP = \rho g dh$ , so that

$$\frac{d \ln T}{dh} = \frac{\alpha g}{C_p}$$

from which the temperature  $T$  at any depth  $h$  may be computed as a function of the temperature at any assigned depth if the ratio  $\alpha/C_p$  is known at all depths ( $g$ , the acceleration of gravity, will be taken as constant in the mantle).

---

From DyMaLi: In the interior of a convecting medium temperatures follow an adiabatic profile. At the top and bottom of a convecting layer thermal boundary layers with large thermal gradients form. The interior is thermally well mixed and therefore essentially isothermal, with a slight increase of temperatures with depth due to the effect of pressure. For example in the Earth's mantle the geothermal gradient  $\partial T / \partial z$  is about 20C/km near the surface and about 0.3C/km in the interior of the mantle. This small gradient in the interior is the adiabatic gradient. If a small volume of material is moved to shallower depth it experiences a slight increase in volume due to the decreasing pressure and associated with this a slight decrease in temperature. This change in temperature is the adiabatic temperature change.

The adiabatic gradient can be determined from the thermodynamics relation between entropy per unit mass  $S$ , temperature  $T$ , and pressure  $P$ :

$$dS = \left( \frac{dS}{dT} \right)_P dT + \left( \frac{dS}{dP} \right)_T dP = \frac{C_p}{T} dT - \frac{\alpha}{\rho} dP$$

In case of a reversible adiabatic process the entropy change is zero, and so the adiabatic gradient is:

$$\left(\frac{dT}{dP}\right)_s = \frac{\alpha T}{\rho C_p}$$

The gradient can also be expressed in terms of depth, remembering that  $dp = \rho g dz$  in a hydrostatic fluid:

$$\left(\frac{dT}{dz}\right)_s = \frac{\alpha g T}{C_p}$$

Thus to determine the adiabatic gradient one needs values of  $\alpha$  and  $C_p$  with depth. These are obtained from laboratory experiments.

One also needs an estimate of density as a function of depth, which is generally determined from seismology. Integration of the adiabatic gradient in terms of pressure then gives temperature as a function of pressure. Temperature as a function of depth is obtained by integrating the density distribution to obtain  $g$  as a function of depth.

not finished —————

Vol07\_02

This adiabaticity hypothesis should, however, not be taken too literally (Jeanloz and Morris, 1987 [641]). In most numerical simulations, the resulting averaged geotherm can be far (a few hundred kelvins) from adiabatic (Bunge *et al.*, 2001 [171]). First, radioactive heating, dissipation, and diffusion are never totally negligible, second, even if each fluid parcel follows its own adiabatic geotherm, the average geotherm may not correspond to any particular adiabat.

————— From Steinberger oct 26th, 2020:

Yes, I think the temperature below the lithosphere is still fairly well-constrained, based on magmas produced on mid-oceanic ridges (away from plumes) if one takes these as representative. Regarding Dannberg and Sobolev (not Solomatov!) I don't know why it is lower. In their supplementary figures, the extrapolation to the surface is actually  $\sim 1250\text{K}$ , not  $1250^\circ\text{C}$ , which is even less. Perhaps best to directly ask Julianne about this.

How the temperature increases with depth is more uncertain; the temperature gradient is often taken as adiabatic; that's what I also assumed in the 2006 paper, then it is defined by a ordinary differential equation (eq. 12 in that paper). I think the main uncertainty is the thermal expansivity. In my model, it strongly decreases with depth, I guess that is why my models have a lower temperature in the lower mantle than other models. I think that strong decrease in expansivity is still the consensus, but I didn't closely follow the literature recently. But the temperature gradient between the thermal boundary layers may actually be subadiabatic, hence even lower, due to cold slabs sinking to and accumulating above the CMB, and hot plume material feeding into the asthenosphere, below the lithosphere. So, yes, I would say the difference in the deep mantle (if not more) reflects the current state in the community.

And I think the uncertainties of CMB temperature are even larger. And I would of course be happy if you host my data on a public github repo. Just for the viscosity profile, there should be some explanation given that it shouldn't necessarily be taken "at face value" but that (like in the 2006 paper) it is possible to multiply different parts of the profile with different "scaling viscosities".



**Relevant Literature:** *On the thermal gradient in the Earth's deep interior*, Tirone (2016) [1270]  
*Is the mantle geotherm subadiabatic*, Jeanloz & Morris (1987) [641]  
*Mantle convection, the asthenosphere, and Earth's thermal history* King (2015) [704]

check section 7.7 in Fowler !

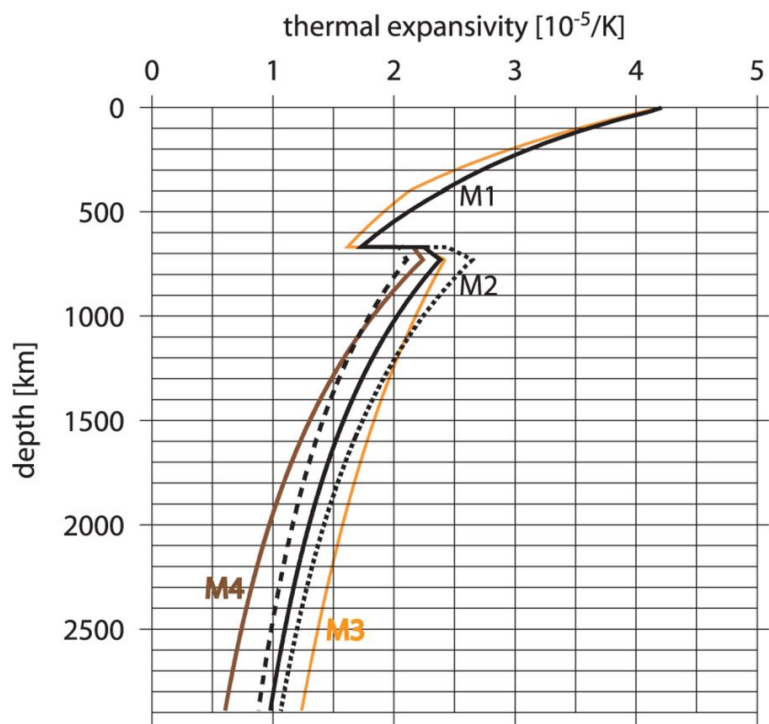
boundary layer jape82

evolution shpe79



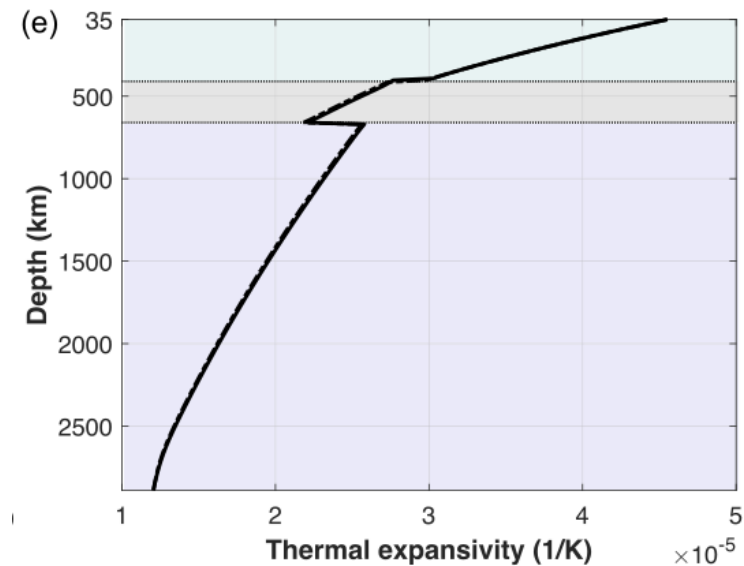
# 18.5 Earth radial thermal expansion profile

thermal\_expansion\_profile.tex



Taken from Steinberger and Calderwood [1203].

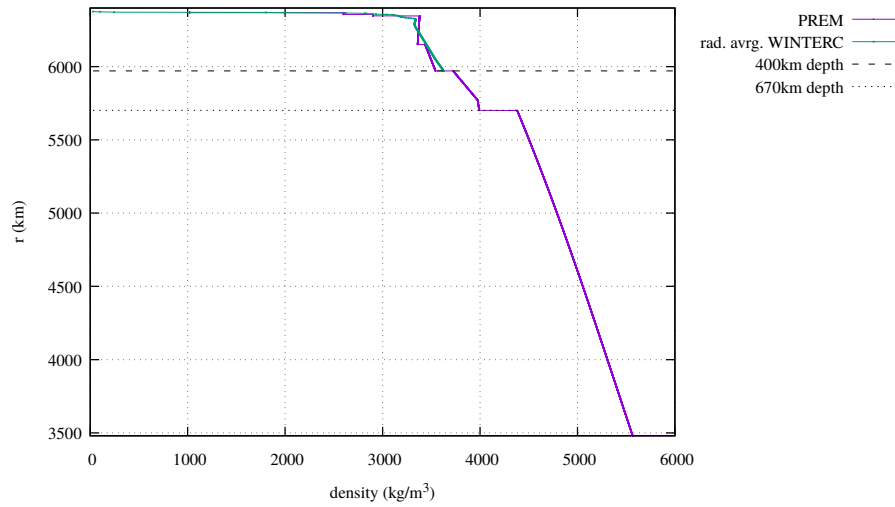
see Matyska *et al.* (2011) [844]  
Mantle convection with internal heating and pressure-dependent thermal expansivity, Leitch *et al.* (1991) [763]  
Eq(8) of Hassan *et al.* [552]




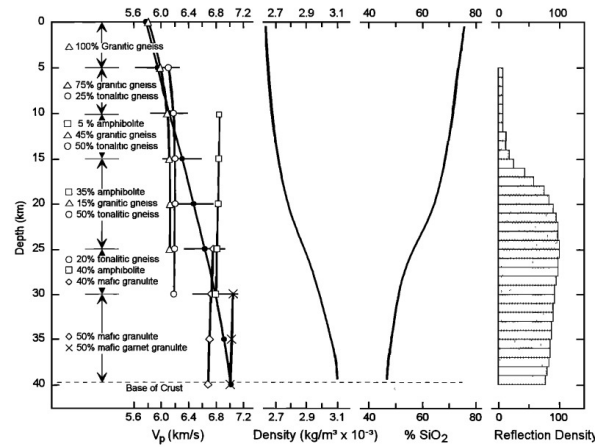
Taken from Neuharth and Mittelstaedt [934].

## 18.6 Earth radial density profile

density\_profile.tex



 Relevant Literature: Kennett (1998) [695]



Taken from Christensen and Mooney [240] (1995). A model for average crustal petrology versus depth consistent with average velocity depth profile (solid circles) and velocity depth curves for common rock types (open symbols). Variations of density and SiO<sub>2</sub> content with depth are from rock percentages shown on the left.

Let us look at the density and pressure profiles in a 1D isothermal 'planet'. We start from

$$-\vec{\nabla} p + \rho \vec{g} = \vec{0}$$

In 1D, and assuming  $\vec{g} = -g\vec{e}_z$ :

$$-\frac{dp}{dz} - \rho g = 0$$

or

$$\frac{dp}{dz} = -\rho g$$

Assuming  $\rho$  and  $g$  constant in the domain  $z \in [0, L]$ , we can solve this ODE and we obtain:

$$p(z) = \rho g(L - z)$$

Let us now turn to the case of an isothermal but compressible fluid. Its density is now given by

$$\rho(p) = \rho_0(1 - \beta p)$$

where  $\beta$  is the compressibility (assumed to be constant in the domain). We must then solve

$$\begin{aligned} \frac{dp}{dz} = -\rho_0(1 - \beta p)g &\Rightarrow \frac{dp}{1 - \beta p} = -\rho_0 g dz \\ &\Rightarrow \int \frac{dp}{1 - \beta p} = - \int \rho_0 g dz \\ &\Rightarrow -\frac{1}{\beta} \ln(1 - \beta p) = -\rho_0 g z + C \\ &\Rightarrow \ln(1 - \beta p) = \beta \rho_0 g z + D \\ &\Rightarrow 1 - \beta p = \exp(\beta \rho_0 g z + D) \\ &\Rightarrow p(z) = \frac{1}{\beta} [1 - \exp(\beta \rho_0 g z + D)] \end{aligned} \tag{18.16}$$

At  $z = L$  we require  $p = 0$  so we obtain

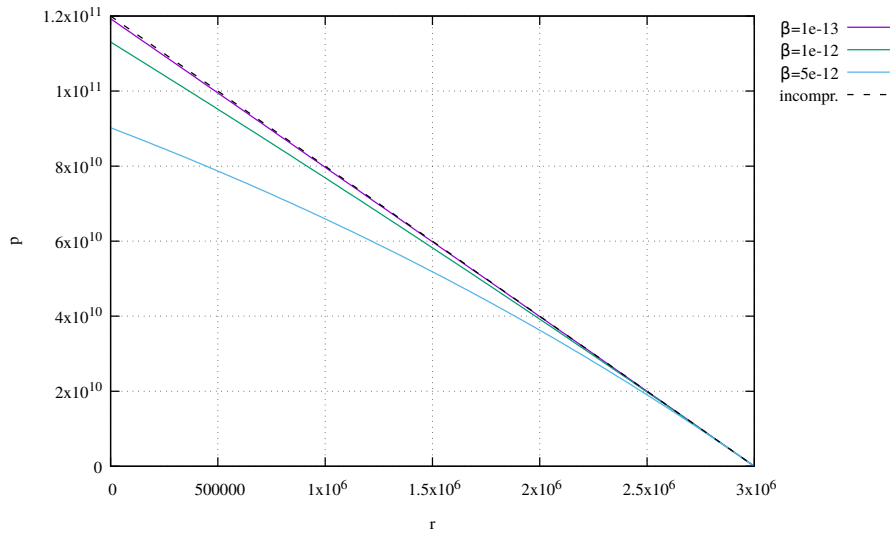
$$p(z) = \frac{1}{\beta} [1 - \exp(\beta \rho_0 g (z - L))]$$

Note that when the compressibility tends to zero, by virtue of

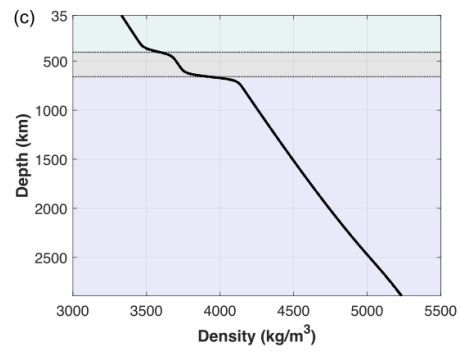
$$\exp x \sim 1 + x + \frac{x^2}{2} + \dots$$

for  $x \rightarrow 0$  we then recover the linear pressure profile above.

Let us now take  $\rho_0 = 4000 \text{ kg m}^{-3}$ ,  $g = 10 \text{ m s}^{-2}$  and  $\beta = 4 \cdot 10^{-12} \text{ Pa}^{-1}$  [439] and  $L = 3000 \text{ km}$ .



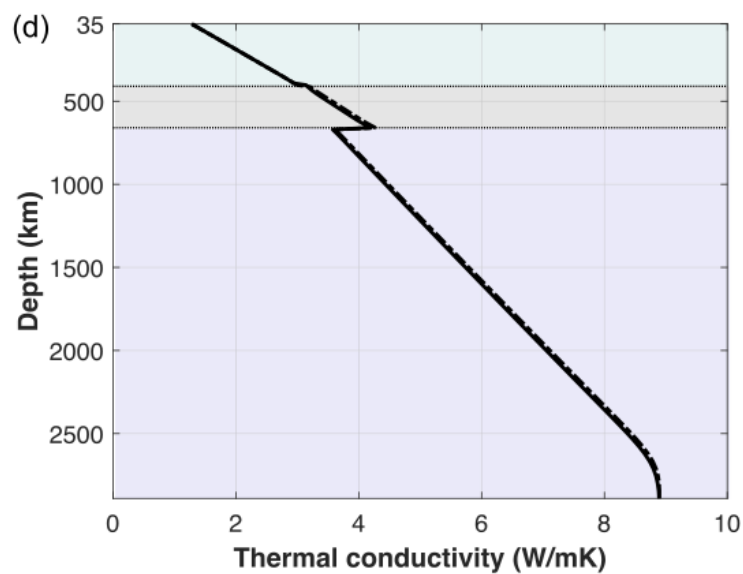
TODO: produce same plot with density



Taken from Neuharth and Mittelstaedt [934].

# 18.7 Earth radial thermal conductivity profile

thermal\_conductivity\_profile.tex



Taken from Neuharth and Mittelstaedt [934].

# Appendix A

## Matrix properties

### A.0.1 Symmetric matrices

Any *symmetric* matrix has only real eigenvalues, is always diagonalizable, and has orthogonal eigenvectors. A symmetric  $N \times N$  real matrix  $\mathbf{M}$  is said to be

- **positive definite** if  $\vec{x} \cdot \mathbf{M} \cdot \vec{x} > 0$  for every non-zero vector  $\vec{x}$  of  $n$  real numbers. All the eigenvalues of a Symmetric Positive Definite (SPD) matrix are positive. If  $A$  and  $B$  are positive definite, then so is  $A+B$ . The matrix inverse of a positive definite matrix is also positive definite. An SPD matrix has a unique Cholesky decomposition. In other words the matrix  $\mathbf{M}$  is positive definite if and only if there exists a unique lower triangular matrix  $\mathbf{L}$ , with real and strictly positive diagonal elements, such that  $\mathbf{M} = \mathbf{L}\mathbf{L}^T$  (the product of a lower triangular matrix and its conjugate transpose). This factorization is called the Cholesky decomposition of  $\mathbf{M}$ .
- **positive semi-definite** if  $\vec{x} \cdot \mathbf{M} \cdot \vec{x} \geq 0$
- **negative definite** if  $\vec{x} \cdot \mathbf{M} \cdot \vec{x} < 0$
- **negative semi-definite** if  $\vec{x} \cdot \mathbf{M} \cdot \vec{x} \leq 0$

The Stokes linear system

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathbf{v} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \mathbf{g} \end{pmatrix}$$

is **indefinite** (i.e. it has positive as well as negative eigenvalues).

A square matrix that is not invertible is called **singular** or degenerate. A square matrix is singular if and only if its determinant is 0. Singular matrices are rare in the sense that if you pick a random square matrix, it will almost surely not be singular.

### A.0.2 Eigenvalues of positive definite matrix

Suppose our matrix  $\mathbf{M}$  has eigenvalue  $\lambda$ .

If  $\lambda = 0$ , then there is some eigenvector  $\vec{x}$  so that  $\mathbf{M} \cdot \vec{x} = \lambda \vec{x} = \vec{0}$ . But then  $\vec{x}^T \cdot \mathbf{M} \cdot \vec{x} = 0$ , and so  $\mathbf{M}$  is not positive definite.

If  $\lambda < 0$ , then there is some eigenvector  $\vec{x}$  so that  $\mathbf{M} \cdot \vec{x} = \lambda \vec{x}$ . But then  $\vec{x}^T \cdot \mathbf{M} \cdot \vec{x} = \lambda |\vec{x}|^2$ , which is negative since  $|\vec{x}|^2 > 0$  and  $\lambda < 0$ . Thus  $\mathbf{M}$  is not positive definite.

And so if  $\mathbf{M}$  is positive definite, it only has positive eigenvalues.

### A.0.3 Schur complement

From wiki. In linear algebra and the theory of matrices, the Schur complement of a matrix block (i.e., a submatrix within a larger matrix) is defined as follows. Suppose  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ ,  $\mathbf{D}$  are respectively  $p \times p$ ,  $p \times q$ ,  $q \times p$  and  $q \times q$  matrices, and  $\mathbb{D}$  is invertible. Let

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}$$

so that  $\mathbf{M}$  is a  $(p+q) \times (p+q)$  matrix. Then the Schur complement of the block  $\mathbf{D}$  of the matrix  $\mathbf{M}$  is the  $p \times p$  matrix

$$\mathbf{S} = \mathbf{A} - \mathbf{B} \cdot \mathbf{D}^{-1} \cdot \mathbf{C}$$

Application to solving linear equations: The Schur complement arises naturally in solving a system of linear equations such as

$$\begin{aligned} \mathbf{A} \cdot \vec{x} + \mathbf{B} \cdot \vec{y} &= \vec{f} \\ \mathbf{C} \cdot \vec{x} + \mathbf{D} \cdot \vec{y} &= \vec{g} \end{aligned}$$

where  $\vec{x}$ ,  $\vec{f}$  are  $p$ -dimensional vectors,  $\vec{y}$ ,  $\vec{g}$  are  $q$ -dimensional vectors, and  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ ,  $\mathbf{D}$  are as above. Multiplying the bottom equation by  $\mathbf{B} \cdot \mathbf{D}^{-1}$  and then subtracting from the top equation one obtains

$$(\mathbf{A} - \mathbf{B} \cdot \mathbf{D}^{-1} \cdot \mathbf{C}) \cdot \vec{x} = \vec{f} - \mathbf{B} \cdot \mathbf{D}^{-1} \cdot \vec{g}$$

Thus if one can invert  $\mathbf{D}$  as well as the Schur complement of  $\mathbf{D}$ , one can solve for  $\vec{x}$ , and then by using the equation  $\mathbf{C} \cdot \vec{x} + \mathbf{D} \cdot \vec{y} = \vec{g}$  one can solve for  $y$ . This reduces the problem of inverting a  $(p+q) \times (p+q)$  matrix to that of inverting a  $p \times p$  matrix and a  $q \times q$  matrix. In practice one needs  $\mathbf{D}$  to be well-conditioned in order for this algorithm to be numerically accurate.

Considering now the Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{v} \\ \vec{p} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{g} \end{pmatrix}$$

Factorising for  $\vec{p}$  we end up with a **velocity-Schur complement**. Solving for  $\vec{p}$  in the second equation and inserting the expression for  $\vec{p}$  into the first equation we have

$$\mathbb{S}_v \cdot \vec{v} = \vec{f} \quad \text{with} \quad \mathbb{S}_v = \mathbb{K} + \mathbb{G} \cdot \mathbb{C}^{-1} \cdot \mathbb{G}^T$$

Factorising for  $\vec{v}$  we get a **pressure-Schur complement**.

$$\mathbb{S}_p \cdot \vec{p} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} \quad \text{with} \quad \mathbb{S}_p = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}$$

# Appendix B

## Don't be a hero - unless you have to

What follows was published online on July 17th, 2017 at <https://blogs.egu.eu/divisions/gd/2017/07/19/dont-be-a-hero-unless-you-have-to/> It was written by me and edited by Iris van Zelst, at the time PhD student at ETH Zürich.

In December 2013, I was invited to give a talk about the ASPECT code [1] at the American Geological Union conference in San Francisco. Right after my talk, Prof. Louis Moresi took the stage and gave a talk entitled: *Underworld: What we set out to do, How far did we get, What did we Learn?*

The abstract went as follows:

“Underworld was conceived as a tool for modelling 3D lithospheric deformation coupled with the underlying / surrounding mantle flow. The challenges involved were to find a method capable of representing the complicated, non-linear, history dependent rheology of the near surface as well as being able to model mantle convection, and, simultaneously, to be able to solve the numerical system efficiently. [...] The elegance of the method is that it can be completely described in a couple of sentences. However, there are some limitations: it is not obvious how to retain this elegance for unstructured or adaptive meshes, arbitrary element types are not sufficiently well integrated by the simple quadrature approach, and swarms of particles representing volumes are usually an inefficient representation of surfaces.”

Aside from the standard numerical modelling jargon, Louis used a term during his talk which I thought at the time had a nice ring to it: hero codes. In short, I believe he meant the codes written essentially by one or two people who at some point in time spent great effort into writing a code (usually choosing a range of applications, a geometry, a number of dimensions, a particular numerical method to solve the relevant PDEs(1), and a tracking method for the various fields of interest).

In the long list of Hero codes, one could cite (in alphabetical order) CITCOMS [1], DOUAR [8], FANTOM [2], IELVIS [5], LaMEM [3], pTatin [4], SLIM3D [10], SOPALE [7], StaggyYY [6], SULEC [11], Underworld [9], and I apologise to all other heroes out there whom I may have overlooked. And who does not want to be a hero? The Spiderman of geodynamics, the Superwoman of modelling?

Louis' talk echoed my thoughts on two key choices we (computational geodynamicists) are facing: Hero or not, and if yes, what type?

### Hero or not?

Speaking from experience, it is an intense source of satisfaction when peer-reviewed published results are obtained with the very code one has painstakingly put together over months, if not years. But is it worth it?

On the one hand, writing one own's code is a source of deep learning, a way to ensure that one understands the tool and knows its limitations, and a way to ensure that the code has the appropriate combination of features which are necessary to answer the research question at hand. On the other hand, it is akin to a journey; a rather long term commitment; a sometimes frustrating endeavour,



with no guarantee of success. Let us not deny it - many a student has started with one code only to switch to plan B sooner or later. Ultimately, this yields a satisfactory tool with often little to no perennial survival over the 5 year mark, a scarce if at all existent documentation, and almost always not compliant with the growing trend of long term repeatability. Furthermore, the resulting code will probably bear the marks of its not-all-knowing creator in its DNA and is likely not to be optimal nor efficient by modern computational standards.

This brings me to the second choice: elegance & modularity or tailored code & raw performance? Should one develop a code in a very broad framework using as much external libraries as possible or is there still space for true heroism?

It is my opinion that the answer to this question is: both. The current form of heroism no more lies in writing one's own FEM(2)/FDM(3) packages, meshers, or solvers from scratch, but in cleverly taking advantage of state-of-the-art packages such as for example p4est [15] for Adaptive Mesh Refinement, PetSc [13] or Trilinos [14] for solvers, Saint Germain [17] for particle tracking, deal.ii [12] or Fenics [16] for FEM, and sharing their codes through platforms such as svn, bitbucket or github.

In reality, the many different ways of approaching the development or usage of a (new) code is linked to the diversity of individual projects, but ultimately anyone who dares to touch a code (let alone write one) is a hero in his/her own right: although (super-)heroes can be awesome on their own, they often complete each other, team up and join forces for maximum efficiency. Let us all be heroes, then, and join efforts to serve Science to the best of our abilities.

#### Abbreviations

- (1) PDE: Partial Differential Equation
- (2) FEM: Finite Element Method
- (3) FDM: Finite Difference Method

#### References

- [1] Zhong *et al.* , JGR 105, 2000;
- [2] Thieulot, PEPI 188, 2011;
- [3] Kaus *et al.* , NIC Symposium proceedings, 2016;
- [4] May *et al.* , CMAME 290, 2015
- [5] Gerya and Yuen, PEPI 163, 2007
- [6] Tackley, PEPI 171, 2008
- [7] Fullsack, GJI 120, 1995
- [8] Braun *et al.* , PEPI 171, 2008
- [9] <http://www.underworldcode.org/>
- [10] Popov and Sobolev, PEPI 171, 2008
- [11] <http://www.geodynamics.no/buiter/sulec.html>
- [12] Bangerth *et al.* , J. Numer. Math., 2016; <http://www.dealii.org/>
- [13] <http://www.mcs.anl.gov/petsc/>
- [14] <https://trilinos.org/>
- [15] Burstedde *et al.* , SIAM journal on Scientific Computing, 2011; <http://www.p4est.org/>
- [16] <https://fenicsproject.org/>
- [17] Quenette *et al.* , Proceedings 19th IEEE, 2007

# Appendix C

## Some useful Python commands

app\_useful\_python.tex

### C.0.1 Sparse matrices

So far, the best way I have found to deal with sparse matrices is to declare the matrix as a 'lil\_matrix' (linked list).

```
from scipy.sparse import csr_matrix, lil_matrix
A_mat = lil_matrix((Nfem,Nfem),dtype=np.float64)
```

One then adds terms to it as if it was a full/dense matrix. Once the assembly is done, the conversion to CSR format is trivial:

```
A_mat=A_mat.tocsr()
```

Finally the solver can be called:

```
sol=sps.linalg.spsolve(A_mat,rhs)
```

### C.0.2 condition number

if the matrix has been declared as lil\_matrix, first convert it to a dense matrix:

```
A_mat=A_mat.dense()
```

The condition number of the matrix is simply obtained as follows:

```
from numpy import linalg as LA
print(LA.cond(A_mat))
```

### C.0.3 Weird stuff

Python is touted as the one language students should learn and master. However it is a language which allows \*way\* too much liberty in its syntax and encourages students to be sloppy.

For instance the following code runs just fine:

```
for k in range(0,5):
 for k in range(0,5):
 for k in range(0,5):
 print(k)
```

This alone should disqualify this language. It is easy to see the obvious problem with this code, but adding a few lines of code in between each 'for' line hides the problem and the absence of any warning makes this code a nightmare to debug.

<https://www.w3schools.com/python/default.asp>

<https://www.codecademy.com/learn/learn-python>

<https://learnpythonthehardway.org/book/>

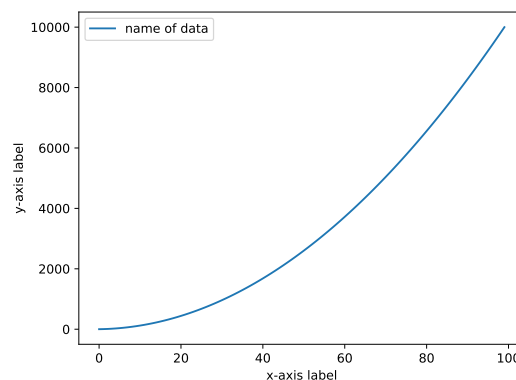
## C.0.4 Making simple 2D plots

```
import matplotlib.pyplot as plt

number of points
N=100

a despicable way of filling two arrays
x_data=[]
y_data=[]
for i in range(0,N):
 x=i
 y=i**2+2*i+1.
 x_data.append(x)
 y_data.append(y)

generating a 2D figure with the data
plt.figure()
plt.plot(x_data,y_data, label = 'name_of_data')
plt.xlabel('x-axis_label')
plt.ylabel('y-axis_label')
plt.legend()
plt.savefig('myplot.pdf', bbox_inches='tight')
plt.show()
```



## C.0.5 Making simple 3D plots of scatter

```
fig = plt.figure()
ax = plt.axes(projection='3d')
ax.set_title("insert_here_text_for_title")
size = ..some value..
ax.scatter3D(x, y, z, s = size)
```

## C.0.6 How to debug your code

Debugging a FE code is by no means trivial. There is (at least) a grid, a connectivity array, basis functions and their derivatives, the elemental matrices and rhs, the assembly, boundary conditions, and a call to a solver before the solution (if the solver returns one!) can be visualised.

- First and foremost, make sure that your grid of points is correct. For instance, you can resort to exporting it to an ascii file as follows:

```
np.savetxt('velocity.ascii',np.array([x,y,u,v]).T,header='#x,y,u,v')
```

In two dimensions, you should set for example  $n_{elx}=3$  and  $n_{ely}=2$ , so that for  $Q_1$  elements the grid counts 12 points. Then make sure the coordinates and the order of the points makes sense. Repeat the process for pressure nodes, temperature nodes, etc ...

- Then it is time to check the connectivity array(s).

```
for iel in range(0,nel):
 print("iel=",iel)
 for k in range(0,m)
 print("node_",icon[0,iel],"at_pos.",x[icon[0,iel]],y[icon[0,iel]])
```

This displays the list of nodes and their positions making each element. Repeat the process for every connectivity array.

- We can go on with testing that the all basis functions are 1 on their node and zero elsewhere:

```
for i in range(0,m):
 print('node',i,':',NNV(rnodes[i],snodes[i]))
```

here the arrays rnodes and snodes contain the (r,s) coordinates of the m nodes

- test jacobian, compute volume of domain
- $\text{sum}(dNNNdx)=0$
- print nodes where bc

## C.0.7 Optional arguments

Courtesy of Henry Brett.

```
def myfunc(a,b, *optional_arguments, **keyword_arguments):
 print(a)
 print(b)
 for ar in optional_arguments:
 print(ar)
 d=keyword_arguments.get("d", None)
 print(d)
```

```
a="dog"
b="cat"
myfunc(a,b,"shrek","fiona",d="donkey")
```

## C.0.8 drawing and filling quadrilaterals

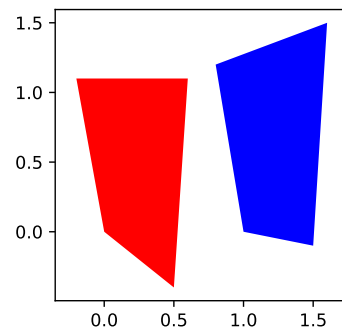
```
import matplotlib.pyplot as plt

plt.figure(figsize=(3, 3))
plt.axis('equal')

x=(1,1.5,1.6,0.8)
y=(0,-0.1,1.5,1.2)
plt.fill(x, y,"b")

x=(0,0.5,0.6,-0.2)
y=(0,-0.4,1.1,1.1)
plt.fill(x, y,"r")

plt.savefig('xxx.pdf')
plt.show()
```



# Appendix D

## Some useful maths

maths.tex

### D.0.1 Inverse of a 3x3 matrix

Let us assume we wish to solve the system  $\mathbf{A} \cdot \vec{X} = \vec{b}$ , with  $\vec{X} = (x, y)$ . Then the solution is given by The solution is given by

$$x = \frac{1}{\det(\mathbf{A})} \begin{vmatrix} b_1 & a_{21} \\ b_2 & a_{22} \end{vmatrix} \quad y = \frac{1}{\det(\mathbf{A})} \begin{vmatrix} a_{11} & b_1 \\ a_{21} & b_2 \end{vmatrix}$$

### D.0.2 Inverse of a 3x3 matrix

Let us consider the  $3 \times 3$  matrix  $\mathbf{M}$

$$\mathbf{M} = \begin{pmatrix} M_{xx} & M_{xy} & M_{xz} \\ M_{yx} & M_{yy} & M_{yz} \\ M_{zx} & M_{zy} & M_{zz} \end{pmatrix}$$

1. Find  $\det(\mathbf{M})$ , the determinant of the Matrix  $\mathbf{M}$ . The determinant will usually show up in the denominator of the inverse. If the determinant is zero, the matrix won't have an inverse.
2. Find  $\mathbf{M}^T$ , the transpose of the matrix. Transposing means reflecting the matrix about the main diagonal.

$$\mathbf{M}^T = \begin{pmatrix} M_{xx} & M_{yx} & M_{zx} \\ M_{xy} & M_{yy} & M_{zy} \\ M_{xz} & M_{yz} & M_{zz} \end{pmatrix}$$

3. Find the determinant of each of the  $2 \times 2$  minor matrices. For instance  $\tilde{M}_{xx} = M_{yy}M_{zz} - M_{yz}M_{zy}$ , or  $\tilde{M}_{xz} = M_{xy}M_{yz} - M_{xz}M_{yy}$ .
4. assemble the  $\tilde{\mathbf{M}}$  matrix:

$$\tilde{\mathbf{M}} = \begin{pmatrix} +\tilde{M}_{xx} & -\tilde{M}_{xy} & +\tilde{M}_{xz} \\ -\tilde{M}_{yx} & +\tilde{M}_{yy} & -\tilde{M}_{yz} \\ +\tilde{M}_{zx} & -\tilde{M}_{zy} & +\tilde{M}_{zz} \end{pmatrix}$$

5. the inverse of  $\mathbf{M}$  is then given by

$$\mathbf{M}^{-1} = \frac{1}{\det(\mathbf{M})} \tilde{\mathbf{M}}$$

Another approach which of course is equivalent to the above is Cramer's rule. Let us assume we wish to solve the system  $\mathbf{A} \cdot \vec{X} = \vec{b}$ , with  $\vec{X} = (x, y, z)$ . Then the solution is given by

$$x = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} b_1 & a_{12} & a_{13} \\ b_2 & a_{22} & a_{23} \\ b_3 & a_{32} & a_{33} \end{vmatrix} \quad y = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} a_{11} & b_1 & a_{13} \\ a_{21} & b_2 & a_{23} \\ a_{31} & b_3 & a_{33} \end{vmatrix} \quad z = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} a_{11} & a_{12} & b_1 \\ a_{21} & a_{22} & b_2 \\ a_{31} & a_{32} & b_3 \end{vmatrix}$$

# Appendix E

## Elemental matrices for simple geometries

app\_elemental\_matrix.tex

In what follows I compute the mass matrix for a variety of reference elements. If you wish to use these in a code, do not forget to take the jacobian of the transformation/mapping into account.

### E.0.1 1D segments

#### Linear basis functions

Let us start with the mass matrix (which we encountered in Section 6.1 – although we leave the  $\rho C_p$  term out):

$$\mathbf{M}_e = \int_{\Omega_e} \vec{N}^T \vec{N} dV = \int_{-1}^{+1} \vec{N}^T \vec{N} dr \quad (\text{E.1})$$

on the reference element, with

$$\vec{N}^T = \begin{pmatrix} N_1(r) & N_2(r) \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1-r & 1+r \end{pmatrix}$$

We have

$$\int_{-1}^{+1} N_1(r) N_1(r) dr = 2/3 \quad (\text{E.2})$$

$$\int_{-1}^{+1} N_1(r) N_2(r) dr = 1/3 \quad (\text{E.3})$$

$$\int_{-1}^{+1} N_2(r) N_2(r) dr = 2/3 \quad (\text{E.4})$$

Following the procedure in Section 6.1 we arrive at

$$\mathbf{M}^e = \frac{1}{3} \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

The lumped mass matrix is then

$$\bar{\mathbf{M}}^e = \frac{1}{3} \begin{pmatrix} 2+1 & 0 \\ 0 & 1+2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (\text{E.5})$$



**Remark.** The sum of all the terms in the mass matrix must be equal to 2. Indeed:

$$\begin{aligned}
 \sum_{ij} M_{ij} &= \sum_{ij} \int_{-1}^{+1} N_i N_j dr \\
 &= \int_{-1}^{+1} (N_1 N_1 + N_1 N_2 + N_2 N_1 + N_2 N_2) dr \\
 &= \int_{-1}^{+1} [N_1(N_1 + N_2) + N_2(N_1 + N_2)] dr \\
 &= \int_{-1}^{+1} (N_1 + N_2) dr \\
 &= 2
 \end{aligned}$$

### Quadratic basis functions

There are now three nodes in the segment so that the mass matrix is now a  $3 \times 3$  matrix. We have (see Section 5.2.1)

$$\vec{N}^T(r) = \begin{pmatrix} N_1(r) \\ N_2(r) \\ N_3(r) \end{pmatrix} = \begin{pmatrix} \frac{1}{2}r(r-1) \\ 1-r^2 \\ \frac{1}{2}r(r+1) \end{pmatrix} \quad (\text{E.6})$$

We then have to compute

$$\begin{aligned}
 \int_{-1}^{+1} N_1(r) N_1(r) dr &= \frac{8}{30} = 0.26666 \\
 \int_{-1}^{+1} N_1(r) N_2(r) dr &= \frac{4}{30} = 0.13333 \\
 \int_{-1}^{+1} N_1(r) N_3(r) dr &= -\frac{2}{30} = -0.06666... \\
 \int_{-1}^{+1} N_2(r) N_2(r) dr &= \frac{16}{15} = 1.06666 \\
 \int_{-1}^{+1} N_2(r) N_3(r) dr &= \frac{4}{30} = 0.13333 \\
 \int_{-1}^{+1} N_3(r) N_3(r) dr &= \frac{8}{30} = 0.26666
 \end{aligned}$$

and finally

$$\mathbf{M}^e = \frac{1}{30} \begin{pmatrix} 8 & 4 & -2 \\ 4 & 32 & 4 \\ -2 & 4 & 8 \end{pmatrix} \quad (\text{E.7})$$

The lumped mass matrix is then

$$\begin{aligned}
 \bar{\mathbf{M}}^e &= \frac{1}{30} \begin{pmatrix} 8+4-2 & 0 & 0 \\ 0 & 4+32+4 & 0 \\ 0 & 0 & -2+4+8 \end{pmatrix} \\
 &= \frac{1}{30} \begin{pmatrix} 10 & 0 & 0 \\ 0 & 40 & 0 \\ 0 & 0 & 10 \end{pmatrix} \\
 &= \frac{1}{3} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (\text{E.8})
 \end{aligned}$$

We can easily verify that

$$\sum_{ij} M_{ij} = 2 \quad \sum_{ij} \bar{M}_{ij} = 2$$

### Cubic basis functions

There are now four nodes in the segment so that the mass matrix is now a  $4 \times 4$  matrix. We have (see Section 5.2.3)

$$\vec{N}^T(r) = \begin{pmatrix} N_1(r) \\ N_2(r) \\ N_3(r) \\ N_4(r) \end{pmatrix} = \frac{1}{16} \begin{pmatrix} -1 + r + 9r^2 - 9r^3 \\ 9 - 27r - 9r^2 + 27r^3 \\ 9 + 27r - 9r^2 - 27r^3 \\ -1 - r + 9r^2 + 9r^3 \end{pmatrix} \quad (\text{E.9})$$

$$\begin{aligned} \int_{-1}^{+1} N_1(r) N_1(r) dr &= \frac{1}{256} \frac{4096}{105} \\ \int_{-1}^{+1} N_1(r) N_2(r) dr &= \frac{1}{256} \frac{1056}{35} \\ \int_{-1}^{+1} N_1(r) N_3(r) dr &= -\frac{1}{256} \frac{384}{35} \\ \int_{-1}^{+1} N_1(r) N_4(r) dr &= \frac{1}{256} \frac{608}{105} \\ \int_{-1}^{+1} N_2(r) N_2(r) dr &= \frac{1}{256} \frac{6912}{35} \\ \int_{-1}^{+1} N_2(r) N_3(r) dr &= -\frac{1}{256} \frac{864}{35} \\ \int_{-1}^{+1} N_2(r) N_4(r) dr &= -\frac{1}{256} \frac{384}{35} \\ \int_{-1}^{+1} N_3(r) N_3(r) dr &= \frac{1}{256} \frac{6912}{35} \\ \int_{-1}^{+1} N_3(r) N_4(r) dr &= \frac{1}{256} \frac{1056}{35} \\ \int_{-1}^{+1} N_4(r) N_4(r) dr &= \frac{1}{256} \frac{4096}{105} \end{aligned}$$

and finally

$$\mathbf{M}^e = \frac{1}{16} \frac{1}{105} \begin{pmatrix} 256 & 198 & -72 & 38 \\ 198 & 1296 & -162 & -72 \\ -72 & -162 & 1296 & 198 \\ 38 & -72 & 198 & 256 \end{pmatrix} \quad (\text{E.10})$$

The lumped mass matrix is then

$$\begin{aligned}
\bar{\mathbf{M}}^e &= \frac{1}{16} \frac{1}{105} \begin{pmatrix} 256 + 198 - 72 + 38 & 0 & 0 & 0 \\ 0 & 198 + 1296 - 162 - 72 & 0 & 0 \\ 0 & 0 & -72 - 162 + 1296 + 198 & 0 \\ 0 & 0 & 0 & 38 - 72 + 198 + \end{pmatrix} \\
&= \frac{1}{16} \frac{1}{105} \begin{pmatrix} 420 & 0 & 0 & 0 \\ 0 & 1260 & 0 & 0 \\ 0 & 0 & 1260 & 0 \\ 0 & 0 & 0 & 420 \end{pmatrix} \\
&= \frac{1}{4} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}
\end{aligned}$$

We can easily verify that

$$\sum_{ij} M_{ij} = 2 \qquad \sum_{ij} \bar{M}_{ij} = 2$$

### Quartic basis functions

There are now five nodes in the segment so that the mass matrix is now a  $5 \times 5$  matrix. We have (see Section 5.2.4)

$$\vec{N}^T(r) = \begin{pmatrix} N_1(r) \\ N_2(r) \\ N_3(r) \\ N_4(r) \\ N_5(r) \end{pmatrix} = \frac{1}{6} \begin{pmatrix} r - r^2 - 4r^3 + 4r^4 \\ -8r + 16r^2 + 8r^3 - 16r^4 \\ 6 - 30r^2 + 24r^4 \\ 8r + 16r^2 - 8r^3 - 16r^4 \\ -r - r^2 + 4r^3 + 4r^4 \end{pmatrix} \quad (\text{E.11})$$

$$\begin{aligned}
\int_{-1}^{+1} N_1(r)N_1(r)dr &= \frac{1}{36} \frac{1168}{315} \\
\int_{-1}^{+1} N_1(r)N_2(r)dr &= \frac{1}{36} \frac{1184}{315} \\
\int_{-1}^{+1} N_1(r)N_3(r)dr &= -\frac{1}{36} \frac{232}{105} \\
\int_{-1}^{+1} N_1(r)N_4(r)dr &= \frac{1}{36} \frac{32}{45} \\
\int_{-1}^{+1} N_1(r)N_5(r)dr &= -\frac{1}{36} \frac{116}{315} \\
\int_{-1}^{+1} N_2(r)N_2(r)dr &= \frac{1}{36} \frac{1024}{45} \\
\int_{-1}^{+1} N_2(r)N_3(r)dr &= -\frac{1}{36} \frac{512}{105} \\
\int_{-1}^{+1} N_2(r)N_4(r)dr &= \frac{1}{36} \frac{1024}{315} \\
\int_{-1}^{+1} N_2(r)N_5(r)dr &= \frac{1}{36} \frac{32}{45} \\
\int_{-1}^{+1} N_3(r)N_3(r)dr &= \frac{1}{36} \frac{832}{35} \\
\int_{-1}^{+1} N_3(r)N_4(r)dr &= -\frac{1}{36} \frac{512}{105} \\
\int_{-1}^{+1} N_3(r)N_5(r)dr &= -\frac{1}{36} \frac{232}{105} \\
\int_{-1}^{+1} N_4(r)N_4(r)dr &= \frac{1}{36} \frac{1024}{45} \\
\int_{-1}^{+1} N_4(r)N_5(r)dr &= \frac{1}{36} \frac{1184}{315} \\
\int_{-1}^{+1} N_5(r)N_5(r)dr &= \frac{1}{36} \frac{1168}{315}
\end{aligned} \tag{E.12}$$

$$\mathbf{M}^e = \frac{1}{36} \frac{1}{315} \begin{pmatrix} 1168 & 1184 & -696 & 224 & -116 \\ 1184 & 7168 & -1536 & 1024 & 224 \\ -696 & -1536 & 7488 & -1536 & -696 \\ 224 & 1024 & -1536 & 7168 & 1184 \\ -116 & 224 & -696 & 1184 & 1168 \end{pmatrix} \tag{E.13}$$

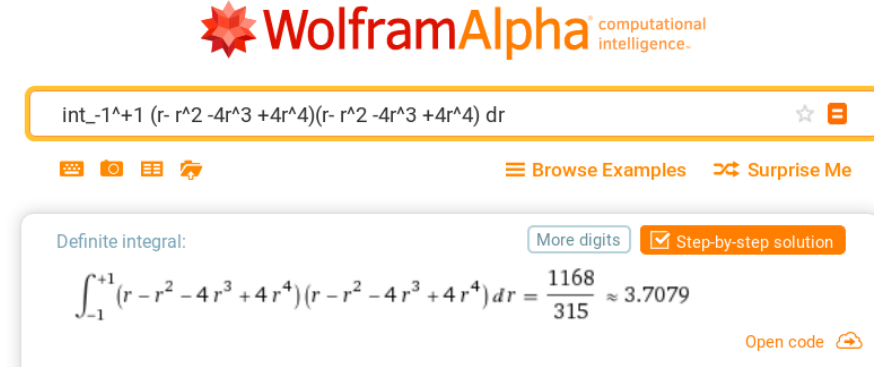
The lumped mass matrix is then

$$\bar{\mathbf{M}}^e = \frac{1}{36} \frac{1}{315} \begin{pmatrix} 1764 & 0 & 0 & 0 & 0 \\ 0 & 8064 & 0 & 0 & 0 \\ 0 & 0 & 3024 & 0 & 0 \\ 0 & 0 & 0 & 8064 & 0 \\ 0 & 0 & 0 & 0 & 1764 \end{pmatrix} = \frac{1}{45} \begin{pmatrix} 7 & 0 & 0 & 0 & 0 \\ 0 & 32 & 0 & 0 & 0 \\ 0 & 0 & 12 & 0 & 0 \\ 0 & 0 & 0 & 32 & 0 \\ 0 & 0 & 0 & 0 & 7 \end{pmatrix} \tag{E.14}$$

We can once again easily verify that

$$\sum_{ij} M_{ij} = 2 \quad \sum_{ij} \bar{M}_{ij} = 2$$

Note that all the integrals above were done very conveniently with the WolframAlpha software/website<sup>1</sup>. Example:



## E.0.2 Quadrilaterals: rectangular linear elements

### Mass matrix

We assume that each element is a rectangle of size  $h_x \times h_y$ . We start from the linear basis functions in the reference element as a function of  $r, s$ :

$$N_1(r, s) = \frac{1}{4}(1-r)(1-s) \quad (\text{E.15})$$

$$N_2(r, s) = \frac{1}{4}(1+r)(1-s) \quad (\text{E.16})$$

$$N_3(r, s) = \frac{1}{4}(1+r)(1+s) \quad (\text{E.17})$$

$$N_4(r, s) = \frac{1}{4}(1-r)(1+s) \quad (\text{E.18})$$

and their derivatives:

$$\partial_r N_1(r, s) = -\frac{1}{4}(1-s)$$

$$\partial_r N_2(r, s) = \frac{1}{4}(1-s)$$

$$\partial_r N_3(r, s) = \frac{1}{4}(1+s)$$

$$\partial_r N_4(r, s) = -\frac{1}{4}(1+s)$$

$$\partial_s N_1(r, s) = -\frac{1}{4}(1-r)$$

$$\partial_s N_2(r, s) = -\frac{1}{4}(1+r)$$

$$\partial_s N_3(r, s) = \frac{1}{4}(1+r)$$

$$\partial_s N_4(r, s) = \frac{1}{4}(1-r)$$

<sup>1</sup><https://www.wolframalpha.com/>

We wish to compute the integral of a function  $f(x, y)$  over the rectangular element:

$$\iint f(x, y) dx dy = \iint f(x(r, s), y(r, s)) \left| \frac{\partial(x, y)}{\partial(r, s)} \right| dr ds \quad (\text{E.19})$$

$$= \iint f(x(r, s), y(r, s)) \begin{vmatrix} \partial x / \partial r & \partial x / \partial s \\ \partial y / \partial r & \partial y / \partial s \end{vmatrix} dr ds \quad (\text{E.20})$$

From

$$x(r, s) = \sum_{i=1}^4 N_i(r, s) x_i \quad \text{and} \quad y(r, s) = \sum_{i=1}^4 N_i(r, s) y_i$$

we can write

$$\begin{aligned} \frac{\partial x}{\partial r}(r, s) &= \frac{\partial N_1}{\partial r} x_1 + \frac{\partial N_2}{\partial r} x_2 + \frac{\partial N_3}{\partial r} x_3 + \frac{\partial N_4}{\partial r} x_4 \\ &= -\frac{1}{4}(1-s)x_1 + \frac{1}{4}(1-s)x_2 + \frac{1}{4}(1+s)x_3 - \frac{1}{4}(1+s)x_4 \\ &= \frac{1}{4}(-x_1 + x_2 + x_3 - x_4 + s(x_1 - x_2 + x_3 - x_4)) \\ &= \frac{1}{4}(h_x + h_x + s(x_1 - x_2 + x_2 - x_1)) \\ &= \frac{1}{2}h_x \\ \frac{\partial x}{\partial s}(r, s) &= \frac{\partial N_1}{\partial s} x_1 + \frac{\partial N_2}{\partial s} x_2 + \frac{\partial N_3}{\partial s} x_3 + \frac{\partial N_4}{\partial s} x_4 \\ &= -\frac{1}{4}(1-r)x_1 - \frac{1}{4}(1+r)x_2 + \frac{1}{4}(1+r)x_3 + \frac{1}{4}(1-r)x_4 \\ &= \frac{1}{4}(-x_1 - x_2 + x_3 + x_4 + r(x_1 - x_2 + x_3 - x_4)) \\ &= \frac{1}{4}(-x_1 - x_2 + x_2 + x_1 + r(x_1 - x_2 + x_2 - x_1)) \\ &= 0 \\ \frac{\partial y}{\partial r}(r, s) &= 0 \\ \frac{\partial y}{\partial s}(r, s) &= \frac{1}{2}h_y \end{aligned} \quad (\text{E.21})$$

Then

$$\begin{vmatrix} \partial x / \partial r & \partial x / \partial s \\ \partial y / \partial r & \partial y / \partial s \end{vmatrix} = \begin{vmatrix} h_x & 0 \\ 0 & h_y \end{vmatrix} = \frac{h_x h_y}{4}$$

and finally

$$\boxed{\iint_{\square} f(x, y) dx dy = \frac{h_x h_y}{4} \int_{-1}^1 \int_{-1}^1 f(x(r, s), y(r, s)) dr ds} \quad (\text{E.22})$$

Then the mass matrix is given by

$$\begin{aligned}
\mathbf{M}_e &= \frac{h_x h_y}{4} \int_{-1}^1 \int_{-1}^1 \begin{pmatrix} N_1(r, s)N_1(r, s) & N_1(r, s)N_2(r, s) & N_1(r, s)N_3(r, s) & N_1(r, s)N_4(r, s) \\ N_2(r, s)N_1(r, s) & N_2(r, s)N_2(r, s) & N_2(r, s)N_3(r, s) & N_2(r, s)N_4(r, s) \\ N_3(r, s)N_1(r, s) & N_3(r, s)N_2(r, s) & N_3(r, s)N_3(r, s) & N_3(r, s)N_4(r, s) \\ N_4(r, s)N_1(r, s) & N_4(r, s)N_2(r, s) & N_4(r, s)N_3(r, s) & N_4(r, s)N_4(r, s) \end{pmatrix} dr ds \\
&= \frac{h_x h_y}{9} \begin{pmatrix} 1 & 1/2 & 1/4 & 1/2 \\ 1/2 & 1 & 1/2 & 1/4 \\ 1/4 & 1/2 & 1 & 1/2 \\ 1/2 & 1/4 & 1/2 & 1 \end{pmatrix}
\end{aligned} \tag{E.23}$$

**Diffusion matrix**

$$\mathbf{K}_d^e = k \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \mathbf{B}^T(r, s) \cdot \mathbf{B}(r, s) dr ds$$

with

$$\mathbf{B}(r, s) = \begin{pmatrix} -\frac{1}{h_x} \frac{1}{2}(1-s) & \frac{1}{h_x} \frac{1}{2}(1-s) & \frac{1}{h_x} \frac{1}{2}(1+s) & -\frac{1}{h_x} \frac{1}{2}(1+s) \\ -\frac{1}{h_y} \frac{1}{2}(1-r) & -\frac{1}{h_y} \frac{1}{2}(1+r) & \frac{1}{h_y} \frac{1}{2}(1+r) & \frac{1}{h_y} \frac{1}{2}(1-r) \end{pmatrix}$$

Then

$$\mathbf{B}^T(r, s) \cdot \mathbf{B}(r, s) \tag{E.24}$$

$$= \begin{pmatrix} -\frac{1}{h_x} \frac{1}{2}(1-s) & -\frac{1}{h_y} \frac{1}{2}(1-r) \\ \frac{1}{h_x} \frac{1}{2}(1-s) & -\frac{1}{h_y} \frac{1}{2}(1+r) \\ \frac{1}{h_x} \frac{1}{2}(1+s) & \frac{1}{h_y} \frac{1}{2}(1+r) \\ -\frac{1}{h_x} \frac{1}{2}(1+s) & \frac{1}{h_y} \frac{1}{2}(1-r) \end{pmatrix} \cdot \begin{pmatrix} -\frac{1}{h_x} \frac{1}{2}(1-s) & \frac{1}{h_x} \frac{1}{2}(1-s) & \frac{1}{h_x} \frac{1}{2}(1+s) & -\frac{1}{h_x} \frac{1}{2}(1+s) \\ -\frac{1}{h_y} \frac{1}{2}(1-r) & -\frac{1}{h_y} \frac{1}{2}(1+r) & \frac{1}{h_y} \frac{1}{2}(1+r) & \frac{1}{h_y} \frac{1}{2}(1-r) \end{pmatrix} \tag{E.25}$$

$$= \frac{1}{4h_x^2} \begin{pmatrix} -(1-s) & -(1-r) \\ (1-s) & -(1+r) \\ (1+s) & (1+r) \\ -(1+s) & (1-r) \end{pmatrix} \cdot \begin{pmatrix} -(1-s) & (1-s) & (1+s) & -(1+s) \\ -(1-r) & -(1+r) & (1+r) & (1-r) \end{pmatrix} \tag{E.26}$$

$$= \frac{1}{4h_x^2} \begin{pmatrix} (1-r)^2 & (1-r^2) & -(1-r^2) & -(1-r)^2 \\ (1-r^2) & (1+r)^2 & -(1+r)^2 & -(1-r^2) \\ -(1-r^2) & -(1+r)^2 & (1+r)^2 & (1-r^2) \\ -(1-r)^2 & -(1-r^2) & (1-r^2) & (1-r)^2 \end{pmatrix} \tag{E.27}$$

$$= \frac{1}{4h_y^2} \begin{pmatrix} (1-s)^2 & -(1-s)^2 & -(1-s^2) & (1-s^2) \\ -(1-s)^2 & (1-s)^2 & (1-s^2) & -(1-s^2) \\ -(1-s^2) & (1-s^2) & (1+s)^2 & -(1+s)^2 \\ (1-s^2) & -(1-s^2) & -(1+s)^2 & (1+s)^2 \end{pmatrix} \tag{E.28}$$

$$\tag{E.29}$$

So in the end

$$\mathbf{K}_d^e = k \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4h_x^2} \begin{pmatrix} (1-r)^2 & (1-r^2) & -(1-r^2) & -(1-r)^2 \\ (1-r^2) & (1+r)^2 & -(1+r)^2 & -(1-r^2) \\ -(1-r^2) & -(1+r)^2 & (1+r)^2 & (1-r^2) \\ -(1-r)^2 & -(1-r^2) & (1-r^2) & (1-r)^2 \end{pmatrix} dr ds \quad (\text{E.30})$$

$$+ \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4h_y^2} \begin{pmatrix} (1-s)^2 & -(1-s)^2 & -(1-s^2) & (1-s^2) \\ -(1-s)^2 & (1-s)^2 & (1-s^2) & -(1-s^2) \\ -(1-s^2) & (1-s^2) & (1+s)^2 & -(1+s)^2 \\ (1-s^2) & -(1-s^2) & -(1+s)^2 & (1+s)^2 \end{pmatrix} dr ds \quad (\text{E.31})$$

$$= k \frac{kh_y}{8h_x} \int_{-1}^{+1} \begin{pmatrix} (1-r)^2 & (1-r^2) & -(1-r^2) & -(1-r)^2 \\ (1-r^2) & (1+r)^2 & -(1+r)^2 & -(1-r^2) \\ -(1-r^2) & -(1+r)^2 & (1+r)^2 & (1-r^2) \\ -(1-r)^2 & -(1-r^2) & (1-r^2) & (1-r)^2 \end{pmatrix} dr \quad (\text{E.32})$$

$$+ \frac{kh_x}{8h_y} \int_{-1}^{+1} \begin{pmatrix} (1-s)^2 & -(1-s)^2 & -(1-s^2) & (1-s^2) \\ -(1-s)^2 & (1-s)^2 & (1-s^2) & -(1-s^2) \\ -(1-s^2) & (1-s^2) & (1+s)^2 & -(1+s)^2 \\ (1-s^2) & -(1-s^2) & -(1+s)^2 & (1+s)^2 \end{pmatrix} ds \quad (\text{E.33})$$

$$= k \frac{kh_y}{8h_x} \frac{4}{3} \begin{pmatrix} 2 & 1 & -1 & 2 \\ 1 & 2 & -2 & -1 \\ -1 & -2 & 2 & 1 \\ -2 & -1 & 1 & 2 \end{pmatrix} + \frac{kh_x}{8h_y} \frac{4}{3} \begin{pmatrix} 2 & -2 & -1 & 1 \\ -2 & 2 & 1 & -1 \\ -1 & 1 & 2 & -2 \\ 1 & -1 & -2 & 2 \end{pmatrix} \quad (\text{E.34})$$

$$= \frac{kh_x h_y}{6} \begin{pmatrix} \frac{2}{h_x^2} + \frac{2}{h_y^2} & -\frac{2}{h_x^2} + \frac{1}{h_y^2} & -\frac{1}{h_x^2} - \frac{1}{h_y^2} & \frac{1}{h_x^2} - \frac{2}{h_y^2} \\ -\frac{2}{h_x^2} + \frac{1}{h_y^2} & \frac{2}{h_x^2} + \frac{2}{h_y^2} & \frac{1}{h_x^2} - \frac{2}{h_y^2} & -\frac{1}{h_x^2} - \frac{1}{h_y^2} \\ -\frac{1}{h_x^2} - \frac{1}{h_y^2} & \frac{1}{h_x^2} - \frac{2}{h_y^2} & \frac{2}{h_x^2} + \frac{2}{h_y^2} & -\frac{2}{h_x^2} + \frac{1}{h_y^2} \\ \frac{1}{h_x^2} - \frac{2}{h_y^2} & -\frac{1}{h_x^2} - \frac{1}{h_y^2} & -\frac{2}{h_x^2} + \frac{1}{h_y^2} & \frac{2}{h_x^2} + \frac{2}{h_y^2} \end{pmatrix} \quad (\text{E.35})$$

## Advection matrix

$$\mathbf{K}_a = \rho C_p \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \mathbf{N}^T(r, s) (\vec{\mathbf{v}} \cdot \mathbf{B}(r, s)) dr ds$$

with

$$\vec{\mathbf{v}} \cdot \mathbf{B}(r, s) = \left( -\frac{u}{2h_x}(1-s) - \frac{v}{2h_y}(1-r) \quad \frac{u}{2h_x}(1-s) - \frac{v}{2h_y}(1+r) \quad \frac{u}{2h_x}(1+s) + \frac{v}{2h_y}(1+r) \quad -\frac{u}{2h_x}(1+s) + \frac{v}{2h_y}(1-r) \right)$$



Assuming that the velocity is constant within the element (which is almost always not true!), we can write:

$$\mathbf{K}_a = \rho C_p \frac{h_x h_y}{16} \frac{v}{2h_y} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} (1-r)(1-s) \\ (1+r)(1-s) \\ (1+r)(1+s) \\ (1-r)(1+s) \end{pmatrix} \begin{pmatrix} -(1-r) & -(1+r) & (1+r) & (1-r) \end{pmatrix} dr ds \quad (\text{E.36})$$

$$+ \rho C_p \frac{h_x h_y}{16} \frac{u}{2h_x} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} (1-r)(1-s) \\ (1+r)(1-s) \\ (1+r)(1+s) \\ (1-r)(1+s) \end{pmatrix} \begin{pmatrix} -(1-s) & (1-s) & (1+s) & -(1+s) \end{pmatrix} dr ds \quad (\text{E.37})$$

$$= \rho C_p \frac{h_x v}{32} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} (1-r)(1-s) \\ (1+r)(1-s) \\ (1+r)(1+s) \\ (1-r)(1+s) \end{pmatrix} \begin{pmatrix} -(1-r) & -(1+r) & (1+r) & (1-r) \end{pmatrix} dr ds \quad (\text{E.38})$$

$$+ \rho C_p \frac{h_y u}{32} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} (1-r)(1-s) \\ (1+r)(1-s) \\ (1+r)(1+s) \\ (1-r)(1+s) \end{pmatrix} \begin{pmatrix} -(1-s) & (1-s) & (1+s) & -(1+s) \end{pmatrix} dr ds \quad (\text{E.39})$$

$$= \rho C_p \frac{1}{12} \left( v h_x \begin{pmatrix} -2 & -1 & 1 & 2 \\ -1 & -2 & 2 & 1 \\ -1 & -2 & 2 & 1 \\ -2 & -1 & 1 & 2 \end{pmatrix} + u h_y \begin{pmatrix} -2 & 2 & 1 & -1 \\ -2 & 2 & 1 & -1 \\ -1 & 1 & 2 & -2 \\ -1 & 1 & 2 & -2 \end{pmatrix} \right) \quad (\text{E.40})$$

and finally

$$\mathbf{K}_a = \frac{\rho C_p}{3} \begin{pmatrix} -\frac{1}{2}u h_y - \frac{1}{2}v h_x & \frac{1}{2}u h_y - \frac{1}{4}v h_x & \frac{1}{4}u h_y + \frac{1}{4}v h_x & -\frac{1}{4}u h_y + \frac{1}{2}v h_x \\ -\frac{1}{2}u h_y - \frac{1}{4}v h_x & \frac{1}{2}u h_y - \frac{1}{2}v h_x & \frac{1}{4}u h_y + \frac{1}{2}v h_x & -\frac{1}{4}u h_y + \frac{1}{4}v h_x \\ -\frac{1}{4}u h_y - \frac{1}{4}v h_x & \frac{1}{4}u h_y - \frac{1}{2}v h_x & \frac{1}{2}u h_y + \frac{1}{2}v h_x & -\frac{1}{2}u h_y + \frac{1}{4}v h_x \\ -\frac{1}{4}u h_y - \frac{1}{2}v h_x & \frac{1}{4}u h_y - \frac{1}{4}v h_x & \frac{1}{2}u h_y + \frac{1}{4}v h_x & -\frac{1}{2}u h_y + \frac{1}{2}v h_x \end{pmatrix}$$

## Matrices for D.G.

In the context of Discontinuous Galerkin methods we will need

$$\begin{aligned} \mathbf{J}_x &= \int_{\square} \partial_x \vec{N}^T(x, y) \vec{N}(x, y) dx dy \\ \mathbf{J}_y &= \int_{\square} \partial_y \vec{N}^T(x, y) \vec{N}(x, y) dx dy \end{aligned} \quad (\text{E.41})$$

We have

$$\partial_x N_i(x, y) = \frac{\partial N_i}{\partial r} \frac{\partial r}{\partial x} \quad \text{and} \quad \partial_y N_i(x, y) = \frac{\partial N_i}{\partial s} \frac{\partial s}{\partial y}$$

Since

$$r = \frac{2}{h_x}(x - x_0) - 1 \quad \text{and} \quad s = \frac{2}{h_y}(y - y_0) - 1$$

then

$$\frac{\partial r}{\partial x} = \frac{2}{h_x} \quad \frac{\partial s}{\partial y} = \frac{2}{h_y}$$

so

$$\begin{aligned}
\mathbf{J}_x &= \int_{\square} \partial_x \vec{N}^T(x, y) \vec{N}(x, y) dx dy \\
&= \frac{2}{h_x} \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \partial_r \vec{N}^T(r, s) \vec{N}(r, s) dr ds \\
&= \frac{h_y}{2} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} -\frac{1}{4}(1-s) \\ +\frac{1}{4}(1-s) \\ +\frac{1}{4}(1+s) \\ -\frac{1}{4}(1+s) \end{pmatrix} \begin{pmatrix} \frac{1}{4}(1-r)(1-s) & \frac{1}{4}(1+r)(1-s) & \frac{1}{4}(1+r)(1+s) & \frac{1}{4}(1-r)(1+s) \end{pmatrix} \\
&= \frac{h_y}{32} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} -(1-r)(1-s)^2 & -(1+r)(1-s)^2 & -(1+r)(1-s^2) & -(1-r)(1-s^2) \\ (1-r)(1-s)^2 & (1+r)(1-s)^2 & (1+r)(1-s^2) & (1-r)(1-s^2) \\ (1-r)(1-s^2) & (1+r)(1-s^2) & (1+r)(1+s)^2 & (1-r)(1+s)^2 \\ -(1-r)(1-s^2) & -(1+r)(1-s^2) & -(1+r)(1+s)^2 & -(1-r)(1+s)^2 \end{pmatrix} dr ds \\
&= \frac{h_y}{32} \begin{pmatrix} -16/3 & -16/3 & -8/3 & -8/3 \\ 16/3 & 16/3 & 8/3 & 8/3 \\ 8/3 & 8/3 & 16/3 & 16/3 \\ -8/3 & -8/3 & -16/3 & -16/3 \end{pmatrix} \\
&= \frac{h_y}{12} \begin{pmatrix} -2 & -2 & -1 & -1 \\ 2 & 2 & 1 & 1 \\ 1 & 1 & 2 & 2 \\ -1 & -1 & -2 & -2 \end{pmatrix} \\
\mathbf{J}_y &= \int_{\square} \partial_y \vec{N}^T(x, y) \vec{N}(x, y) dx dy \\
&= \frac{2}{h_y} \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} -\frac{1}{4}(1-r) \\ -\frac{1}{4}(1+r) \\ +\frac{1}{4}(1+r) \\ +\frac{1}{4}(1-r) \end{pmatrix} \begin{pmatrix} \frac{1}{4}(1-r)(1-s) & \frac{1}{4}(1+r)(1-s) & \frac{1}{4}(1+r)(1+s) & \frac{1}{4}(1-r)(1+s) \end{pmatrix} \\
&= \frac{h_x}{32} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} -(1-r)^2(1-s) & -(1-r^2)(1-s) & -(1-r^2)(1+s) & -(1-r)^2(1+s) \\ -(1-r^2)(1-s) & -(1+r)^2(1-s) & -(1+r)^2(1+s) & -(1-r^2)(1+s) \\ (1-r^2)(1-s) & (1+r)^2(1-s) & (1+r)^2(1+s) & (1-r^2)(1+s) \\ (1-r)^2(1-s) & (1-r^2)(1-s) & (1-r^2)(1+s) & (1-r)^2(1+s) \end{pmatrix} dr ds \\
&= \frac{h_x}{32} \begin{pmatrix} -16/3 & -8/3 & -8/3 & -16/3 \\ -8/3 & -16/3 & -16/3 & -8/3 \\ 8/3 & 16/3 & 16/3 & 8/3 \\ 16/3 & 8/3 & 8/3 & 16/3 \end{pmatrix} \\
&= \frac{h_x}{12} \begin{pmatrix} -2 & -1 & -1 & -2 \\ -1 & -2 & -2 & -1 \\ 1 & 2 & 2 & 1 \\ 2 & 1 & 1 & 2 \end{pmatrix}
\end{aligned}$$

## Computing matrix $C_1$

$$C_1 = \int_{\partial\Omega_1} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma$$

The edge  $\partial\Omega_1$  is bounded by the coordinates of nodes  $x_1, y_1$  and  $x_2, y_2$ . This segment can be parameterised by  $t \in [0, 1]$ :

$$\vec{r}(t) = (1-t) \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} + t \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} (x_2 - x_1)t + x_1 \\ (y_2 - y_1)t + y_1 \end{pmatrix}$$

Let us assume that  $C$  is a smooth curve and that it is given by the parametric equations  $x = h(t)$ ,  $y = g(t)$  and  $a \leq t \leq b$ . The line integral of a function  $f(x, y)$  over  $C$  is computed as follows.

$$\int_C f(x, y) ds = \int_a^b f(h(t), g(t)) \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt$$

In our case  $dx/dt = x_2 - x_1$  and  $dy/dt = y_2 - y_1$  so

$$\sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} = h_x$$

Then

$$\begin{aligned} C_1 &= \int_{\partial\Omega_1} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma \\ &= \begin{pmatrix} \int_{\partial\Omega_1} N_1(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_1(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_1(x, y) N_3(x, y) d\Gamma & \int_{\partial\Omega_1} N_1(x, y) N_4(x, y) d\Gamma \\ \int_{\partial\Omega_1} N_2(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_2(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_2(x, y) N_3(x, y) d\Gamma & \int_{\partial\Omega_1} N_2(x, y) N_4(x, y) d\Gamma \\ \int_{\partial\Omega_1} N_3(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_3(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_3(x, y) N_3(x, y) d\Gamma & \int_{\partial\Omega_1} N_3(x, y) N_4(x, y) d\Gamma \\ \int_{\partial\Omega_1} N_4(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_4(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_4(x, y) N_3(x, y) d\Gamma & \int_{\partial\Omega_1} N_4(x, y) N_4(x, y) d\Gamma \end{pmatrix} \\ &= h_x \begin{pmatrix} \int_0^1 N_1(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_1(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_1(x(t), y(t)) N_3(x(t), y(t)) dt & \int_0^1 N_1(x(t), y(t)) N_4(x(t), y(t)) dt \\ \int_0^1 N_2(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_2(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_2(x(t), y(t)) N_3(x(t), y(t)) dt & \int_0^1 N_2(x(t), y(t)) N_4(x(t), y(t)) dt \\ \int_0^1 N_3(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_3(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_3(x(t), y(t)) N_3(x(t), y(t)) dt & \int_0^1 N_3(x(t), y(t)) N_4(x(t), y(t)) dt \\ \int_0^1 N_4(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_4(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_4(x(t), y(t)) N_3(x(t), y(t)) dt & \int_0^1 N_4(x(t), y(t)) N_4(x(t), y(t)) dt \end{pmatrix} \end{aligned}$$

On this edge  $y_1 = y_2$  so  $y(t) = y_1$ .

$$\begin{aligned}
N_1(x(t), y(t)) &= \frac{x_3 - x(t)}{x_3 - x_1} \frac{y_3 - y(t)}{y_3 - y_1} \\
&= \frac{x_3 - (x_2 - x_1)t - x_1}{x_3 - x_1} \frac{y_3 - y_1}{y_3 - y_1} \\
&= \frac{x_3 - x_1 - (x_3 - x_1)t}{x_3 - x_1} \frac{y_3 - y_1}{y_3 - y_1} \\
&= (1 - t)
\end{aligned} \tag{E.44}$$

$$\begin{aligned}
N_2(x(t), y(t)) &= \frac{x(t) - x_1}{x_2 - x_1} \frac{y_3 - y(t)}{y_3 - y_2} \\
&= \frac{(x_2 - x_1)t + x_1 - x_1}{x_2 - x_1} \frac{y_3 - y_1}{y_3 - y_2} \\
&= t
\end{aligned} \tag{E.45}$$

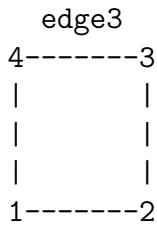
$$N_3(x(t), y(t)) = 0 \quad \text{by construction on edge 1} \tag{E.46}$$

$$N_4(x(t), y(t)) = 0 \quad \text{by construction on edge 1} \tag{E.47}$$

so that

$$\begin{aligned}
\mathbf{C}_1 &= h_x \begin{pmatrix} \int_0^1 (1-t)^2 dt & \int_0^1 (1-t)t dt & 0 & 0 \\ \int_0^1 t(1-t) dt & \int_0^1 t^2 dt & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
&= h_x \begin{pmatrix} 1/3 & 1/6 & 0 & 0 \\ 1/6 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
&= \frac{h_x}{6} \begin{pmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}
\end{aligned} \tag{E.48}$$

Computing matrix  $\mathbf{C}_3$



$$\mathbf{C}_3 = \int_{\partial\Omega_3} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma = \int_{3 \rightarrow 4} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma$$

The edge  $\partial\Omega_3$  is bounded by the coordinates of nodes  $x_3, y_3$  and  $x_4, y_4$ . This segment can be parameterised by  $t \in [0, 1]$ :

$$\vec{r}(t) = (1-t) \begin{pmatrix} x_3 \\ y_3 \end{pmatrix} + t \begin{pmatrix} x_4 \\ y_4 \end{pmatrix} = \begin{pmatrix} (x_4 - x_3)t + x_3 \\ (y_4 - y_3)t + y_3 \end{pmatrix} = \begin{pmatrix} (x_4 - x_3)t + x_3 \\ y_3 \end{pmatrix}$$

since  $y_3 = y_4$ . Here again the jacobian of the transformation is  $h_x$ .

$$\begin{aligned}
N_1(x(t), y(t)) &= 0 && \text{by construction on edge 3} \\
N_2(x(t), y(t)) &= 0 && \text{by construction on edge 3} \\
N_3(x(t), y(t)) &= \frac{x(t) - x_1}{x_3 - x_1} \frac{y(t) - y_1}{y_3 - y_1} \\
&= \frac{(x_4 - x_3)t + x_3 - x_1}{x_3 - x_1} \frac{y_3 - y_1}{y_3 - y_1} \\
&= \frac{(x_4 - x_3)t + x_3 - x_1}{x_3 - x_1} \\
&= 1 - t \\
N_4(x(t), y(t)) &= \frac{x(t) - x_3}{x_4 - x_3} \frac{y(t) - y_1}{y_4 - y_1} \\
&= \frac{(x_4 - x_3)t + x_3 - x_3}{x_4 - x_3} \frac{y_4 - y_1}{y_4 - y_1} \\
&= t
\end{aligned} \tag{E.49}$$

Then

$$\begin{aligned}
\mathbf{C}_3 &= h_x \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \int_0^1 (1-t)^2 dt & \int_0^1 (1-t)t dt \\ 0 & 0 & \int_0^1 t(1-t) dt & \int_0^1 t^2 dt \end{pmatrix} \\
&= h_x \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1/3 & 1/6 \\ 0 & 0 & 1/6 & 1/3 \end{pmatrix} \\
&= \frac{h_x}{6} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{pmatrix}
\end{aligned} \tag{E.50}$$

### The $\mathbb{K}$ and $\mathbb{G}$ matrices for $Q_1 \times P_0$ in 2D

Let us consider a regular grid composed of  $nelx \times nely$  rectangular linear elements on a domain of dimensions  $L_x \times L_y$ . We are here interested in the elemental matrices  $\mathbb{K}_e$  and  $\mathbb{G}_e$ . Borrowing (for example) from stone 1 we can write a simple code which compute these matrices by means of numerical integration. This code is available there: `python_codes/Gel/compute_K_G_S_q1p0.py`.

We start with square elements and a constant viscosity  $\eta = 1$ . We find that  $\mathbb{K}_e$  is independent of the resolution (i.e. independent of  $nelx = nely$ ):

$$\mathbb{K}_e = \begin{pmatrix} 1 & 0.25 & -0.5 & -0.25 & -0.5 & -0.25 & 0 & 0.25 \\ 0.25 & 1 & 0.25 & 0 & -0.25 & -0.5 & -0.25 & -0.5 \\ -0.5 & 0.25 & 1 & -0.25 & 0 & -0.25 & -0.5 & 0.25 \\ -0.25 & 0 & -0.25 & 1 & 0.25 & -0.5 & 0.25 & -0.5 \\ -0.5 & -0.25 & 0 & 0.25 & 1 & 0.25 & -0.5 & -0.25 \\ -0.25 & -0.5 & -0.25 & -0.5 & 0.25 & 1 & 0.25 & 0 \\ 0 & -0.25 & -0.5 & 0.25 & -0.5 & 0.25 & 1 & -0.25 \\ 0.25 & -0.5 & 0.25 & -0.5 & -0.25 & 0 & -0.25 & 1 \end{pmatrix}$$

and its lumped version  $\tilde{\mathbb{K}}_{i,j} = \sum_j |\mathbb{K}_{i,j}|$  is:

$$\tilde{\mathbb{K}}_e = \begin{pmatrix} 3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{pmatrix}$$

Let us now turn to non-square elements. Let us print  $\mathbb{K}_e$  and  $\mathbb{G}_e$  for 2 resolutions, 4x8 and 8x4:

$$\mathbb{K}_e = \begin{pmatrix} 1 & 0.25 & 0 & -0.25 & -0.5 & -0.25 & -0.5 & 0.25 \\ 0.25 & 1.5 & 0.25 & 0.5 & -0.25 & -0.75 & -0.25 & -1.25 \\ 0 & 0.25 & 1 & -0.25 & -0.5 & -0.25 & -0.5 & 0.25 \\ -0.25 & 0.5 & -0.25 & 1.5 & 0.25 & -1.25 & 0.25 & -0.75 \\ -0.5 & -0.25 & -0.5 & 0.25 & 1 & 0.25 & 0 & -0.25 \\ -0.25 & -0.75 & -0.25 & -1.25 & 0.25 & 1.5 & 0.25 & 0.5 \\ -0.5 & -0.25 & -0.5 & 0.25 & 0 & 0.25 & 1 & -0.25 \\ 0.25 & -1.25 & 0.25 & -0.75 & -0.25 & 0.5 & -0.25 & 1.5 \end{pmatrix} \quad \mathbb{G}_e = \begin{pmatrix} 0.0625 \\ 0.125 \\ -0.0625 \\ 0.125 \\ -0.0625 \\ -0.125 \\ 0.0625 \\ -0.125 \end{pmatrix}$$

$$\mathbb{K}_e = \begin{pmatrix} 1.5 & 0.25 & -1.25 & -0.25 & -0.75 & -0.25 & 0.5 & 0.25 \\ 0.25 & 1 & 0.25 & -0.5 & -0.25 & -0.5 & -0.25 & 0 \\ -1.25 & 0.25 & 1.5 & -0.25 & 0.5 & -0.25 & -0.75 & 0.25 \\ -0.25 & -0.5 & -0.25 & 1 & 0.25 & 0 & 0.25 & -0.5 \\ -0.75 & -0.25 & 0.5 & 0.25 & 1.5 & 0.25 & -1.25 & -0.25 \\ -0.25 & -0.5 & -0.25 & 0 & 0.25 & 1 & 0.25 & -0.5 \\ 0.5 & -0.25 & -0.75 & 0.25 & -1.25 & 0.25 & 1.5 & -0.25 \\ 0.25 & 0 & 0.25 & -0.5 & -0.25 & -0.5 & -0.25 & 1 \end{pmatrix} \quad \mathbb{G}_e = \begin{pmatrix} 0.125 \\ 0.0625 \\ -0.125 \\ 0.0625 \\ -0.125 \\ -0.0625 \\ 0.125 \\ -0.0625 \end{pmatrix}$$

We find that both matrices  $K_e$  and  $\mathbb{G}_e$  are different from each other and different from the ones obtained with square elements. We have no other choice than computing these by hand in order to express these as a function of  $h_x$  and  $h_y$ . Let us start with  $\mathbb{K}_e$ :

$$\mathbb{K}_e = \iint_{\Omega_e} \mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B} \, dV = \int_{x_1}^{x_3} \int_{y_1}^{y_3} \mathbf{B}^T(x, y) \cdot \mathbf{C}_\eta \cdot \mathbf{B}(x, y) \, dx dy \quad \text{with} \quad \mathbf{C}_\eta = \eta \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

In a rectangle bounded by  $[x_1, x_3] \times [y_1, y_3]$  the basis functions are given by:

$$\begin{aligned} \mathcal{N}_1^\gamma(x, y) &= \left( \frac{x_3 - x}{h_x} \right) \left( \frac{y_3 - y}{h_y} \right) \\ \mathcal{N}_2^\gamma(x, y) &= \left( \frac{x - x_1}{h_x} \right) \left( \frac{y_3 - y}{h_y} \right) \\ \mathcal{N}_3^\gamma(x, y) &= \left( \frac{x - x_1}{h_x} \right) \left( \frac{y - y_1}{h_y} \right) \\ \mathcal{N}_4^\gamma(x, y) &= \left( \frac{x_3 - x}{h_x} \right) \left( \frac{y - y_1}{h_y} \right) \end{aligned}$$

so that

$$\begin{aligned}
\frac{\partial \mathcal{N}_1^\gamma}{\partial x} &= -\frac{1}{h_x} \frac{y_3 - y}{h_y} = -\frac{1}{h_x} \frac{1}{2}(1 - s) \\
\frac{\partial \mathcal{N}_2^\gamma}{\partial x} &= \frac{1}{h_x} \frac{y_3 - y}{h_y} = +\frac{1}{h_x} \frac{1}{2}(1 - s) \\
\frac{\partial \mathcal{N}_3^\gamma}{\partial x} &= \frac{1}{h_x} \frac{y - y_1}{h_y} = +\frac{1}{h_x} \frac{1}{2}(1 + s) \\
\frac{\partial \mathcal{N}_4^\gamma}{\partial x} &= -\frac{1}{h_x} \frac{y - y_1}{h_y} = -\frac{1}{h_x} \frac{1}{2}(1 + s) \\
\frac{\partial \mathcal{N}_1^\gamma}{\partial y} &= -\frac{1}{h_y} \frac{x_3 - x}{h_x} = -\frac{1}{h_y} \frac{1}{2}(1 - r) \\
\frac{\partial \mathcal{N}_2^\gamma}{\partial y} &= -\frac{1}{h_y} \frac{x - x_1}{h_x} = -\frac{1}{h_y} \frac{1}{2}(1 + r) \\
\frac{\partial \mathcal{N}_3^\gamma}{\partial y} &= \frac{1}{h_y} \frac{x - x_1}{h_x} = +\frac{1}{h_y} \frac{1}{2}(1 + r) \\
\frac{\partial \mathcal{N}_4^\gamma}{\partial y} &= \frac{1}{h_y} \frac{x_3 - x}{h_x} = +\frac{1}{h_y} \frac{1}{2}(1 - r)
\end{aligned}$$

The matrix  $\mathbf{B}$  is given by

$$\mathbf{B}(x, y) = \begin{pmatrix} \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & 0 & \frac{\partial \mathcal{N}_2^\gamma}{\partial x} & 0 & \frac{\partial \mathcal{N}_3^\gamma}{\partial x} & 0 & \frac{\partial \mathcal{N}_4^\gamma}{\partial x} & 0 \\ 0 & \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & 0 & \frac{\partial \mathcal{N}_2^\gamma}{\partial y} & 0 & \frac{\partial \mathcal{N}_3^\gamma}{\partial y} & 0 & \frac{\partial \mathcal{N}_4^\gamma}{\partial y} \\ \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & \frac{\partial \mathcal{N}_2^\gamma}{\partial y} & \frac{\partial \mathcal{N}_2^\gamma}{\partial x} & \frac{\partial \mathcal{N}_3^\gamma}{\partial y} & \frac{\partial \mathcal{N}_3^\gamma}{\partial x} & \frac{\partial \mathcal{N}_4^\gamma}{\partial y} & \frac{\partial \mathcal{N}_4^\gamma}{\partial x} \end{pmatrix}$$

so that

$$\mathbf{C}_\eta \cdot \mathbf{B} = \eta \begin{pmatrix} 2\frac{\partial \mathcal{N}_1^\gamma}{\partial x} & 0 & 2\frac{\partial \mathcal{N}_2^\gamma}{\partial x} & 0 & 2\frac{\partial \mathcal{N}_3^\gamma}{\partial x} & 0 & 2\frac{\partial \mathcal{N}_4^\gamma}{\partial x} & 0 \\ 0 & 2\frac{\partial \mathcal{N}_1^\gamma}{\partial y} & 0 & 2\frac{\partial \mathcal{N}_2^\gamma}{\partial y} & 0 & 2\frac{\partial \mathcal{N}_3^\gamma}{\partial y} & 0 & 2\frac{\partial \mathcal{N}_4^\gamma}{\partial y} \\ \frac{\partial \mathcal{N}_1^\gamma}{\partial y} & \frac{\partial \mathcal{N}_1^\gamma}{\partial x} & \frac{\partial \mathcal{N}_2^\gamma}{\partial y} & \frac{\partial \mathcal{N}_2^\gamma}{\partial x} & \frac{\partial \mathcal{N}_3^\gamma}{\partial y} & \frac{\partial \mathcal{N}_3^\gamma}{\partial x} & \frac{\partial \mathcal{N}_4^\gamma}{\partial y} & \frac{\partial \mathcal{N}_4^\gamma}{\partial x} \end{pmatrix}$$

Let us start with the diagonal terms

$$\begin{aligned}
\mathbb{K}_{11} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{11} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \left( \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \right)^2 + \left( \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \right)^2 \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ 2 \left( -\frac{1}{h_x} \frac{1}{2} (1-s) \right)^2 + \left( -\frac{1}{h_y} \frac{1}{2} (1-r) \right)^2 \right] dr ds \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ \frac{1}{2h_x^2} (1-s)^2 + \frac{1}{4h_y^2} (1-r)^2 \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ \frac{1}{h_x^2} \frac{8}{3} + \frac{1}{2h_y^2} \frac{8}{3} \right] \\
&= \frac{1}{3} \left( \frac{2h_y}{h_x} + \frac{h_x}{h_y} \right) \\
\mathbb{K}_{22} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{22} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ \left( \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \right)^2 + 2 \left( \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \right)^2 \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ \left( -\frac{1}{h_x} \frac{1}{2} (1-s) \right)^2 + 2 \left( -\frac{1}{h_y} \frac{1}{2} (1-r) \right)^2 \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ \frac{1}{2h_x^2} \frac{8}{3} + \frac{1}{h_y^2} \frac{8}{3} \right] \\
&= \frac{1}{3} \left( \frac{h_y}{h_x} + \frac{2h_x}{h_y} \right) \\
\mathbb{K}_{33} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{33} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \left( \frac{\partial \mathcal{N}_2^\gamma}{\partial x} \right)^2 + \left( \frac{\partial \mathcal{N}_2^\gamma}{\partial y} \right)^2 \right] dx dy \\
&= \mathbb{K}_{11} \\
\mathbb{K}_{44} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{44} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ \left( \frac{\partial \mathcal{N}_2^\gamma}{\partial x} \right)^2 + 2 \left( \frac{\partial \mathcal{N}_2^\gamma}{\partial y} \right)^2 \right] dx dy \\
&= \mathbb{K}_{22} \\
\mathbb{K}_{55} &= \mathbb{K}_{11} \\
\mathbb{K}_{66} &= \mathbb{K}_{22} \\
\mathbb{K}_{77} &= \mathbb{K}_{11} \\
\mathbb{K}_{88} &= \mathbb{K}_{22}
\end{aligned}$$



Let us now focus on the first row of  $\mathbb{K}$ :

$$\begin{aligned}
\mathbb{K}_{12} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{11} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -\frac{1}{h_x} \frac{1}{2} (1-s) \cdot -\frac{1}{h_y} \frac{1}{2} (1-r) \right] dr ds \\
&= \frac{1}{16} \int_{-1}^{+1} \int_{-1}^{+1} (1-s)(1-r) dr ds \\
&= \frac{1}{4} \\
\mathbb{K}_{13} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{13} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_2^\gamma}{\partial x} + \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \frac{\partial \mathcal{N}_2^\gamma}{\partial y} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -2 \frac{1}{h_x} \frac{1}{2} (1-s) \cdot \frac{1}{h_x} \frac{1}{2} (1-s) - \frac{1}{h_y} \frac{1}{2} (1-r) \cdot -\frac{1}{h_y} \frac{1}{2} (1+r) \right] dr ds \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -\frac{1}{2h_x^2} (1-s)(1-s) + \frac{1}{4h_y^2} (1-r)(1+r) \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ -\frac{1}{2h_x^2} 2 \frac{8}{3} + \frac{1}{4h_y^2} \frac{4}{3} 2 \right] \\
&= \frac{1}{3} \left( -\frac{2h_y}{h_x} + \frac{h_x}{2h_y} \right) \\
\mathbb{K}_{14} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{11} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ \frac{\partial \mathcal{N}_2^\gamma}{\partial x} \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \right] dx dy \\
&= -\frac{1}{4} \\
\mathbb{K}_{15} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{15} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_3^\gamma}{\partial x} + \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \frac{\partial \mathcal{N}_3^\gamma}{\partial y} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -2 \frac{1}{h_x} \frac{1}{2} (1-s) \cdot \frac{1}{h_x} \frac{1}{2} (1+s) - \frac{1}{h_y} \frac{1}{2} (1-r) \cdot \frac{1}{h_y} \frac{1}{2} (1+r) \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ -\frac{1}{2h_x^2} \frac{4}{3} 2 - \frac{1}{4h_y^2} \frac{4}{3} 2 \right] \\
&= \frac{1}{3} \left( -\frac{h_y}{h_x} - \frac{h_x}{2h_y} \right) \\
\mathbb{K}_{16} &= -\frac{1}{4} \\
\mathbb{K}_{17} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{17} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_4^\gamma}{\partial x} + \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \frac{\partial \mathcal{N}_4^\gamma}{\partial y} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -2 \frac{1}{h_x} \frac{1}{2} (1-s) \cdot \frac{1}{h_x} \frac{1}{2} (1+s) - \frac{1}{h_y} \frac{1}{2} (1-r) \cdot \frac{1}{h_y} \frac{1}{2} (1-r) \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ \frac{1}{2h_x^2} 2 \frac{4}{3} - \frac{1}{4h_y^2} \frac{8}{3} 2 \right]
\end{aligned}$$

And now the other rows:

$$\begin{aligned}
\mathbb{K}_{23} &= \frac{1}{4} \\
\mathbb{K}_{24} &= -\mathbb{K}_{17} \\
\mathbb{K}_{25} &= -\frac{1}{4} \\
\mathbb{K}_{26} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{26} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \frac{\partial \mathcal{N}_3^\gamma}{\partial y} + \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_3^\gamma}{\partial x} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -2 \frac{1}{h_y} \frac{1}{2} (1-r) \cdot \frac{1}{h_y} \frac{1}{2} (1+r) - \frac{1}{h_x} \frac{1}{2} (1-s) \cdot \frac{1}{h_x} \frac{1}{2} (1+s) \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ -\frac{1}{2h_y^2} \frac{4}{3} 2 - \frac{1}{4h_x^2} 2 \frac{4}{3} \right] \\
&= \frac{1}{3} \left( -\frac{h_x}{h_y} - \frac{h_y}{2h_x} \right) \\
\mathbb{K}_{27} &= -\frac{1}{4} \\
\mathbb{K}_{28} &= \int_{x_1}^{x_3} \int_{y_1}^{y_3} (\mathbf{B}^T \cdot \mathbf{C}_\eta \cdot \mathbf{B})_{17} dx dy \\
&= \int_{x_1}^{x_3} \int_{y_1}^{y_3} \left[ 2 \frac{\partial \mathcal{N}_1^\gamma}{\partial y} \frac{\partial \mathcal{N}_4^\gamma}{\partial y} + \frac{\partial \mathcal{N}_1^\gamma}{\partial x} \frac{\partial \mathcal{N}_4^\gamma}{\partial x} \right] dx dy \\
&= \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \left[ -2 \frac{1}{h_y} \frac{1}{2} (1-r) \cdot \frac{1}{h_y} \frac{1}{2} (1-r) - \frac{1}{h_x} \frac{1}{2} (1-s) \cdot -\frac{1}{h_x} \frac{1}{2} (1+s) \right] dr ds \\
&= \frac{h_x h_y}{4} \left[ -\frac{1}{2h_y^2} \frac{8}{3} 2 + \frac{1}{4h_x^2} 2 \frac{4}{3} \right] \\
&= \frac{1}{3} \left( -\frac{2h_x}{h_y} + \frac{h_y}{2h_x} \right) \\
\mathbb{K}_{34} &= -\frac{1}{4} \\
\mathbb{K}_{35} &= \mathbb{K}_{17} \\
\mathbb{K}_{36} &= -\frac{1}{4} \\
\mathbb{K}_{37} &= \mathbb{K}_{15} \\
\mathbb{K}_{38} &= +\frac{1}{4} \\
\mathbb{K}_{45} &= \frac{1}{4} \\
\mathbb{K}_{46} &= \mathbb{K}_{28} \\
\mathbb{K}_{47} &= \frac{1}{4} \\
\mathbb{K}_{48} &= \mathbb{K}_{26} \\
\mathbb{K}_{56} &= \frac{1}{4} \\
\mathbb{K}_{57} &= \mathbb{K}_{13} \\
\mathbb{K}_{58} &= -\frac{1}{4} \\
\mathbb{K}_{67} &= \frac{1}{4} \\
\mathbb{K}_{68} &= -\mathbb{K}_{17} \\
\mathbb{K}_{78} &= -\frac{1}{4}
\end{aligned}$$

so that (matrix  $\mathbb{K}_e$  is symmetric so only half is shown here)<sup>2</sup>:

$$\begin{pmatrix} \frac{1}{3}\left(\frac{2h_y}{h_x} + \frac{h_x}{h_y}\right) & \frac{1}{4} & \frac{1}{3}\left(-\frac{2h_y}{h_x} + \frac{h_x}{2h_y}\right) & -\frac{1}{4} & \frac{1}{3}\left(-\frac{h_y}{h_x} - \frac{h_x}{2h_y}\right) & -\frac{1}{4} & \frac{1}{3}\left(\frac{h_y}{h_x} - \frac{h_x}{h_y}\right) & \frac{1}{4} \\ \cdot & \frac{1}{3}\left(\frac{h_y}{h_x} + \frac{2h_x}{h_y}\right) & \frac{1}{4} & -\frac{1}{3}\left(\frac{h_y}{h_x} - \frac{h_x}{h_y}\right) & -\frac{1}{4} & \frac{1}{3}\left(-\frac{h_x}{h_y} - \frac{h_y}{2h_x}\right) & -\frac{1}{4} & \frac{1}{3}\left(-\frac{2h_x}{h_y} + \frac{h_y}{2h_x}\right) \\ \cdot & \cdot & \frac{1}{3}\left(\frac{2h_y}{h_x} + \frac{h_x}{h_y}\right) & -\frac{1}{4} & \frac{1}{3}\left(\frac{h_y}{h_x} - \frac{h_x}{h_y}\right) & -\frac{1}{4} & \frac{1}{3}\left(-\frac{h_y}{h_x} - \frac{h_x}{2h_y}\right) & \frac{1}{4} \\ \cdot & \cdot & \cdot & \frac{1}{3}\left(\frac{h_y}{h_x} + \frac{2h_x}{h_y}\right) & \frac{1}{4} & \frac{1}{3}\left(-\frac{2h_x}{h_y} + \frac{h_y}{2h_x}\right) & \frac{1}{4} & \frac{1}{3}\left(-\frac{h_x}{h_y} - \frac{h_y}{2h_x}\right) \\ \cdot & \cdot & \cdot & \cdot & \frac{1}{3}\left(\frac{2h_y}{h_x} + \frac{h_x}{h_y}\right) & \frac{1}{4} & \frac{1}{3}\left(-\frac{2h_y}{h_x} + \frac{h_x}{2h_y}\right) & -\frac{1}{4} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \frac{1}{3}\left(\frac{h_y}{h_x} + \frac{2h_x}{h_y}\right) & \frac{1}{4} & -\frac{1}{3}\left(\frac{h_y}{h_x} - \frac{h_x}{h_y}\right) \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \frac{1}{3}\left(\frac{2h_y}{h_x} + \frac{h_x}{h_y}\right) & -\frac{1}{4} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \frac{1}{3}\left(\frac{h_y}{h_x} + \frac{2h_x}{h_y}\right) \end{pmatrix}$$

For reference I present here under 3 matrices obtained with three coarse resolutions:

- $4 \times 4$  mesh

$$\begin{bmatrix} 1.0000 & 0.2500 & -0.5000 & -0.2500 & -0.5000 & -0.2500 & 0.0000 & 0.2500 \\ 0.2500 & 1.0000 & 0.2500 & -0.0000 & -0.2500 & -0.5000 & -0.2500 & -0.5000 \\ -0.5000 & 0.2500 & 1.0000 & -0.2500 & 0.0000 & -0.2500 & -0.5000 & 0.2500 \\ -0.2500 & -0.0000 & -0.2500 & 1.0000 & 0.2500 & -0.5000 & 0.2500 & -0.5000 \\ -0.5000 & -0.2500 & 0.0000 & 0.2500 & 1.0000 & 0.2500 & -0.5000 & -0.2500 \\ -0.2500 & -0.5000 & -0.2500 & -0.5000 & 0.2500 & 1.0000 & 0.2500 & -0.0000 \\ 0.0000 & -0.2500 & -0.5000 & 0.2500 & -0.5000 & 0.2500 & 1.0000 & -0.2500 \\ 0.2500 & -0.5000 & 0.2500 & -0.5000 & -0.2500 & -0.0000 & -0.2500 & 1.0000 \end{bmatrix}$$

- $7 \times 5$  mesh

$$\begin{bmatrix} 1.1714 & 0.2500 & -0.8143 & -0.2500 & -0.5857 & -0.2500 & 0.2286 & 0.2500 \\ 0.2500 & 0.9429 & 0.2500 & -0.2286 & -0.2500 & -0.4714 & -0.2500 & -0.2429 \\ -0.8143 & 0.2500 & 1.1714 & -0.2500 & 0.2286 & -0.2500 & -0.5857 & 0.2500 \\ -0.2500 & -0.2286 & -0.2500 & 0.9429 & 0.2500 & -0.2429 & 0.2500 & -0.4714 \\ -0.5857 & -0.2500 & 0.2286 & 0.2500 & 1.1714 & 0.2500 & -0.8143 & -0.2500 \\ -0.2500 & -0.4714 & -0.2500 & -0.2429 & 0.2500 & 0.9429 & 0.2500 & -0.2286 \\ 0.2286 & -0.2500 & -0.5857 & 0.2500 & -0.8143 & 0.2500 & 1.1714 & -0.2500 \\ 0.2500 & -0.2429 & 0.2500 & -0.4714 & -0.2500 & -0.2286 & -0.2500 & 0.9429 \end{bmatrix}$$

- $5 \times 7$  mesh

$$\begin{bmatrix} 0.9429 & 0.2500 & -0.2429 & -0.2500 & -0.4714 & -0.2500 & -0.2286 & 0.2500 \\ 0.2500 & 1.1714 & 0.2500 & 0.2286 & -0.2500 & -0.5857 & -0.2500 & -0.8143 \\ -0.2429 & 0.2500 & 0.9429 & -0.2500 & -0.2286 & -0.2500 & -0.4714 & 0.2500 \\ -0.2500 & 0.2286 & -0.2500 & 1.1714 & 0.2500 & -0.8143 & 0.2500 & -0.5857 \\ -0.4714 & -0.2500 & -0.2286 & 0.2500 & 0.9429 & 0.2500 & -0.2429 & -0.2500 \\ -0.2500 & -0.5857 & -0.2500 & -0.8143 & 0.2500 & 1.1714 & 0.2500 & 0.2286 \\ -0.2286 & -0.2500 & -0.4714 & 0.2500 & -0.2429 & 0.2500 & 0.9429 & -0.2500 \\ 0.2500 & -0.8143 & 0.2500 & -0.5857 & -0.2500 & 0.2286 & -0.2500 & 1.1714 \end{bmatrix}$$

<sup>2</sup>Values above are fully checked, values in matrix below should be re-checked to be sure

Turning now to the lumped version of  $\mathbb{K}_e$ , and because the sum of  $\pm 1/4$  terms always add up to 1, we find:

$$\begin{aligned}\tilde{\mathbb{K}}_{11} &= |\mathbb{K}_{11}| + |\mathbb{K}_{13}| + |\mathbb{K}_{15}| + |\mathbb{K}_{17}| + 1 \\ \tilde{\mathbb{K}}_{22} &= |\mathbb{K}_{22}| + |\mathbb{K}_{24}| + |\mathbb{K}_{26}| + |\mathbb{K}_{28}| + 1 \\ \tilde{\mathbb{K}}_{33} &= |\mathbb{K}_{31}| + |\mathbb{K}_{33}| + |\mathbb{K}_{35}| + |\mathbb{K}_{37}| + 1 \\ \tilde{\mathbb{K}}_{44} &= |\mathbb{K}_{42}| + |\mathbb{K}_{44}| + |\mathbb{K}_{46}| + |\mathbb{K}_{48}| + 1 \\ \tilde{\mathbb{K}}_{55} &= |\mathbb{K}_{51}| + |\mathbb{K}_{53}| + |\mathbb{K}_{55}| + |\mathbb{K}_{57}| + 1 \\ \tilde{\mathbb{K}}_{66} &= |\mathbb{K}_{62}| + |\mathbb{K}_{64}| + |\mathbb{K}_{66}| + |\mathbb{K}_{68}| + 1 \\ \tilde{\mathbb{K}}_{77} &= |\mathbb{K}_{71}| + |\mathbb{K}_{73}| + |\mathbb{K}_{75}| + |\mathbb{K}_{77}| + 1 \\ \tilde{\mathbb{K}}_{88} &= |\mathbb{K}_{82}| + |\mathbb{K}_{84}| + |\mathbb{K}_{86}| + |\mathbb{K}_{88}| + 1\end{aligned}$$

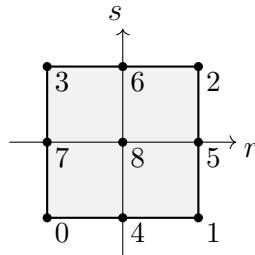
Turning now to the  $\mathbb{G}_e$  matrix, we find:

$$\mathbb{G}_e = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathcal{N}^p dV = - \int_{\Omega_e} \begin{pmatrix} \frac{\partial \mathcal{N}_1}{\partial x} \\ \frac{\partial \mathcal{N}_1}{\partial y} \\ \frac{\partial \mathcal{N}_2}{\partial x} \\ \frac{\partial \mathcal{N}_2}{\partial y} \\ \frac{\partial \mathcal{N}_3}{\partial x} \\ \frac{\partial \mathcal{N}_3}{\partial y} \\ \frac{\partial \mathcal{N}_4}{\partial x} \\ \frac{\partial \mathcal{N}_4}{\partial y} \end{pmatrix} dV = - \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} \begin{pmatrix} -\frac{1}{h_x} \frac{1}{2} (1-s) \\ -\frac{1}{h_y} \frac{1}{2} (1-r) \\ +\frac{1}{h_x} \frac{1}{2} (1-s) \\ -\frac{1}{h_y} \frac{1}{2} (1+r) \\ +\frac{1}{h_x} \frac{1}{2} (1+s) \\ +\frac{1}{h_y} \frac{1}{2} (1+r) \\ -\frac{1}{h_x} \frac{1}{2} (1+s) \\ +\frac{1}{h_y} \frac{1}{2} (1-r) \end{pmatrix} dr ds = \begin{pmatrix} h_y/2 \\ h_x/2 \\ -h_y/2 \\ +h_x/2 \\ -h_y/2 \\ -h_x/2 \\ +h_y/2 \\ -h_x/2 \end{pmatrix}$$

which is Eq. (3.65) in Elman, Silvester, and Wathen [371]. Since the pressure is constant inside each element, then  $\mathbb{G}_{el}$  is  $(ndof_V * m_V, m_P) = (8 \times 1)$ .

### E.0.3 Quadrilaterals: rectangular quadratic elements

(tikz-q22d.tex)



$$\vec{\mathcal{N}}(r, s) = \begin{pmatrix} \mathcal{N}_1(r, s) \\ \mathcal{N}_2(r, s) \\ \mathcal{N}_3(r, s) \\ \mathcal{N}_4(r, s) \\ \mathcal{N}_5(r, s) \\ \mathcal{N}_6(r, s) \\ \mathcal{N}_7(r, s) \\ \mathcal{N}_8(r, s) \\ \mathcal{N}_9(r, s) \end{pmatrix} = \begin{pmatrix} \mathcal{N}_1(r) \mathcal{N}_1(s) \\ \mathcal{N}_2(r) \mathcal{N}_1(s) \\ \mathcal{N}_3(r) \mathcal{N}_1(s) \\ \mathcal{N}_1(r) \mathcal{N}_2(s) \\ \mathcal{N}_2(r) \mathcal{N}_2(s) \\ \mathcal{N}_3(r) \mathcal{N}_2(s) \\ \mathcal{N}_1(r) \mathcal{N}_3(s) \\ \mathcal{N}_2(r) \mathcal{N}_3(s) \\ \mathcal{N}_3(r) \mathcal{N}_3(s) \end{pmatrix} = \begin{pmatrix} \frac{1}{2} r(r-1) \frac{1}{2} s(s-1) \\ (1-r^2) \frac{1}{2} s(s-1) \\ \frac{1}{2} r(r+1) \frac{1}{2} s(s-1) \\ \frac{1}{2} r(r-1) (1-s^2) \\ (1-r^2) (1-s^2) \frac{1}{2} r(r+1) (1-s^2) \\ \frac{1}{2} r(r-1) \frac{1}{2} s(s+1) \\ (1-r^2) \frac{1}{2} s(s+1) \\ \frac{1}{2} r(r+1) \frac{1}{2} s(s+1) \end{pmatrix}$$

The mass matrix on the reference element is then

$$\mathbf{M}^e = \iint_{\square} \begin{pmatrix} \mathcal{N}_1\mathcal{N}_1 & \mathcal{N}_1\mathcal{N}_2 & \mathcal{N}_1\mathcal{N}_3 & \mathcal{N}_1\mathcal{N}_4 & \mathcal{N}_1\mathcal{N}_5 & \mathcal{N}_1\mathcal{N}_6 & \mathcal{N}_1\mathcal{N}_7 & \mathcal{N}_1\mathcal{N}_8 & \mathcal{N}_1\mathcal{N}_9 \\ \mathcal{N}_2\mathcal{N}_1 & \mathcal{N}_2\mathcal{N}_2 & \mathcal{N}_2\mathcal{N}_3 & \mathcal{N}_2\mathcal{N}_4 & \mathcal{N}_2\mathcal{N}_5 & \mathcal{N}_2\mathcal{N}_6 & \mathcal{N}_2\mathcal{N}_7 & \mathcal{N}_2\mathcal{N}_8 & \mathcal{N}_2\mathcal{N}_9 \\ \vdots & & & & & & & & \\ \mathcal{N}_9\mathcal{N}_1 & \mathcal{N}_9\mathcal{N}_2 & \mathcal{N}_9\mathcal{N}_3 & \mathcal{N}_9\mathcal{N}_4 & \mathcal{N}_9\mathcal{N}_5 & \mathcal{N}_9\mathcal{N}_6 & \mathcal{N}_9\mathcal{N}_7 & \mathcal{N}_9\mathcal{N}_8 & \mathcal{N}_9\mathcal{N}_9 \end{pmatrix} drds$$

with for example<sup>3</sup>

$$\iint_{\square} \mathcal{N}_1(r, s)\mathcal{N}_1(r, s) drds = \frac{16}{225} \quad (\text{E.51})$$

$$\iint_{\square} \mathcal{N}_1(r, s)\mathcal{N}_2(r, s) drds = \frac{8}{225} \quad (\text{E.52})$$

$$\iint_{\square} \mathcal{N}_1(r, s)\mathcal{N}_3(r, s) drds = -\frac{4}{225} \quad (\text{E.53})$$

$$\iint_{\square} \mathcal{N}_1(r, s)\mathcal{N}_8(r, s) drds = \frac{-2}{225} \quad (\text{E.54})$$

$$\iint_{\square} \mathcal{N}_1(r, s)\mathcal{N}_9(r, s) drds = \frac{1}{225} \quad (\text{E.55})$$

$$\iint_{\square} \mathcal{N}_5(r, s)\mathcal{N}_5(r, s) drds = \frac{256}{225} \quad (\text{E.56})$$

$$(\text{E.57})$$

In [STONE](#) 107, we find for a 3x2 mesh on domain 6x4:

$$\begin{aligned} \mathbf{M} = & \\ & \begin{bmatrix} 16. & 8. & -4. & 8. & 4. & -2. & -4. & -2. & 1. \\ 8. & 64. & 8. & 4. & 32. & 4. & -2. & -16. & -2. \\ -4. & 8. & 16. & -2. & 4. & 8. & 1. & -2. & -4. \\ 8. & 4. & -2. & 64. & 32. & -16. & 8. & 4. & -2. \\ 4. & 32. & 4. & 32. & 256. & 32. & 4. & 32. & 4. \\ -2. & 4. & 8. & -16. & 32. & 64. & -2. & 4. & 8. \\ -4. & -2. & 1. & 8. & 4. & -2. & 16. & 8. & -4. \\ -2. & -16. & -2. & 4. & 32. & 4. & 8. & 64. & 8. \\ 1. & -2. & -4. & -2. & 4. & 8. & -4. & 8. & 16. \end{bmatrix} / 225 \end{aligned}$$

---

<sup>3</sup>Thank you WolframAlpha again!

## E.0.4 Hexahedra: cuboid elements

We here assume that each element is a cuboid<sup>4</sup>. We set the domain size to  $L_x = 4$ ,  $L_y = 3$  and  $L_z = 2$ , with  $nelx = 4$ ,  $nely = 3$  and  $nelz = 2$ . Here again the viscosity is set to  $\eta = 1$  so that we find that

$$\mathbb{K}_{el} = \frac{1}{8 \cdot 9} \begin{pmatrix} 32 & 6 & 6 & -8 & -6 & -6 & -10 & -6 & -3 & 4 & 6 & 3 & 4 & 3 & 6 & -10 & -3 & -6 \\ 6 & 32 & 6 & 6 & 4 & 3 & -6 & -10 & -3 & -6 & -8 & -6 & 3 & 4 & 6 & 3 & -4 & 3 \\ 6 & 6 & 32 & 6 & 3 & 4 & 3 & 3 & -4 & 3 & 6 & 4 & -6 & -6 & -8 & -6 & -3 & -10 \\ -8 & 6 & 6 & 32 & -6 & -6 & 4 & -6 & -3 & -10 & 6 & 3 & -10 & 3 & 6 & 4 & -3 & -6 \\ -6 & 4 & 3 & -6 & 32 & 6 & 6 & -8 & -6 & 6 & -10 & -3 & -3 & -4 & 3 & -3 & 4 & 6 \\ -6 & 3 & 4 & -6 & 6 & 32 & -3 & 6 & 4 & -3 & 3 & -4 & 6 & -3 & -10 & 6 & -6 & -8 \\ 10 & -6 & 3 & 4 & 6 & -3 & 32 & 6 & -6 & -8 & -6 & 6 & -8 & -3 & 3 & -4 & 3 & -3 \\ -6 & -10 & 3 & -6 & -8 & 6 & 6 & 32 & -6 & 6 & 4 & -3 & -3 & -8 & 3 & -3 & -10 & 6 \\ -3 & -3 & -4 & -3 & -6 & 4 & -6 & -6 & 32 & -6 & -3 & 4 & 3 & 3 & -8 & 3 & 6 & -10 \\ 4 & -6 & 3 & -10 & 6 & -3 & -8 & 6 & -6 & 32 & -6 & 6 & -4 & -3 & 3 & -8 & 3 & -3 \\ 6 & -8 & 6 & 6 & -10 & 3 & -6 & 4 & -3 & -6 & 32 & -6 & 3 & -10 & 6 & 3 & -8 & 3 \\ 3 & -6 & 4 & 3 & -3 & -4 & 6 & -3 & 4 & 6 & -6 & 32 & -3 & 6 & -10 & -3 & 3 & -8 \\ 4 & 3 & -6 & -10 & -3 & 6 & -8 & -3 & 3 & -4 & 3 & -3 & 32 & 6 & -6 & -8 & -6 & 6 \\ 3 & 4 & -6 & 3 & -4 & -3 & -3 & -8 & 3 & -3 & -10 & 6 & 6 & 32 & -6 & 6 & 4 & -3 \\ 6 & 6 & -8 & 6 & 3 & -10 & 3 & 3 & -8 & 3 & 6 & -10 & -6 & -6 & 32 & -6 & -3 & 4 \\ 10 & 3 & -6 & 4 & -3 & 6 & -4 & -3 & 3 & -8 & 3 & -3 & -8 & 6 & -6 & 32 & -6 & 6 \\ -3 & -4 & -3 & -3 & 4 & -6 & 3 & -10 & 6 & 3 & -8 & 3 & -6 & 4 & -3 & -6 & 32 & -6 \\ -6 & 3 & -10 & -6 & 6 & -8 & -3 & 6 & -10 & -3 & 3 & -8 & 6 & -3 & 4 & 6 & -6 & 32 \\ -8 & -3 & -3 & -4 & 3 & 3 & 4 & 3 & 6 & -10 & -3 & -6 & -10 & -6 & -3 & 4 & 6 & 3 \\ -3 & -8 & -3 & -3 & -10 & -6 & 3 & 4 & 6 & 3 & -4 & 3 & -6 & -10 & -3 & -6 & -8 & -6 \\ -3 & -3 & -8 & -3 & -6 & -10 & -6 & -6 & -8 & -6 & -3 & -10 & 3 & 3 & -4 & 3 & 6 & 4 \\ -4 & -3 & -3 & -8 & 3 & 3 & -10 & 3 & 6 & 4 & -3 & -6 & 4 & -6 & -3 & -10 & 6 & 3 \\ 3 & -10 & -6 & 3 & -8 & -3 & -3 & -4 & 3 & -3 & 4 & 6 & 6 & -8 & -6 & 6 & -10 & -3 \\ 3 & -6 & -10 & 37 & -3 & -8 & 6 & -3 & -10 & 6 & -6 & -8 & -3 & 6 & 4 & -3 & 3 & -4 \end{pmatrix}$$

and

$$\mathbb{G}_{el} = \frac{1}{2 \cdot 9} \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \\ 1 \\ -1 \\ -1 \\ 1 \\ 1 \\ -1 \\ 1 \\ 1 \\ -1 \\ -1 \\ 1 \\ -1 \\ -1 \\ -1 \\ 1 \\ -1 \\ -1 \\ -1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$$

<sup>4</sup><https://en.wikipedia.org/wiki/Cuboid>

### E.0.5 Triangles: linear elements

We start from the linear basis functions in the reference triangle given as a function of  $r, s$ :

$$\mathcal{N}_1(r, s) = 1 - r - s \quad (\text{E.58})$$

$$\mathcal{N}_2(r, s) = r \quad (\text{E.59})$$

$$\mathcal{N}_3(r, s) = s \quad (\text{E.60})$$

and their derivatives:

$$\partial_r \mathcal{N}_1(r, s) = -1$$

$$\partial_r \mathcal{N}_2(r, s) = 1$$

$$\partial_r \mathcal{N}_3(r, s) = 0$$

$$\partial_s \mathcal{N}_1(r, s) = -1$$

$$\partial_s \mathcal{N}_2(r, s) = 0$$

$$\partial_s \mathcal{N}_3(r, s) = 1$$

We wish to compute the integral of a function  $f(x, y)$  over the triangle by means of a change of variables  $(x, y) \rightarrow (r, s)$ :

$$\begin{aligned} \iint f(x, y) dx dy &= \iint f(x(r, s), y(r, s)) \left| \frac{\partial(x, y)}{\partial(r, s)} \right| dr ds \\ &= \iint f(x(r, s), y(r, s)) \begin{vmatrix} \partial x / \partial r & \partial x / \partial s \\ \partial y / \partial r & \partial y / \partial s \end{vmatrix} dr ds \end{aligned} \quad (\text{E.61})$$

From

$$x(r, s) = \sum_{i=1}^3 \mathcal{N}_i(r, s) x_i \quad \text{and} \quad y(r, s) = \sum_{i=1}^3 \mathcal{N}_i(r, s) y_i$$

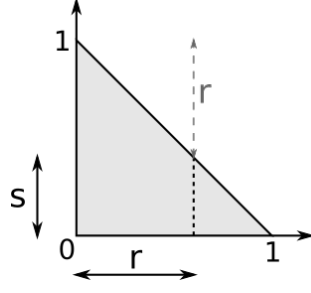
we can write

$$\begin{aligned} \frac{\partial x}{\partial r}(r, s) &= \sum_{i=1}^3 \frac{\partial \mathcal{N}_i}{\partial r} x_i = \frac{\partial \mathcal{N}_1}{\partial r} x_1 + \frac{\partial \mathcal{N}_2}{\partial r} x_2 + \frac{\partial \mathcal{N}_3}{\partial r} x_3 = -x_1 + x_2 \\ \frac{\partial x}{\partial s}(r, s) &= \sum_{i=1}^3 \frac{\partial \mathcal{N}_i}{\partial s} x_i = \frac{\partial \mathcal{N}_1}{\partial s} x_1 + \frac{\partial \mathcal{N}_2}{\partial s} x_2 + \frac{\partial \mathcal{N}_3}{\partial s} x_3 = -x_1 + x_3 \\ \frac{\partial y}{\partial r}(r, s) &= \sum_{i=1}^3 \frac{\partial \mathcal{N}_i}{\partial r} y_i = \frac{\partial \mathcal{N}_1}{\partial r} y_1 + \frac{\partial \mathcal{N}_2}{\partial r} y_2 + \frac{\partial \mathcal{N}_3}{\partial r} y_3 = -y_1 + y_2 \\ \frac{\partial y}{\partial s}(r, s) &= \sum_{i=1}^3 \frac{\partial \mathcal{N}_i}{\partial s} y_i = \frac{\partial \mathcal{N}_1}{\partial s} y_1 + \frac{\partial \mathcal{N}_2}{\partial s} y_2 + \frac{\partial \mathcal{N}_3}{\partial s} y_3 = -y_1 + y_3 \end{aligned} \quad (\text{E.62})$$

Then

$$\begin{vmatrix} \partial x / \partial r & \partial x / \partial s \\ \partial y / \partial r & \partial y / \partial s \end{vmatrix} = \begin{vmatrix} -x_1 + x_2 & -x_1 + x_3 \\ -y_1 + y_2 & -y_1 + y_3 \end{vmatrix} = (x_2 - x_1)(y_3 - y_1) - (y_2 - y_1)(x_3 - x_1) = 2S$$

where  $S$  is the area of the triangle and which is independent of  $(r, s)$ . Looking at the reference element, we find that when  $r$  goes from 0 to 1,  $s$  can only take values between 0 and  $1 - r$ .



Then the bounds of the integrals are simply:

$$\iint_{\triangle} f(x,y) dx dy = 2S \int_0^1 \left( \int_0^{1-r} f(x(r,s), y(r,s)) ds \right) dr \quad (\text{E.63})$$

and the mass matrix is given by

$$\mathbf{M}_e = 2S \int_0^1 \left[ \int_0^{1-r} \begin{pmatrix} (1-r-s)^2 & (1-r-s)r & (1-r-s)s \\ (1-r-s)r & r^2 & rs \\ (1-r-s)s & rs & s^2 \end{pmatrix} ds \right] dr \quad (\text{E.64})$$

$$= 2S \int_0^1 \begin{pmatrix} \int_0^{1-r} (1-r-s)^2 ds & \int_0^{1-r} (1-r-s)r ds & \int_0^{1-r} (1-r-s)s ds \\ \int_0^{1-r} (1-r-s)r ds & \int_0^{1-r} r^2 ds & \int_0^{1-r} rs ds \\ \int_0^{1-r} (1-r-s)s ds & \int_0^{1-r} rs ds & \int_0^{1-r} s^2 ds \end{pmatrix} dr \quad (\text{E.65})$$

$$= 2S \begin{pmatrix} 1/12 & 1/24 & 1/24 \\ 1/24 & 1/12 & 1/24 \\ 1/24 & 1/24 & 1/12 \end{pmatrix} \quad (\text{E.66})$$

$$= \frac{S}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} \quad (\text{E.67})$$

This is Eq.(4.10e) of Li [779]. Also note that in the context of the heat transport equation this matrix is multiplied by  $\rho C_p$ .

We will then compute the  $\mathbf{J}_x$  and  $\mathbf{J}_y$  matrices. We start from the basis functions expressed in the  $(x,y)$  coordinate system:

$$\begin{aligned} \mathcal{N}_1(x,y) &= \frac{1}{2S} (x_2 y_3 - x_3 y_2 + (y_2 - y_3)x + (x_3 - x_2)y) \\ \mathcal{N}_2(x,y) &= \frac{1}{2S} (x_3 y_1 - x_1 y_3 + (y_3 - y_1)x + (x_1 - x_3)y) \\ \mathcal{N}_3(x,y) &= \frac{1}{2S} (x_1 y_2 - x_2 y_1 + (y_1 - y_2)x + (x_2 - x_1)y) \end{aligned}$$



where  $S$  is the area of the element. We then have

$$\begin{aligned}
\partial_x \mathcal{N}_1(x, y) &= \frac{1}{2S}(y_2 - y_3) \\
\partial_x \mathcal{N}_2(x, y) &= \frac{1}{2S}(y_3 - y_1) \\
\partial_x \mathcal{N}_3(x, y) &= \frac{1}{2S}(y_1 - y_2) \\
\partial_y \mathcal{N}_1(x, y) &= \frac{1}{2S}(x_3 - x_2) \\
\partial_y \mathcal{N}_2(x, y) &= \frac{1}{2S}(x_1 - x_3) \\
\partial_y \mathcal{N}_3(x, y) &= \frac{1}{2S}(x_2 - x_1)
\end{aligned}$$

We start with

$$\begin{aligned}
\mathbf{J}_x &= \iint_{\Delta} \partial_x \vec{\mathcal{N}}^T \vec{\mathcal{N}} dV \\
&= \iint_{\Delta} \begin{pmatrix} \frac{1}{2S}(y_2 - y_3) \\ \frac{1}{2S}(y_3 - y_1) \\ \frac{1}{2S}(y_1 - y_2) \end{pmatrix} \begin{pmatrix} \mathcal{N}_1(x, y) & \mathcal{N}_2(x, y) & \mathcal{N}_3(x, y) \end{pmatrix} dxdy \\
&= \frac{1}{2S} \begin{pmatrix} y_{23} \iint_{\Delta} \mathcal{N}_1 dxdy & y_{23} \iint_{\Delta} \mathcal{N}_2 dxdy & y_{23} \iint_{\Delta} \mathcal{N}_3 dxdy \\ y_{31} \iint_{\Delta} \mathcal{N}_1 dxdy & y_{31} \iint_{\Delta} \mathcal{N}_2 dxdy & y_{31} \iint_{\Delta} \mathcal{N}_3 dxdy \\ y_{12} \iint_{\Delta} \mathcal{N}_1 dxdy & y_{12} \iint_{\Delta} \mathcal{N}_2 dxdy & y_{12} \iint_{\Delta} \mathcal{N}_3 dxdy \end{pmatrix} \quad (\text{E.68})
\end{aligned}$$

where we have introduced the notation  $x_{ij} = x_i - x_j$ . We then need to compute

$$\begin{aligned}
\iint_{\Delta} \mathcal{N}_1(x, y) dxdy &= 2S \int_0^1 \left( \int_0^{1-r} \mathcal{N}_1(x(r, s), y(r, s)) ds \right) dr \\
&= 2S \int_0^1 \left( \int_0^{1-r} (1 - r - s) ds \right) dr \\
&= 2S \frac{1}{6} \quad (\text{E.69})
\end{aligned}$$

$$\begin{aligned}
\iint_{\Delta} \mathcal{N}_2(x, y) dxdy &= 2S \int_0^1 \left( \int_0^{1-r} \mathcal{N}_2(x(r, s), y(r, s)) ds \right) dr \\
&= 2S \int_0^1 \left( \int_0^{1-r} r ds \right) dr \\
&= 2S \frac{1}{6} \quad (\text{E.70})
\end{aligned}$$

$$\begin{aligned}
\iint_{\Delta} \mathcal{N}_3(x, y) dxdy &= 2S \int_0^1 \left( \int_0^{1-r} \mathcal{N}_3(x(r, s), y(r, s)) ds \right) dr \\
&= 2S \int_0^1 \left( \int_0^{1-r} s ds \right) dr \\
&= 2S \frac{1}{6} \quad (\text{E.71})
\end{aligned}$$

verify!!

Finally:

$$\mathbf{J}_x = \frac{1}{6} \begin{pmatrix} y_{23} & y_{23} & y_{23} \\ y_{31} & y_{31} & y_{31} \\ y_{12} & y_{12} & y_{12} \end{pmatrix}$$

Likewise

$$\begin{aligned}
\mathbf{J}_y &= \iint_{\Delta} \partial_y \vec{\mathcal{N}}^T \vec{\mathcal{N}} dV \\
&= \iint_{\Delta} \begin{pmatrix} \frac{1}{2S}(x_3 - x_2) \\ \frac{1}{2S}(x_1 - x_3) \\ \frac{1}{2S}(x_2 - x_1) \end{pmatrix} \begin{pmatrix} \mathcal{N}_1(x, y) & \mathcal{N}_2(x, y) & \mathcal{N}_3(x, y) \end{pmatrix} dx dy \\
&= \frac{1}{6} \begin{pmatrix} x_{32} & x_{32} & x_{32} \\ x_{13} & x_{13} & x_{13} \\ x_{21} & x_{21} & x_{21} \end{pmatrix}
\end{aligned} \tag{E.72}$$

We now turn to the other two matrices, the advection  $\mathbb{K}_a$  and diffusion  $\mathbb{K}_d$  matrices. The gradient matrix  $\mathbf{B}$  is given by

$$\mathbf{B} = \begin{pmatrix} \partial_x \mathcal{N}_1 & \partial_x \mathcal{N}_2 & \partial_x \mathcal{N}_3 \\ \partial_y \mathcal{N}_1 & \partial_y \mathcal{N}_2 & \partial_y \mathcal{N}_3 \end{pmatrix} = \frac{1}{2S} \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix}$$

then

$$\mathbf{K}_d = \iint_{\Delta} \mathbf{B}^T k \mathbf{B} dV = \iint_{\Delta} \frac{k}{4S^2} \begin{pmatrix} y_{23} & x_{32} \\ y_{31} & x_{13} \\ y_{12} & x_{21} \end{pmatrix} \cdot \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix} dV$$

If  $k$  is constant within the element, then

$$\mathbf{K}_d = \frac{k}{4S} \begin{pmatrix} y_{23} & x_{32} \\ y_{31} & x_{13} \\ y_{12} & x_{21} \end{pmatrix} \cdot \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix}$$

Turning now to the advection matrix

$$\begin{aligned}
\mathbf{K}_a &= \iint_{\Delta} \vec{\mathcal{N}}^T \vec{\mathbf{v}} \cdot \mathbf{B} dV \\
&= \iint_{\Delta} \vec{\mathcal{N}}^T(x, y) \vec{\mathbf{v}}(x, y) \cdot \mathbf{B}(x, y) dx dy \\
&= 2S \iint_{\Delta} \vec{\mathcal{N}}^T(x(r, s), y(r, s)) \vec{\mathbf{v}}(x(r, s), y(r, s)) \cdot \mathbf{B}(x(r, s), y(r, s)) dr ds \\
&= 2S \iint_{\Delta} \begin{pmatrix} 1-r-s \\ r \\ s \end{pmatrix} \vec{\mathbf{v}}(x(r, s), y(r, s)) \cdot \frac{1}{2S} \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix} dr ds \\
&= \iint_{\Delta} \begin{pmatrix} \mathcal{N}_1(r, s) \\ \mathcal{N}_2(r, s) \\ \mathcal{N}_3(r, s) \end{pmatrix} \vec{\mathbf{v}}(x(r, s), y(r, s)) \cdot \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix} dr ds
\end{aligned}$$

If the velocity is constant within the element (rather rare case) then this can be integrated exactly. If not, a quadrature rule must be used.

Let us assume that indeed velocity is constant inside the element. Then  $\vec{\mathbf{v}}(x(r, s), y(r, s)) = (u_0, v_0)$  and then

$$\begin{aligned}
\mathbf{K}_a &= \iint_{\Delta} \begin{pmatrix} \mathcal{N}_1(r, s) \\ \mathcal{N}_2(r, s) \\ \mathcal{N}_3(r, s) \end{pmatrix} (u_0, v_0) \cdot \begin{pmatrix} y_{23} & y_{31} & y_{12} \\ x_{32} & x_{13} & x_{21} \end{pmatrix} dr ds \\
&= \iint_{\Delta} \begin{pmatrix} \mathcal{N}_1(r, s) \\ \mathcal{N}_2(r, s) \\ \mathcal{N}_3(r, s) \end{pmatrix} \begin{pmatrix} u_0 y_{23} + v_0 x_{32} & u_0 y_{31} + v_0 x_{13} & u_0 y_{12} + v_0 x_{21} \end{pmatrix} dr ds
\end{aligned} \tag{E.73}$$

Using Eqs. (E.69),(E.70),(E.71) we arrive at

$$\mathbf{K}_a = \frac{S}{3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} u_0 y_{23} + v_0 x_{32} & u_0 y_{31} + v_0 x_{13} & u_0 y_{12} + v_0 x_{21} \end{pmatrix} \quad (\text{E.74})$$

$$\mathbf{C}_1 = \int_{\partial\Omega_1} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma$$

The edge  $\partial\Omega_1$  is bounded by the coordinates of nodes  $x_1, y_1$  and  $x_2, y_2$ . This segment can be parameterised by  $t \in [0, 1]$ :

$$\vec{r}(t) = (1-t) \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} + t \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} (x_2 - x_1)t + x_1 \\ (y_2 - y_1)t + y_1 \end{pmatrix}$$

Let us assume that  $C$  is a smooth curve and that it is given by the parametric equations  $x = h(t)$ ,  $y = g(t)$  and  $a \leq t \leq b$ . The line integral of a function  $f(x, y)$  over  $C$  is computed as follows.

$$\int_C f(x, y) ds = \int_a^b f(h(t), g(t)) \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt$$

In our case  $dx/dt = x_2 - x_1$  and  $dy/dt = y_2 - y_1$  so

$$\sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} = L_1$$

Then

$$\begin{aligned} \mathbf{C}_1 &= \int_{\partial\Omega_1} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma \\ &= \begin{pmatrix} \int_{\partial\Omega_1} N_1(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_1(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_1(x, y) N_3(x, y) d\Gamma \\ \int_{\partial\Omega_1} N_2(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_2(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_2(x, y) N_3(x, y) d\Gamma \\ \int_{\partial\Omega_1} N_3(x, y) N_1(x, y) d\Gamma & \int_{\partial\Omega_1} N_3(x, y) N_2(x, y) d\Gamma & \int_{\partial\Omega_1} N_3(x, y) N_3(x, y) d\Gamma \end{pmatrix} \\ &= L_1 \begin{pmatrix} \int_0^1 N_1(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_1(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_1(x(t), y(t)) N_3(x(t), y(t)) dt \\ \int_0^1 N_2(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_2(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_2(x(t), y(t)) N_3(x(t), y(t)) dt \\ \int_0^1 N_3(x(t), y(t)) N_1(x(t), y(t)) dt & \int_0^1 N_3(x(t), y(t)) N_2(x(t), y(t)) dt & \int_0^1 N_3(x(t), y(t)) N_3(x(t), y(t)) dt \end{pmatrix} \end{aligned}$$

We are about to compute the individual terms of the matrix one by one but we will need:

$$\begin{aligned} S &= \frac{1}{2} [(x_1 - x_3)(y_2 - y_3) - (x_2 - x_3)(y_1 - y_3)] \\ &= \frac{1}{2} [x_1 y_2 - x_1 y_3 - x_3 y_2 + x_3 y_3 - x_2 y_1 + x_2 y_3 + x_3 y_1 - x_3 y_3] \\ &= \frac{1}{2} [x_1 y_2 - x_1 y_3 - x_3 y_2 - x_2 y_1 + x_2 y_3 + x_3 y_1] \end{aligned} \tag{E.76}$$

and

$$\begin{aligned}
N_1(x(t), y(t)) &= \frac{1}{2S} [x_2y_3 - x_3y_2 + (y_2 - y_3)x(t) + (x_3 - x_2)y(t)] \\
&= \frac{1}{2S} [x_2y_3 - x_3y_2 + y_{23}(x_{21}t + x_1) + x_{32}(y_{21}t + y_1)] \\
&= \frac{1}{2S} [x_2y_3 - x_3y_2 + y_{23}x_1 + x_{32}y_1 + (y_{23}x_{21} + x_{32}y_{21})t] \\
&= \frac{1}{2S} \underbrace{[x_2y_3 - x_3y_2 + x_1y_2 - x_1y_3 + x_3y_1 - x_2y_1]}_{2S} \\
&\quad + (x_2y_2 - x_1y_2 - x_2y_3 + x_1y_3 + x_3y_2 - x_3y_1 - x_2y_2 + x_2y_1)t] \\
&= \frac{1}{2S} [2S - \underbrace{(x_1y_2 + x_2y_3 - x_1y_3 - x_3y_2 + x_3y_1 - x_2y_1)}_{2S}t] \\
&= \frac{1}{2S} [2S - 2St] \\
&= 1 - t
\end{aligned} \tag{E.77}$$

$$\begin{aligned}
N_2(x(t), y(t)) &= \frac{1}{2S} [x_3y_1 - x_1y_3 + (y_3 - y_1)(x_{21}t + x_1) + (x_1 - x_3)(y_{21}t + y_1)] \\
&= \frac{1}{2S} [x_3y_1 - x_1y_3 + (y_3 - y_1)x_1 + (x_1 - x_3)y_1 + (y_{31}x_{21} + x_{13}y_{21})t] \\
&= \frac{1}{2S} [x_3y_1 - x_1y_3 + x_1y_3 - x_1y_1 + x_1y_1 - x_3y_1 + (y_{31}x_{21} + x_{13}y_{21})t] \\
&= \frac{1}{2S} (y_{31}x_{21} + x_{13}y_{21})t \\
&= t
\end{aligned} \tag{E.78}$$

$$\begin{aligned}
N_3(x(t), y(t)) &= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)x(t) + (x_2 - x_1)y(t)) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)(x_{21}t + x_1) + (x_2 - x_1)(y_{21}t + y_1)) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)x_1 + (x_2 - x_1)y_1 + (y_{12}x_{21} + x_{21}y_{12})t) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + x_1y_1 - x_1y_2 + x_2y_1 - x_1y_1 + (y_{12}x_{21} - x_{21}y_{12})t) \\
&= 0
\end{aligned} \tag{E.79}$$

$$\begin{aligned}
N_3(x(t), y(t)) &= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)x(t) + (x_2 - x_1)y(t)) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)(x_{21}t + x_1) + (x_2 - x_1)(y_{21}t + y_1)) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + (y_1 - y_2)x_1 + (x_2 - x_1)y_1 + (y_{12}x_{21} + x_{21}y_{12})t) \\
&= \frac{1}{2S} (x_1y_2 - x_2y_1 + x_1y_1 - x_1y_2 + x_2y_1 - x_1y_1 + (y_{12}x_{21} - x_{21}y_{12})t) \\
&= 0
\end{aligned} \tag{E.80}$$

then

$$\int_0^1 N_1(x(t), y(t)) N_1(x(t), y(t)) dt = \int_0^1 (1 - t)^2 dt = 1/3 \tag{E.81}$$

$$\int_0^1 N_1(x(t), y(t)) N_2(x(t), y(t)) dt = \int_0^1 (1 - t)t dt = 1/6 \tag{E.82}$$

$$\int_0^1 N_1(x(t), y(t)) N_3(x(t), y(t)) dt = 0 \tag{E.83}$$

$$\int_0^1 N_2(x(t), y(t)) N_2(x(t), y(t)) dt = \int_0^1 t^2 dt = 1/3 \tag{E.84}$$

$$\int_0^1 N_2(x(t), y(t)) N_3(x(t), y(t)) dt = 0 \tag{E.85}$$

$$\int_0^1 N_3(x(t), y(t)) N_3(x(t), y(t)) dt = 0 \tag{E.86}$$

and finally

$$\mathbf{C}_1 = \int_{\partial\Omega_1} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma = \frac{L_1}{6} \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (\text{E.87})$$

$$\mathbf{C}_2 = \int_{\partial\Omega_2} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma = \frac{L_2}{6} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad (\text{E.88})$$

$$\mathbf{C}_3 = \int_{\partial\Omega_3} \vec{N}^T(x, y) \vec{N}(x, y) d\Gamma = \frac{L_3}{6} \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 2 \end{pmatrix} \quad (\text{E.89})$$

# Appendix F

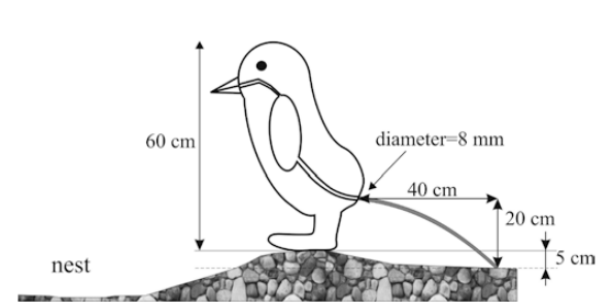
## Finite element terminology in various languages

| English                            | French                                          | Dutch                              |
|------------------------------------|-------------------------------------------------|------------------------------------|
| Finite Element Method              | Méthode des éléments finis                      | Eindige-elementenmethode           |
| Finite Difference Method           | Méthode des différences finies                  | Eindige-differentiemethode         |
| Finite Volume Method               | Méthode des volumes finis                       |                                    |
| Matrix                             | Matrice                                         |                                    |
| Heat transport eq.                 | Equation de transport de la chaleur             | Warmtetransport vergelijking       |
| Momentum conservation eq.          | équation de conservation du moment              | Wet van behoud van impuls          |
| Mass conservation / continuity eq. |                                                 | Continuïteitsvergelijking          |
| Iterative solver                   | solveur itératif                                |                                    |
| Elemental matrix                   |                                                 |                                    |
| Boundary conditions                | conditions aux limites                          | randvoorwaarden                    |
| (In)compressible                   | (in)compressible                                |                                    |
| Surface processes                  | processus de surface                            |                                    |
| an element                         | un élément                                      |                                    |
| Computational geodynamics          | géodynamique numérique                          |                                    |
| Assembly                           | assemblage                                      |                                    |
| Strong form                        |                                                 |                                    |
| Weak form                          | formulation variationnelle / formulation faible |                                    |
| Basis function                     |                                                 |                                    |
| Shape function                     |                                                 |                                    |
| Partial differential eq. (PDE)     | équation aux dérivées partielles (EDP)          | partiële differentiaalvergelijking |
| Node                               | noeud                                           | knooppunt                          |
| Grid, mesh                         | (la) maille / (le) maillage                     | rooster                            |
| Stiffness matrix                   | matrice de raideur                              | stijfheidsmatrix                   |
| Displacement vector                | vecteur déplacement                             | verplaatsingsvector                |
| Tessellation                       | pavage                                          | Betegeling                         |
| Mass matrix                        | matrice de masse                                |                                    |
| Classical mechanics                | mécanique Newtonienne                           | (de) klassieke mechanica           |
| Momentum                           | (le) moment                                     | (de) impuls                        |
| Perimeter                          | (le) perimetre                                  | (de) omtrek                        |
| Wavelength                         | (la) longueur d'onde                            | de golflengte                      |

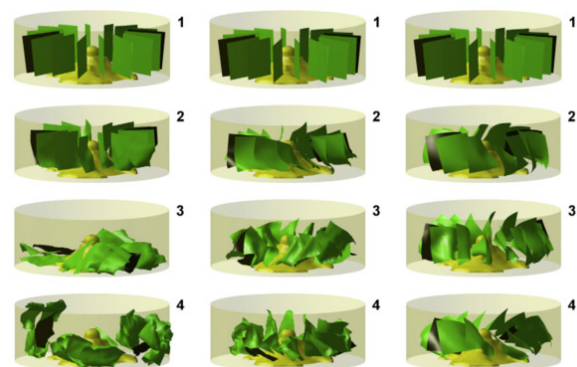
# Appendix G

## Fun modelling

Because sometimes numerical modelling *is* fun ...

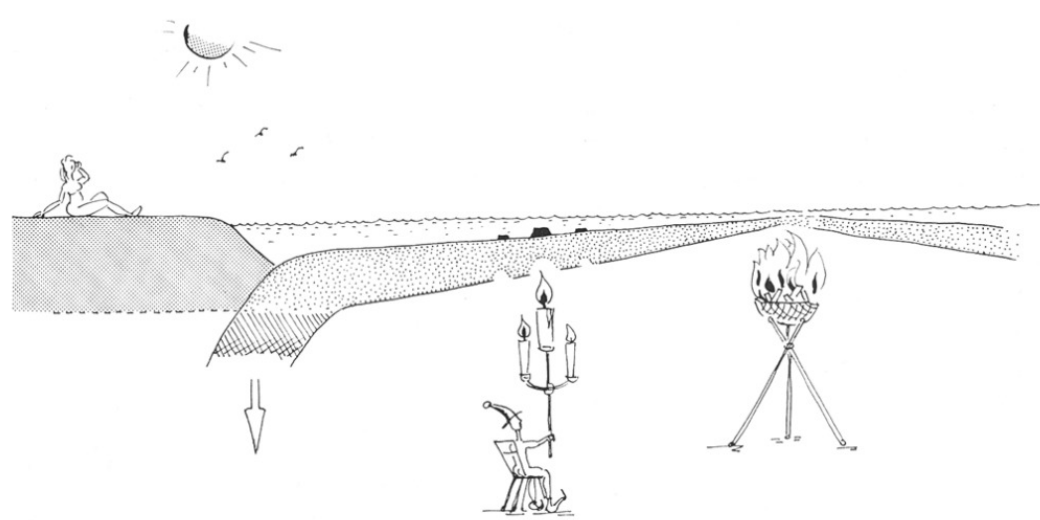


Pressures produced when penguins pooh - calculations on avian defaecation [869]



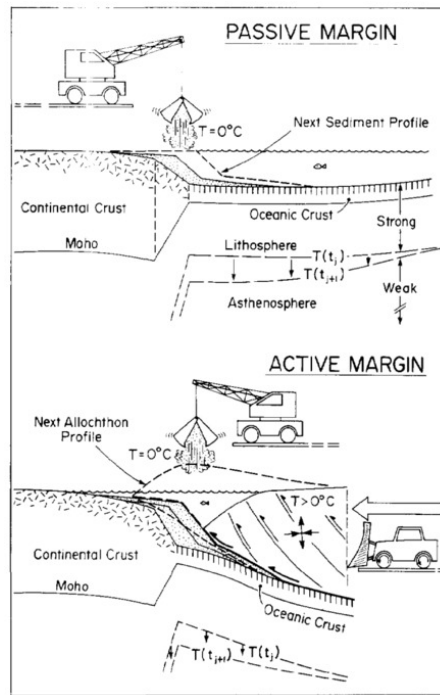
Clothes washing simulations [3]

|                                                                                                                                                                                                                                                                                               |                                                                                                                                                                                                                                                                                          |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><small>Computers and Fluids Vol. 9, pp. 223-253<br/>Pergamon Press Ltd., 1991. Printed in Great Britain</small></p> <p>0045-7950/91/0061-0223\$02.00/0</p> <p><b>DON'T SUPPRESS THE WIGGLES—THEY'RE TELLING YOU SOMETHING!†</b></p> <p>PHILIP M. GRESHO and ROBERT L. LEE</p> <p>[486]</p> | <p>INTERNATIONAL JOURNAL FOR NUMERICAL METHODS IN FLUIDS, VOL. 9, 99–112 (1989)</p> <p><b>ARE FEM SOLUTIONS OF INCOMPRESSIBLE FLOWS REALLY INCOMPRESSIBLE? (OR HOW SIMPLE FLOWS CAN CAUSE HEADACHES!)</b></p> <p>DOMINIQUE PELLETIER, ANDRE FORTIN AND RICARDO CAMARERO</p> <p>[984]</p> |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|



Lithospheric thickness anomaly near the trench and possible driving force of subduction [422]





Convergent margin tectonics [1215]

INTERNATIONAL JOURNAL FOR NUMERICAL METHODS IN ENGINEERING, VOL. 32, 1189–1203 (1991)

## INCOMPRESSIBILITY WITHOUT TEARS—HOW TO AVOID RESTRICTIONS OF MIXED FORMULATION

O. C. ZIENKIEWICZ AND J. WU

*Institute of Numerical Methods in Engineering, University College of Swansea, Swansea, U.K.*

[1435]

## Collective Motion of Humans in Mosh and Circle Pits at Heavy Metal Concerts

Jesse L. Silverberg, <sup>\*</sup> Matthew Bierbaum, James P. Sethna, and Itai Cohen

*Department of Physics and Laboratory of Atomic and Solid State Physics, Cornell University, Ithaca, New York 14853, USA*  
(Received 13 February 2013; published 29 May 2013)

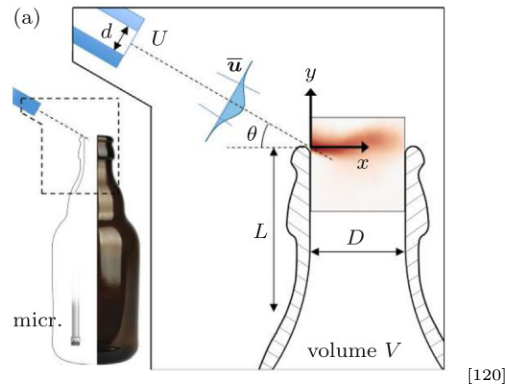
Human collective behavior can vary from calm to panicked depending on social context. Using videos publicly available online, we study the highly energized collective motion of attendees at heavy metal concerts. We find these extreme social gatherings generate similarly extreme behaviors: a disordered gaslike state called a *mosh pit* and an ordered vortexlike state called a *circle pit*. Both phenomena are reproduced in flocking simulations demonstrating that human collective behavior is consistent with the predictions of simplified models.

[1166]

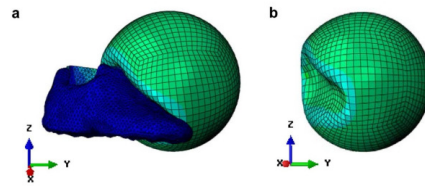
## Toasting the jelly sandwich: The effect of shear heating on lithospheric geotherms and strength

Ebbe H. Hartz *Physics of Geological Processes, University of Oslo, 0316 Oslo, Norway, and Aker Exploration, Haakon VII's gt. 9, P.O. Box 580, Sentrum, NO, 4003 Stavanger, Norway*  
Yuri Y. Podladchikov *Physics of Geological Processes, University of Oslo, 0316 Oslo, Norway*

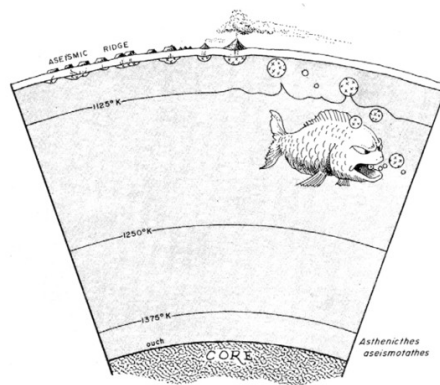
[550]



[120]



Effect of soccer shoe upper on ball behaviour in curve kicks [624]



[1321]

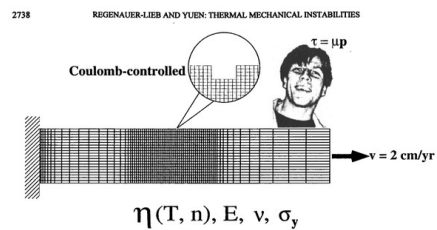
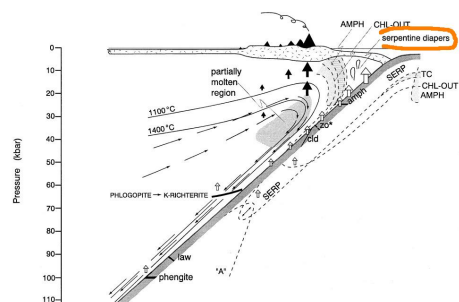


Figure 1. Schematic diagram of a lithosphere cross section with elasto(E)-visco( $\eta$ )-plastic( $\sigma_y$ ) rheology with shear heating. The section has dimensions 800 km on cartesian axis 1, 100 km on cartesian axis 2 (corresponding to depth) and infinite on cartesian axis 3. A small imperfection is present in the center of the plate. No slip is allowed on the left boundary and free slip on the top and bottom boundaries. A constant extension velocity (20 mm/yr) is applied on the right boundary. At  $t=0$  the initial Temperature  $T$  is 987° Kelvin and shear heating is added, conducted and advected during deformation (only source term in the coupled heat equation). Thermal mechanical feedback due to temperature dependent rheology beneath a layer of pressure  $p$  dependent Mohr-Coulomb material (internal friction  $\mu$ , shear stress  $\tau$ ) is analyzed. There is no heat flow on the outer boundaries. The finite element calculations are using a Lagrangian framework. Bobby Poliakov is pictured on the right.

[1053]

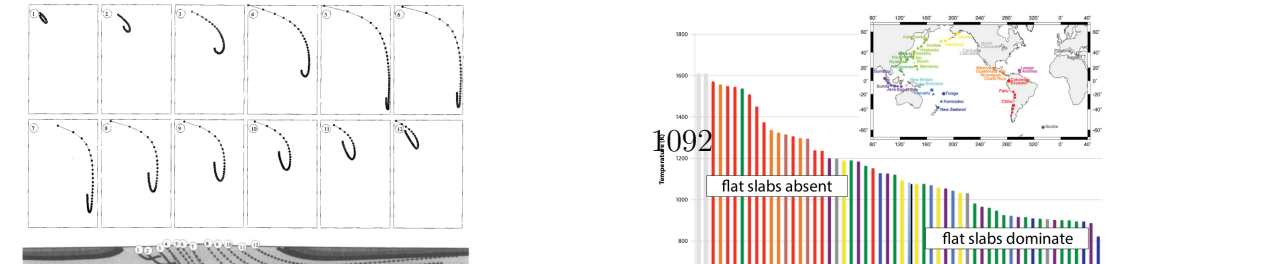
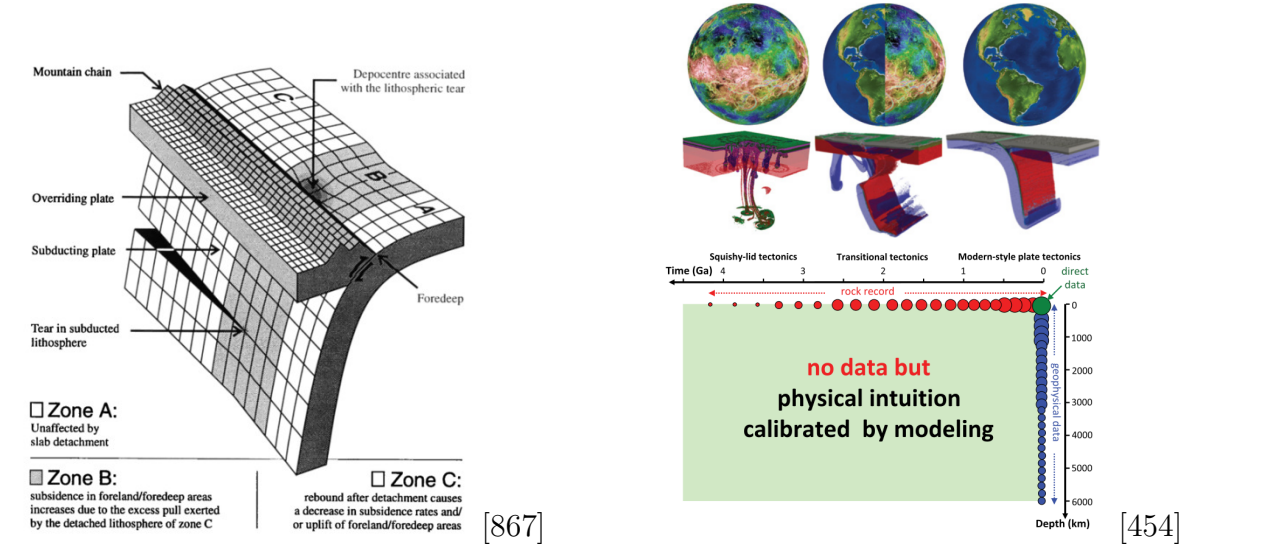
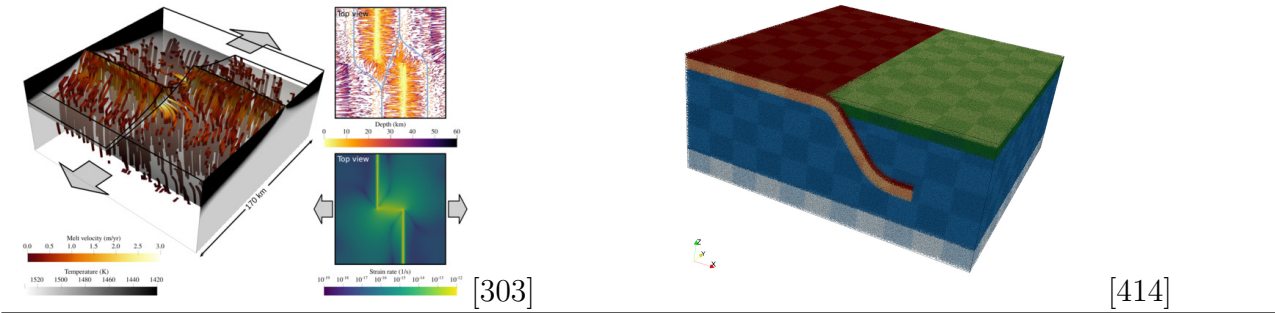
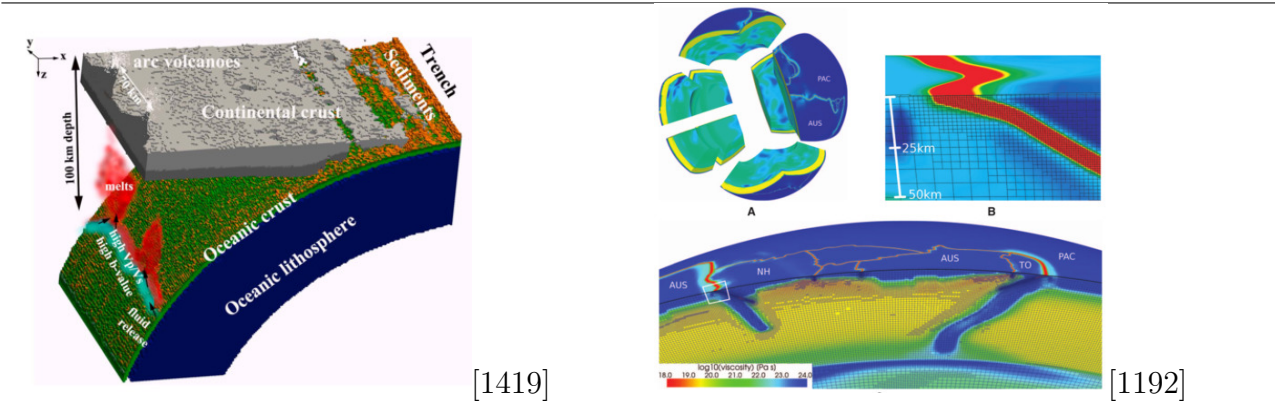


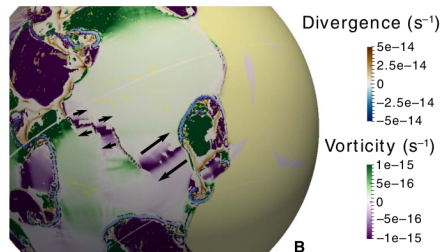
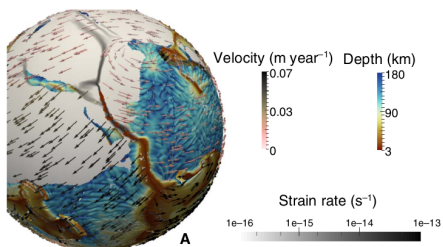
[1129]



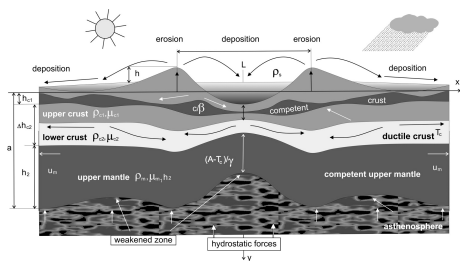
# Appendix H

## Beautiful/interesting images from computational geodynamics

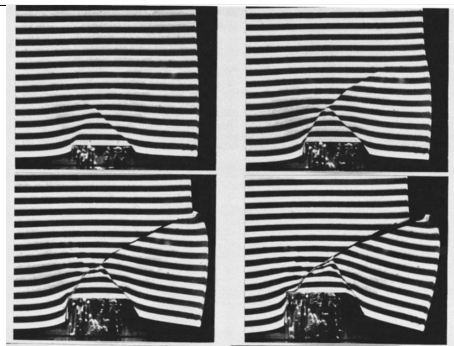




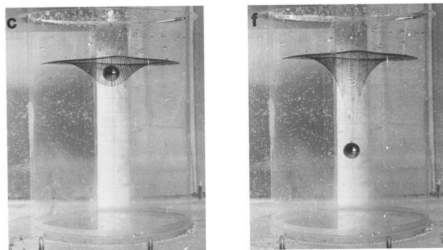
[272]



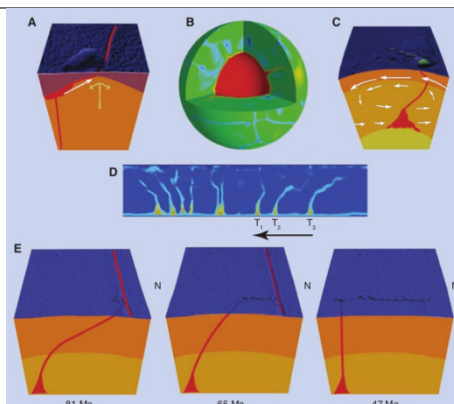
[183]



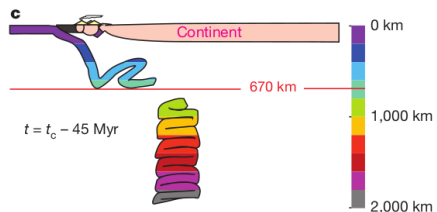
[990]



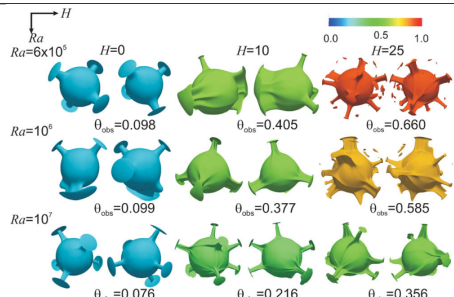
[291]



[1239]



[1164]



[948]





# Appendix I

## Working with Git from the terminal

### I.0.1 Contributing to Aspect

1. make sure that you have an account on GitHub (say, <https://github.com/myusername>) and carry out a proper setup on your local computer as follows:

```
$ git config --global user.name "firstname lastname"
$ git config --global user.email your@email
```

2. On github.com, fork the official ASPECT repository to your repository:  
go to <https://github.com/geodynamics/aspect> and click on the 'fork' button on the upper right corner of the screen.

3. On your machine, in a terminal, clone your repository with  
`$ git clone git@github.com:myusername/aspect.git`. This is now your main<sup>1</sup> branch.

4. Follow the instructions at <https://docs.github.com/en/authentication/connecting-to-github-with-ssh> adding-a-new-ssh-key-to-your-github-account to add a new SSH key to your GitHub account.

5. create a remote of your (online) repository  
`$ git remote add origin git@github.com:myusername/aspect.git` and in order to avoid potential confusion later on, we shall rename our github repo as follows:  
`$ git remote rename origin myusername`

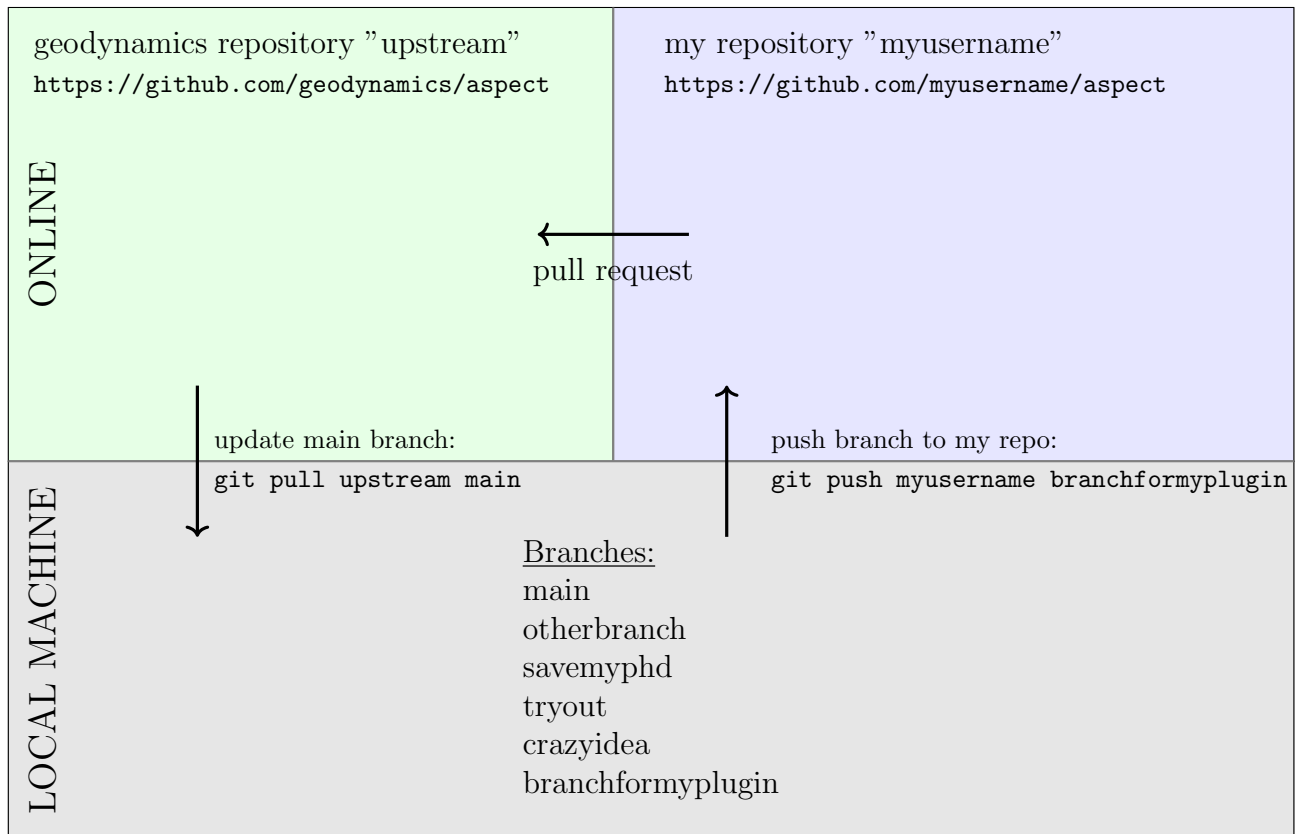
6. Also create a remote of the official version with  
`$ git remote add upstream https://github.com/geodynamics/aspect`

7. Do `$git remote -v` which shows you the URLs that git has stored for the shortname to be used when reading and writing to that remote.

```
$ git remote -v
myusername git@github.com:myusername/aspect.git (fetch)
myusername git@github.com:myusername/aspect.git (push)
upstream https://github.com/geodynamics/aspect.git (fetch)
upstream https://github.com/geodynamics/aspect.git (push)
```

---

<sup>1</sup>previously 'master', although this term should not be used anymore



**Very concretely, if you wish to contribute:** Let us assume that you have found a typo in the file “physics.cc” and you wish to fix this problem.

1. If this is the very first time you use git, the following instruction is not needed. If not, make sure that your terminal points at your main branch: `$ git checkout main`
2. Then, make sure you update your local version: `$ git pull upstream main` and also update your own repo online: `$ git push username main`
3. Create a new branch with a self explanatory name: `$ git branch fix_typo` Then do `$ git branch` to see all your branches. The one that is coloured is the active branch. In order to switch to this new branch, do `$ git checkout fix_typo`. Redo `$ git branch` to verify that the branch 'fix\_typo' is highlighted.
4. Edit the file physics.cc, correct the typo, save and exit. Do `$ git status` and you should see the filename highlighted next to 'modified'.
5. before we go any further, we need to run the indenting script with astyle. In the build directory run `$ make indent`. Note that you need a specific version of astyle, see <https://github.com/geodynamics/aspect/wiki/Indenting-code-and-installing-the-correct-version-of-astyle>
6. add a changelog entry in doc/modules/changes. Look at the entries and find the one that resembles your contribution the most. Use it to write your entry.
7. If this is the only modification you wish to communicate, you then need to add and commit it as follows:

```
$ git add physics.cc
```

If you do `$ git status` again, the file should have changed colour. Then do



```
$ git commit -m "message"
```

where 'message' is a very short description of the modification (e.g. 'I fixed a typo').

8. We now need to bring this modification online

```
$ git push myusername fix_typo
```

9. You are nearly done. The last step takes place on github.com: when you log onto your own github fieldstone page, you should see a large green rectangle "Compare and pull request". Click on this button, follow the instructions.
10. a new page opens, entitled "Open a pull request" (PR). If necessary, add a detailed description of the pull request (this makes sense when you contribute a piece of code, or a whole new section, etc ...). Later you will be able to review the requested changes at the bottom. In the end, simply click "Create pull request".
11. Once you have done so, ASPECT developers will then review it. If they have no comment, the PR will be accepted and your modification will then be incorporated in the main repo. If they have comments, you will be notified via github and a back-and-forth discussion will ensue until the PR is accepted.
12. if the reviewers have comments. Edit/change your file(s). Then add these (`git add ...`), and commit again (`git commit -m "msg"`) and push as before.  
OR: better:  
`git commit --amend -a` so that you don't have to squash/fixup afterwards
13. After the PR has been accepted, the branch is no longer needed. Switch back to your local main branch: `$ git checkout main` . Update your local main and online repo (see step # 2) and then delete the no-longer-needed branch as follows:

```
$ git branch -d fix_typo
```

---

If there is a pb with your PR and u need to rebase. For example if your 2 PRs modify the same line (say for example reference.bib - in that case better spread your changes to different locations in the reference.bib)

- carry out modifications as required by reviewers
- `git add` files, and `git commit`
- `git checkout main`
- `git branch`, make sure you are indeed back on main
- `git pull upstream main`. depending how much happened in the last hours/days it will display a bunch of files/updates
- `git checkout my_branch`
- `git rebase main`. Follow instructions, resolve conflicts in indicated files. `git add problematic_file`. Finish with `git rebase --continue`.

- `git push -f cedriect my_branch`

To configure the editor used by git (do it once):

```
git config --global core.editor "vim"
```

```
git stash git stash pop
```

```
git log
```

how to squash/fixup:

git rebase -i main ->

A similar window as this will open:

```
redhat@batman: ~/ASPECT/pr1/aspect/build
1 pick 6c56e1c54 add lower crustal flow cookbook
2 pick 74e4894ce corrections
3
4 # Rebase d45857ddf..74e4894ce onto d45857ddf (2 commands)
5 #
6 # Commands:
7 # p, pick <commit> = use commit
8 # r, reword <commit> = use commit, but edit the commit message
9 # e, edit <commit> = use commit, but stop for amending
10 # s, squash <commit> = use commit, but meld into previous commit
11 # f, fixup [-C | -c] <commit> = like "squash" but keep only the previous
12 # commit's log message, unless -C is used, in which case
13 # keep only this commit's message; -c is same as -C but
14 # opens the editor
15 # x, exec <command> = run command (the rest of the line) using shell
16 # b, break = stop here (continue rebase later with 'git rebase --continue')
17 # d, drop <commit> = remove commit
18 # l, label <label> = label current HEAD with a name
19 # t, reset <label> = reset HEAD to a label
20 # m, merge [-C <commit> | -c <commit>] <label> [# <oneline>]
21 # create a merge commit using the original merge commit's
22 # message (or the oneline, if no original merge commit was
23 # specified); use -c <commit> to reword the commit message
24 #
25 # These lines can be re-ordered; they are executed from top to bottom.
26 #
27 # If you remove a line here THAT COMMIT WILL BE LOST.
28 #
29 # However, if you remove everything, the rebase will be aborted.
30 #
--
~/ASPECT/pr1/aspect/.git/rebase-merge/git-rebase-todo" 30L, 1358B 2,5 All
```

Since we wish to squash (or rather fixup!) we do:

```
redhat@batman: ~/ASPECT/pr1/aspect/build
1 pick 6c56e1c54 add lower crustal flow cookbook
2 f 74e4894ce corrections
3
4 # Rebase d45857ddf..74e4894ce onto d45857ddf (2 commands)
5 #
6 # Commands:
7 # p, pick <commit> = use commit
8 # r, reword <commit> = use commit, but edit the commit message
9 # e, edit <commit> = use commit, but stop for amending
10 # s, squash <commit> = use commit, but meld into previous commit
11 # f, fixup [-C | -c] <commit> = like "squash" but keep only the previous
12 # commit's log message, unless -C is used, in which case
13 # keep only this commit's message; -c is same as -C but
14 # opens the editor
15 # x, exec <command> = run command (the rest of the line) using shell
16 # b, break = stop here (continue rebase later with 'git rebase --continue')
17 # d, drop <commit> = remove commit
18 # l, label <label> = label current HEAD with a name
19 # t, reset <label> = reset HEAD to a label
20 # m, merge [-C <commit> | -c <commit>] <label> [# <oneline>]
21 # create a merge commit using the original merge commit's
22 # message (or the oneline, if no original merge commit was
23 # specified); use -c <commit> to reword the commit message
24 #
25 # These lines can be re-ordered; they are executed from top to bottom.
26 #
27 # If you remove a line here THAT COMMIT WILL BE LOST.
28 #
29 # However, if you remove everything, the rebase will be aborted.
30 #
-- INSERT --
-- INSERT -- 2,2 All
```

save and exit

do git diff

## I.0.2 Contributing a cookbook in Aspect

1. Download/update/install ASPECT
2. In the cookbooks folder, create a new folder: `$ mkdir my_cool_setup`
3. In this folder, place your .prm file, say my\_cool\_setup.prm
4. Make sure your .prm file is clean, commented, and contains a header with a concise description of what the experiment is, and/or in which publication it originates.
5. In the folder, create a new folder: `$ mkdir doc`
6. In this folder create the file my\_cool\_setup.md file which contains the text for the cookbook. Look at other cookbooks md files for examples of how to include a figure, an equation, cite publications, etc ...
7. Place in this same folder all figures pertaining to the cookbook entry in the manual.
8. go to `/doc/sphinx/user/cookbooks/` and add your cookbook to the list in (for example) `geophysical-setups.md`
9. if you wish to cite publications, add them to `/doc/sphinx/references.bib`
10. In order to generate the manual, go to `/doc/sphinx` and do `$ make html`
11. if you wish to re-generate the part of the manual that comes from the documentation of .cc files, then go to `\build` and make the code, then do in `/doc/`:  
  
`./update_parameters.sh /home/absolute/path/aspect/build/aspect`  
  
and then do `make html`. If you re-modify the .cc file, you need to redo all 3 steps.
12. If there is no error, you should be able to open the file `/doc/sphinx/\_build/html/index.html` with firefox
13. if the referencing of the figures does not work correctly, simply do `make clean` and then make `html` again.
14. before you make a pull request, make sure you run `make indent`

### I.0.3 Contributing to fieldstone

This appendix was contributed by E. van der Wiel.

1. make sure that you have an account on GitHub (say, <https://github.com/myusername>) and carry out a proper setup on your local computer as follows:  

```
$ git config --global user.name "firstname lastname"
$ git config --global user.email your@email
```
2. On github.com, fork the official fieldstone repository to your repository:  
go to <https://github.com/cedrict/fieldstone> and click on the 'fork' button on the upper right corner of the screen.
3. On your machine, in a terminal, clone your repository with  

```
$ git clone git@github.com:myusername/fieldstone.git
```

. This is now your main<sup>2</sup> branch.
4. On your machine, find your security key <sup>3</sup> with 

```
$ less ~/.ssh/id_dsa.pub
```

 and copy this into github.com so you can push to your repository. See <https://help.github.com/en/articles/connecting-to-github-with-ssh> on how to configure github with ssh support (no more login/password to type – if you have cloned the repository with ssh, not html). Please also check the instructions at <https://help.github.com/en/articles/connecting-to-github-with-ssh>.
5. create a remote of your (online) repository  

```
git remote add origin git@github.com:myusername/fieldstone.git
```

 and in order to avoid potential confusion later on, we shall rename our github repo as follows:  

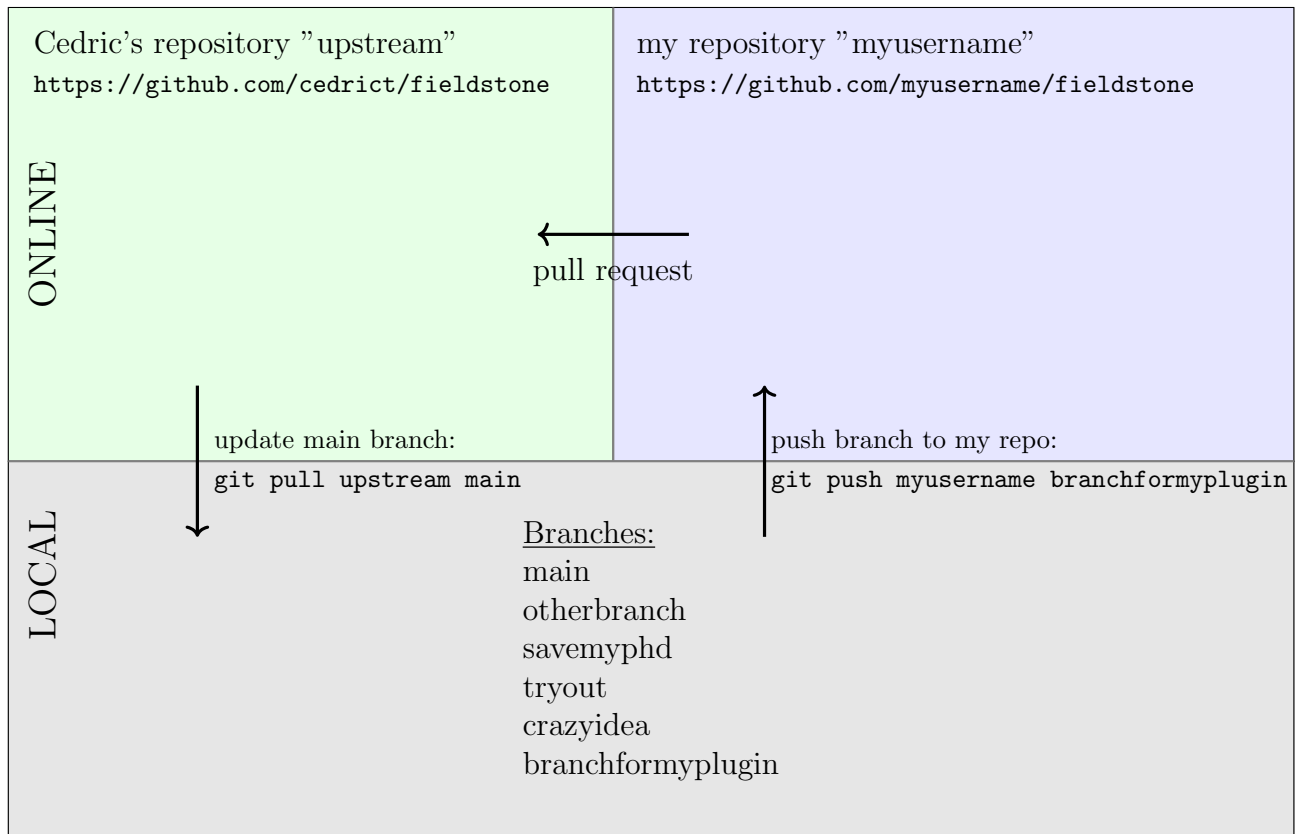
```
$ git remote rename origin myusername
```
6. Also create a remote of the official version with  

```
$ git remote add upstream https://github.com/cedrict/fieldstone
```

---

<sup>2</sup>also 'master', although this term should not be used anymore

<sup>3</sup>Note that you can create a public key as follows: <https://help.github.com/articles/generating-ssh-keys/>



**Very concretely, if you wish to contribute:** Let us assume that you have found a typo in the file "physics.tex" and you wish to report the problem to me. The easiest way is to write me an email with the nature of the problem and the proposed fix. Although I will be grateful for your contribution, this approach can be improved upon by using git's "pull request" command.

1. If this is the very first time you use git, the following instruction is not needed. If not, make sure that your terminal points at your main branch: `$ git checkout master`
2. Then, make sure you update your local version: `$ git pull upstream master` and also update your own repo online: `$ git push`
3. Create a new branch with a self explanatory name: `$ git branch fix_typo` Then do `$ git branch` to see all your branches. The one that is coloured is the active branch. In order to switch to this new branch, do `$ git checkout fix_typo`. Redo `$ git branch` to verify that the branch 'fix\_typo' is highlighted.
4. Edit the file and correct the typo, save and exit. Do `$ git status` and you should see the filename highlighted next to 'modified'.
5. If this is the only modification you wish to communicate, you then need to add and commit it as follows:

```
$ git add physics.tex
```

If you do git status again, the file should have changed colour (?). Then do

```
$ git commit -m "message"
```

where 'message' is a very short description of the modification (e.g. 'I fixed a typo').

6. We now need to bring this modification online

```
$ git push myusername fix_typo
```

7. You are nearly done. The last step takes place on github.com: when you log onto your own github fieldstone page, you should see a large green rectangle "Compare and pull request". Click on this button.
8. a new page opens, entitled "Open a pull request". If necessary, add a detailed description of the pull request (this makes sense when you contribute a piece of code, or a whole new section, etc ...). You can review the requested changes at the bottom. In the end, simply click "Create pull request".
9. Once you have done so, I will receive an email which notifies me of the pull request. I will then review it. If I have no comment, I will accept the PR and your modification will then be incorporated in the master repo. If I have comments, you will be notified via github and a back-and-forth discussion will ensue until I accept the PR.
10. After the PR has been accepted, the branch is no longer needed. Switch back to your local master branch: `$ git checkout master` . Update your local master and online repo (see step # 2) and then delete the no-longer-needed branch as follows:

```
$ git branch -d fix_typo
```

```

(base) UU062931:fieldstone vanderWiel$ git branch
 master
* my_work
(base) UU062931:fieldstone vanderWiel$ git checkout master
Switched to branch 'master'
Your branch is up to date with 'origin/master'.
(base) UU062931:fieldstone vanderWiel$ git pull upstream master
From https://github.com/cedrict/fieldstone
 * branch master -> FETCH_HEAD
Already up to date.
(base) UU062931:fieldstone vanderWiel$ git push
Everything up-to-date
(base) UU062931:fieldstone vanderWiel$

(base) UU062931:fieldstone vanderWiel$ git branch fix_typo
(base) UU062931:fieldstone vanderWiel$ git branch
 fix_typo
* master
 my_work
(base) UU062931:fieldstone vanderWiel$ git checkout fix_typo
Switched to branch 'fix_typo'
(base) UU062931:fieldstone vanderWiel$ git branch
* fix_typo
 master
 my_work
(base) UU062931:fieldstone vanderWiel$

(base) UU062931:fieldstone vanderWiel$ vi physics.tex
(base) UU062931:fieldstone vanderWiel$ git status
On branch fix_typo
Changes not staged for commit:
 (use "git add <file>..." to update what will be committed)
 (use "git restore <file>..." to discard changes in working directory)
 modified: physics.tex

no changes added to commit (use "git add" and/or "git commit -a")
(base) UU062931:fieldstone vanderWiel$ git add physics.tex
(base) UU062931:fieldstone vanderWiel$ git status
On branch fix_typo
Changes to be committed:
 (use "git restore --staged <file>..." to unstage)
 modified: physics.tex

(base) UU062931:fieldstone vanderWiel$ git commit -m "I fixed a typo"
[fix_typo 85beb3e] I fixed a typo
1 file changed, 1 insertion(+), 1 deletion(-)
(base) UU062931:fieldstone vanderWiel$ git push origin fix_typo
Enumerating objects: 5, done.
Counting objects: 100% (5/5), done.
Delta compression using up to 4 threads
Compressing objects: 100% (3/3), done.
Writing objects: 100% (3/3), 294 bytes | 294.00 KiB/s, done.
Total 3 (delta 2), reused 0 (delta 0)
remote: Resolving deltas: 100% (2/2), completed with 2 local objects.
remote:
remote: Create a pull request for 'fix_typo' on GitHub by visiting:
remote: https://github.com/e-wiel/fieldstone/pull/new/fix_typo
remote:
To https://github.com/e-wiel/fieldstone.git
 * [new branch] fix_typo -> fix_typo
(base) UU062931:fieldstone vanderWiel$

235 commits 3 branches 0 releases 3 contributors View license

Your recently pushed branches:
 fix_typo (1 minute ago) [Compare & pull request]

Branch: master New pull request Create new file Upload files Find file Clone or download

This branch is even with cedrict:master. [Pull request] [Compare]
cedrict update:h Latest commit e6faecc 41 minutes ago
images bibi243 3 hours ago

(base) UU062931:fieldstone vanderWiel$ git branch
 fix_typo
* master
 my_work
(base) UU062931:fieldstone vanderWiel$ git pull upstream master
remote: Enumerating objects: 6, done.
remote: Counting objects: 100% (6/6), done.
remote: Compressing objects: 100% (3/3), done.
remote: Total 6 (delta 3), reused 4 (delta 3), pack-reused 0
Unpacking objects: 100% (6/6), done.
From https://github.com/cedrict/fieldstone
 * branch master -> FETCH_HEAD
 e6faecc..cb3a879 master -> upstream/master
Updating e6faecc..cb3a879
Fast-forward
 github.tex | 147 ++++++
 manual.pdf | Bin 65759829 -> 65769253 bytes
 physics.tex | 2 +-
 3 files changed, 143 insertions(+), 6 deletions(-)
(base) UU062931:fieldstone vanderWiel$ git branch -d fix_typo
Deleted branch fix_typo (was 85beb3e).
(base) UU062931:fieldstone vanderWiel$ git push
Total 0 (delta 0), reused 0 (delta 0)
To https://github.com/e-wiel/fieldstone.git
 e6faecc..cb3a879 master -> master
(base) UU062931:fieldstone vanderWiel$

```

Screen captures of the procedure described above as carried out by E. van der Wiel on his Apple laptop.



In what follows we summarize the most important commands one should and remember while working with github. After creating an account one can 'fork' a repository (repo) in the online environment. This repository is a copy from the master directory of the developer and should not be used to adapt or change, as changes from the developer (updates) should be obtained in this 'fork', or as it could also be called; your master branch.

In order to be able to work within a repository, for instance, to run and compile different programs, you should have you own branch of the repository in which YOU CAN make changes. The following commands should be used to make, copy and publish your own version of the repo to your local device and the online github environment.

| command                         | what it does                                                                                                                        |
|---------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|
| git branch                      | shows all branches of your repository and highlights the one you're in.                                                             |
| git checkout -b <my_own_branch> | This makes your own branch called "my_own_branch".                                                                                  |
| git push origin <my_own_branch> | This pushes your own, local, branch to as a second branch in the online repo of github.                                             |
| git checkout <name>             | changing the branch your working in (e.g. master or my_own_branch). Or replace the name with a hyphen to switch to the last branch. |
| git branch -d <my_own_branch>   | Delete your local branch.                                                                                                           |

The following commands should be used in order to update your own local branches from updates made by somewhere else (upstream/master is the main repository). One should do this for the local master branch and, where possible as well for the different local branches you have committed changes to already.

| command                      | what it does                                                                                                                                                                                                                    |
|------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| git checkout master          | To make sure you are in the right branch                                                                                                                                                                                        |
| git fetch upstream           | to fetch updates from upstream repositories to you own local branch (e.g. to update your master branch).                                                                                                                        |
| git merge upstream/master    | Command to update the branch with the fetched repo from 'upstream'.                                                                                                                                                             |
| git push origin master       | To level your own online repository again with the one on your local drive (and thus the one upstream).                                                                                                                         |
| git checkout <my_own_branch> | To switch to your own adapted branch of the repo.                                                                                                                                                                               |
| git merge master             | Used from another branch working directory to combine the new released version of the master repo with the one where all your own changes are put. ->Then git finds all conflicts in different files which you need to resolve. |
| git add .                    | This adds the resolved issues in your own local branch (not master). After which you are able to commit and push your changes back to the online repository.                                                                    |

While you are working in your own branch you can change, add or delete files in any amount you want. However, always check whether your changes do not inflict the outcome of for instance your code. And when uploading from your terminal: if you commit and then push from master branch

your changes will automatically be inserted in the online version of your master branch, when done from another branch it will be shown as a pull request towards your master branch. This request can than, for instance be forwarded to the main repo.

| <b>command</b>            | <b>what it does</b>                                                                                                                                                                                                                                                                           |
|---------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| git commit -a             | This will send your changes/updates from your branch as a commit to your own local branch.                                                                                                                                                                                                    |
| git push origin <changes> | To update the remote repository (on Github) from you local repository (in this case the 'changes' branch). (Actually upload the new version). On-line one can then judge what to do with it. !! this is a pull request towards your own fork/local_branch                                     |
| git status                | Showing the status of your current branch; it shows which files are different between the master file and your adapted branch.                                                                                                                                                                |
| git diff <changes>        | This shows the exact differences between the different branches; one can simply ask for the difference between two branches when pwd in one branch ask for the other branch.                                                                                                                  |
| git merge <my_own_branch> | When used from the master branch (or any other???) this accepts the changes made in your branch and puts them in your local(!) master branch.                                                                                                                                                 |
| git pull origin master    | if the main repository changes, one can pull the newest version towards it's own master file. While keeping your own branches alive with you own changes and vica versa: by running this command the origin/master (remote file) will be cloned and updated to the working branch you are in. |
| git stash (apply)         | ?? While updating your local branch, sometimes git wants to overrule your own changes, with this command you can 'stash' them to look at the differences later. ??                                                                                                                            |

# Appendix J

## Writing a report as homework

app-grading.tex

- The document should contain your full name and student number on the first page.
- The file should be a pdf which name contains your family name
- Layout: is the document visually pleasing? Is it well structured?
- Is there a complete bibliography (when applicable)?
- Does the structure follows this: Introduction - Methods - Results - Discussion - Conclusion - Appendix ?
- Figures: Are they properly numbered? captioned? all figures must be referenced in the text. Are they of good enough quality (no visible pixels)? are they readable? are all axis labelled?
- Text: Overall quality of the language. Are there still typos? Do all sentence make sense?
- if you wish to show lines of code, use verbatim or lstlisting<sup>1</sup>
- Discussion: are the results properly discussed, analyzed? are potential problems, errors, limitations discussed?
- Conclusion: Are the findings/results summarized and generalized?

| No                | Yes                     |
|-------------------|-------------------------|
| $6.67 * 10^{-11}$ | $6.67 \times 10^{-11}$  |
| $kg/m^3$          | $kg/m^3$ or $kg.m^{-3}$ |
| 1x1               | $1 \times 1$            |
| $\cos$            | cos                     |
| docx file         | pdf file                |
| 'if you do this'  | passive form            |



No grey background

<sup>1</sup>[https://en.wikibooks.org/wiki/LaTeX/Source\\_Code\\_Listings](https://en.wikibooks.org/wiki/LaTeX/Source_Code_Listings)

---

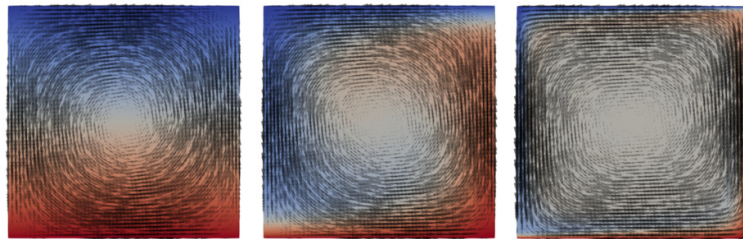
```

U = [-41346074.59639838, -41346106.14861266, -41346322.12423281, -41347034.8219668,
-41348728.85336469, -41352079.33063214, -41357972.05539607, -41367516.98181634
-41382034.29417382, -41402966.8138888, -41431618.19900567, -41468508.42097412
-41511962.23285898, -41555426.6099448, -41583595.20822245, -41570144.20903767
-41484197.88710789, -41308649.41597945, -41054103.1845507, -40748933.68245689
-40415944.45883375, -40057457.53015289, -39655212.74099992, -39185002.81130297
-38639470.67788275, -38036907.55824716, -37405072.1363776, -36759267.24928432
-36092497.74498677, -35380493.83127169, -34600914.0153607, -33754968.22552883
-32868819.25937854, -31975267.12561081, -31098665.54160668, -30252235.42279971
-29441444.64112505, -28667591.04175556, -27930018.48500311, -27227239.55989363
-26557456.02105067, -25918791.91839744, -25309397.59712495, -24727496.32338618
-24171404.76853348, -23639540.9187365, -23130425.50643475, -22642679.89858946
-22175021.99222236, -21726261.0159415, -21295291.79559261, -20881088.84329066
-20482700.50219241, -20099243.29398487, -19729896.55730399, -19373897.42481487
-19030536.15912975, -18699151.84940752, -18379128.45870995, -18069891.20501603
-17770903.2548138, -17481662.7063611, -17201699.83931143, -16930574.60791429
-16667874.35607352, -16413211.73392739, -16166222.7971429, -15926565.27168286
-15693916.96834184, -15467974.33280896, -15248451.11838825, -15035077.16976939
-14827597.30739603, -14625770.30303001, -14429367.93805357, -14238174.13690581
-14051984.16881316, -13870603.91165919, -13693849.17245316, -13521545.05940309
-13353525.40109213, -13189632.20869422, -13029715.17755753, -12873631.22483632
-12721244.06016446, -12572423.78664866, -12427046.52970981, -12284994.09152793

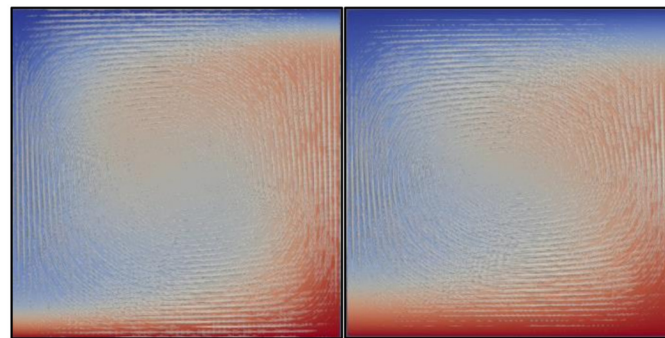
```

No lists/arrays with numbers

---

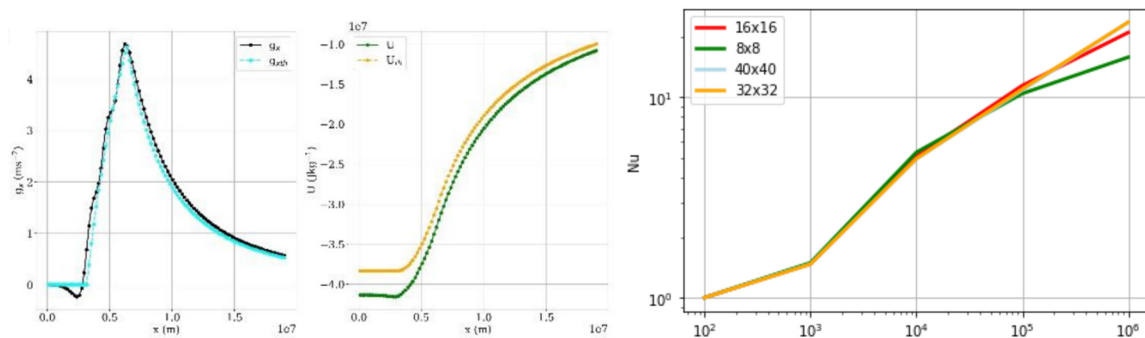


Too many arrows



Poor choice of arrow colour

---



Be careful about how you export your figures. These are unreadable.

$$vrms = (\frac{1}{V} \int_{\Omega} |\mathbf{v}|^2 dV)^{\frac{1}{2}}$$

$$= (\frac{1}{V_{\Omega}} \int_{\Omega} (u^2 + v^2) dV)^{\frac{1}{2}} \quad (9)$$

Paranthesis too small

Figure 1, steady state for the purpose of this work is to be considered within the interval corresponding to Ra values between  $1E+04$  and  $1E+05$ . Hence, Nu values for a steady state range from  $4.92E+00$  to  $1.14E+01$  among all tested resolutions. For the model's standard  $32 \times 32$  resolution, the Nu value rises from  $4.94E+00$  to  $1.10E+01$  within the Ra interval of  $1E+04$  to  $1E+05$ . The critical Ra number ( $Ra_c$ ) for all four resolutions seems to coincide at the value of  $5E+03$ .

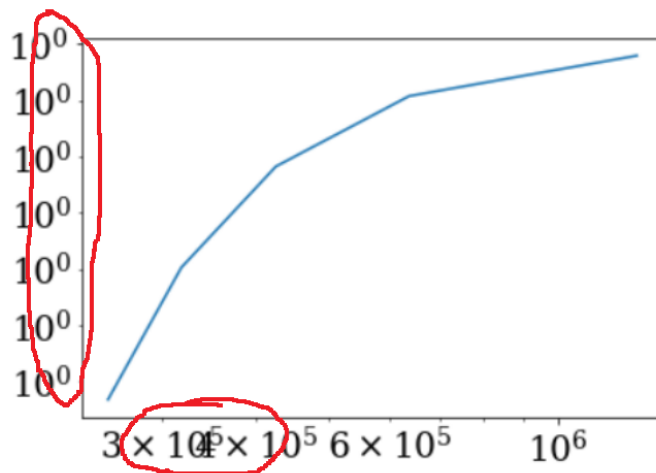
$1.6E+10$  is not acceptable. Replace by  $1.6 \cdot 10^{10}$

$$\nabla \cdot [\eta(\nabla \mathbf{v} + \nabla \mathbf{v}^T)] - \nabla p + \rho(T)\mathbf{g} = \mathbf{0} \quad (5)$$

$$\nabla \cdot \mathbf{v} = 0 \quad (6)$$

$$\rho_0 C_p \mathbf{v} \cdot \nabla T = k \Delta T \quad (7)$$

Equation number is too close to the equation itself. Use labels, do not number equations by hand.



Formatting of both axis lead to unreadable figure.

$$g_x(x,y,z) = \mathcal{G} \sum_{e=1}^{N_e} \rho_e V_e \frac{x-x_e}{|\vec{r}-\vec{r}_e|^3}$$

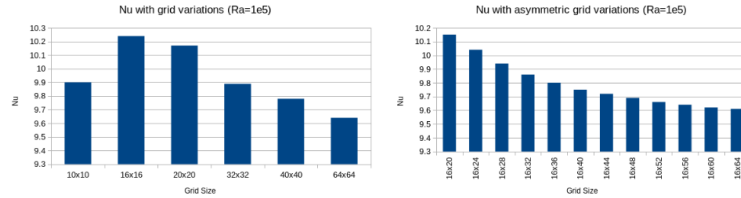
In L<sup>A</sup>T<sub>E</sub>X use `\sum\limits`

The computed mass, using the PREM density distribution, the mass of the Earth, and the discrepancy between the two are  $5.967863262629439e+24$  kg.

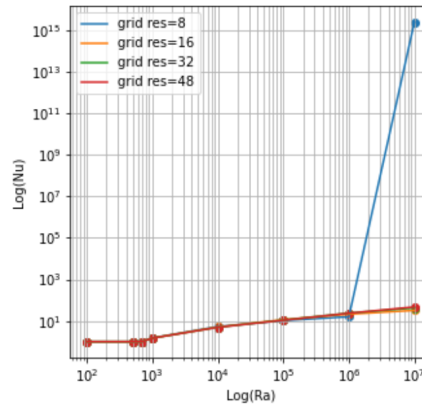
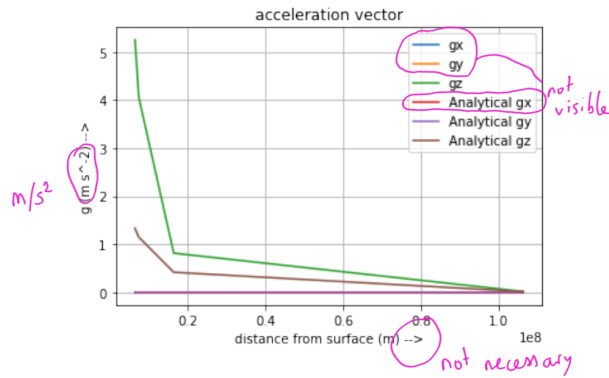
Are so many digits necessary?

With this new density a new mass for the sphere can be calculated. The mass calculated with the PREM density is  $5.138471643448518 \times 10^{24} \text{ kg}$ , the mass of the Earth is  $5.972 \times 10^{24} \text{ kg}$ , the difference between the calculates mass of the sphere and the mass of the Earth is  $8.33528365514822 \times 10^{23} \text{ kg}$ . From the mass calculated with PREM density the gravitational acceleration at the surface can also be calculated and compared between the sphere and the Earth.  $g_{\text{sphere}} = 8.38566314096455 \text{ m/s}^2$ ,  $g_{\text{Earth}} = 9.78291192747942 \text{ m/s}^2$ ,  $g_{\text{difference}} = 1.39724878651487 \text{ m/s}^2$ .

use `\usepackage[cm]{fullpage}` to allow for wider text.



This figure style is to be avoided. Simply use dots and/or lines.



Here Ra and Nu are plotted in log-log scale, not log(Ra) and log(Nu).

The analytical solution for the gravitational field of a **solid sphere** with radius  $R$  is as follows:

$$g(r)_{\text{int}} = -\frac{4\pi}{3} G \rho_0 r; \quad (1)$$

$$U(r)_{\text{int}} = \frac{2\pi}{3} G \rho_0 r^2 + \frac{3}{2} \frac{GM}{r} \quad (2)$$

... for the sphere's interior ( $r \leq R$ ), and ...

$$g(r)_{\text{ext}} = \frac{MG}{r^2}; \quad (3)$$

$$U(r)_{\text{ext}} = -\frac{GM}{r} \quad (4)$$

... for the sphere's exterior ( $r \geq R$ ), where  $g(r)$  is the gravitational acceleration vector norm and  $U(r)$  the gravita-

The dots at the beginning and end of the lines are not necessary.

---

object and the moment of inertia depends quadratically on the size. This makes the error for calculating the moment of inertia slightly higher. The trend that each three graphs follow is slightly unexpected because you would expect the relative error to keep decreasing at higher values of  $N$ , but at a couple of places in the data we see an increase in the

Never *never* use 'you'.

---

For both the full and hollow spheres, a domain of size  $2R \times 2R \times 2R$  centered the radius of the sphere. It was partitioned in  $N \times N \times N$  cells, where  $N$  was

Do not use 'X' but  $\times$  (`\times`).

## J.0.1 Computational Geodynamics Report

All the comments above apply, with additional instructions:

- report should be in  $\text{\LaTeX}$
- The document should contain your full name and student number on the first page.
- The report file should be a pdf which name contains your family name
- not more than 25 pages. If more, use appendices wisely.
- document should be structured in two main parts: FDM and FEM.
- no equations unless necessary to the discussion (still mention the equation that you are solving but refer to an external document/article/book for example).
- use `lstlisting` package to include code
- use `\usepackage[cm]{fullpage}` to format your document
- all codes either in appendix or in zip file (bearing your name).
- a decent introduction (half page to one page) which links the topic of this course to geosciences.
- discussion of results (stability, convergence, influence of resolution, remarks of all kinds).
- if you did not succeed in doing a particular exercise, please explain what you think the problem is, how you know it is not working, etc ...
- think about colormaps, image compression
- DEADLINE: July 16th, 2023, 23:59

I will use this table to grade your reports:

|                                  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|----------------------------------|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
|                                  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| title, names, student nb         |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| L <sup>A</sup> T <sub>E</sub> X? |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| document layout                  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| equations look                   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| equations numbered               |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| use of equations                 |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| figs: caption                    |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| figs: pixels?                    |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| figs: correct?                   |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| English grammar                  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| Typos Introduction               |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| methods/results                  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| Discussion/Conclusion            |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| Extra work?                      |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| Bibliography                     |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| code layout                      |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| code style                       |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| code accuracy                    |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |



# Appendix K

## Analytical expressions for $\mathbb{G}_{el}$

The elemental matrix  $\mathbb{G}_{el}$  is given by (see Section 7.5):

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_x N_1^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_1^\gamma & \partial_x N_1^\gamma \\ \partial_x N_2^\gamma & 0 & \partial_y N_2^\gamma \\ 0 & \partial_y N_2^\gamma & \partial_x N_2^\gamma \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \partial_x N_{m_\gamma}^\gamma & 0 & \partial_y N_{m_\gamma}^\gamma \\ 0 & \partial_y N_{m_\gamma}^\gamma & \partial_x N_{m_\gamma}^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p & N_2^p & \cdots & N_{m_p}^p \\ N_1^p & N_2^p & \cdots & N_{m_p}^p \\ 0 & 0 & \cdots & 0 \end{pmatrix} d\Omega$$

In what follows I set out to compute this elemental matrix for the reference element. All of the integrals were computed with WolframAlpha since it allowed me to copy-paste the L<sup>A</sup>T<sub>E</sub>Xcode directly into the website prompt area and obtain the value of these integrals.

### K.0.1 $Q_1 \times P_0$ element - 2D

For this element,  $m_\gamma = 4$  and  $m_p = 1$  so  $\mathbb{G}_{el}$  is a  $8 \times 1$  matrix:

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\gamma & 0 & \partial_s N_1^\gamma \\ 0 & \partial_s N_1^\gamma & \partial_r N_1^\gamma \\ \partial_r N_2^\gamma & 0 & \partial_s N_2^\gamma \\ 0 & \partial_s N_2^\gamma & \partial_r N_2^\gamma \\ \partial_r N_3^\gamma & 0 & \partial_s N_3^\gamma \\ 0 & \partial_s N_3^\gamma & \partial_r N_3^\gamma \\ \partial_r N_4^\gamma & 0 & \partial_s N_4^\gamma \\ 0 & \partial_s N_4^\gamma & \partial_r N_4^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p \\ N_1^p \\ 0 \end{pmatrix} d\Omega$$

also, since  $N_1^p = 1$  then

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\gamma \\ \partial_s N_1^\gamma \\ \partial_r N_2^\gamma \\ \partial_s N_2^\gamma \\ \partial_r N_3^\gamma \\ \partial_s N_3^\gamma \\ \partial_r N_4^\gamma \\ \partial_s N_4^\gamma \end{pmatrix} d\Omega = \begin{pmatrix} 1 \\ 1 \\ -1 \\ 1 \\ -1 \\ -1 \\ 1 \\ -1 \end{pmatrix}$$

A macro element made of a single element makes no sense since if velocity is prescribed on all sides there is not a single velocity dof left.

We then consider the following macro-element:

| velocity    | pressure    |      |
|-------------|-------------|------|
| 7====8====9 | .====.====. |      |
|             | 3     4     |      |
| 4====5====6 | .====.====. | NV=9 |
|             | 1     2     |      |
| 1====2====3 | .====.====. | NP=4 |

The assembled  $\mathbb{G}$  matrix is then  $18 \times 4$ :

$$\mathbb{G} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ -1 & 0 & 1 & 0 \\ -1 & 1 & -1 & 1 \\ -1 & -1 & 1 & 1 \\ 0 & -1 & 0 & -1 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & -1 & 1 \\ 0 & 0 & -1 & -1 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & -1 \end{pmatrix}$$

After applying boundary conditions on nodes 1,2,3,4,6,7,8,9:

$$\mathbb{G} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 1 & -1 & 1 \\ -1 & -1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

or,

$$\tilde{\mathbb{G}} = \begin{pmatrix} -1 & 1 & -1 & 1 \\ -1 & -1 & 1 & 1 \end{pmatrix}$$

The null space is of size two, spawn by the two vectors:

```
[[-0.1 0.7]
 [0.7 0.1]
 [0.7 0.1]
 [-0.1 0.7]]
```

In the book the authors proceed to show that any such macroelement made of rectangles has a spurious mode.

See code `python_codes/Gel/macro_element_q1p0.py`

### K.0.2 $Q_1 \times P_0$ element - 3D

For this element,  $m_v = 8$  and  $m_p = 1$  so  $\mathbb{G}_{el}$  is a  $3 * 8 \times 1$  matrix:

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^v & 0 & 0 & \partial_s N_1^v & \partial_t N_1^v & 0 \\ 0 & \partial_s N_1^v & 0 & \partial_r N_1^v & 0 & \partial_t N_1^v \\ 0 & 0 & \partial_t N_1^v & 0 & \partial_r N_1^v & \partial_s N_1^v \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \partial_r N_8^v & 0 & 0 & \partial_s N_8^v & \partial_t N_8^v & 0 \\ 0 & \partial_s N_8^v & 0 & \partial_r N_8^v & 0 & \partial_t N_8^v \\ 0 & 0 & \partial_t N_8^v & 0 & \partial_r N_8^v & \partial_s N_8^v \end{pmatrix} \cdot \begin{pmatrix} N_1^p \\ N_1^p \\ N_1^p \\ 0 \\ 0 \\ 0 \end{pmatrix} d\Omega$$

also, since  $N_1^p = 1$  then

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\vee \\ \partial_s N_1^\vee \\ \partial_t N_1^\vee \\ \partial_r N_2^\vee \\ \partial_s N_2^\vee \\ \partial_t N_2^\vee \\ \partial_r N_3^\vee \\ \partial_s N_3^\vee \\ \partial_t N_3^\vee \\ \partial_r N_4^\vee \\ \partial_s N_4^\vee \\ \partial_t N_4^\vee \\ \partial_r N_5^\vee \\ \partial_s N_5^\vee \\ \partial_t N_5^\vee \\ \partial_r N_6^\vee \\ \partial_s N_6^\vee \\ \partial_t N_6^\vee \\ \partial_r N_7^\vee \\ \partial_s N_7^\vee \\ \partial_t N_7^\vee \\ \partial_r N_8^\vee \\ \partial_s N_8^\vee \\ \partial_t N_8^\vee \end{pmatrix} d\Omega = - \begin{pmatrix} \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (-1)(1-s)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(-1)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(1-s)(-1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (+1)(1-s)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(-1)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(1-s)(-1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (+1)(1+s)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(+1)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(1+s)(-1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (-1)(1+s)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(+1)(1-t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(1+s)(-1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (-1)(1-s)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(-1)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(1-s)(+1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (+1)(1-s)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(-1)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(1-s)(+1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (+1)(1+s)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(+1)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1+r)(1+s)(+1) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (-1)(1+s)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(+1)(1+t) dr ds dt \\ \int_{-1}^{+1} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{8} (1-r)(1+s)(+1) dr ds dt \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \\ 1 \\ -1 \\ -1 \\ 1 \\ 1 \\ -1 \\ 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ 1 \\ -1 \\ -1 \\ -1 \\ -1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$$

If we consider a macro-element 2x2x2 of size  $L_x=L_y=L_z=4$ , apply velocity b.c on the all sides we are left with

$$\tilde{G} = \begin{pmatrix} -1. & -1. & -1. & -1. & 1. & 1. & 1. & 1. \\ -1. & -1. & 1. & 1. & -1. & -1. & 1. & 1. \\ -1. & 1. & -1. & 1. & -1. & 1. & -1. & 1. \end{pmatrix}$$

Null space has dimension 5:

```
[[-0.35355339 0.35355339 0.35355339 0.35355339 0.35355339]
 [0.35355339 -0.35355339 -0.35355339 0.35355339 0.35355339]
 [0.35355339 -0.35355339 0.35355339 -0.35355339 0.35355339]
 [0.41990569 0.58009431 0.19336477 0.19336477 -0.19336477]
 [0.58009431 0.41990569 -0.19336477 -0.19336477 0.19336477]
 [0.19336477 -0.19336477 0.74028293 0.03317615 -0.03317615]
 [0.19336477 -0.19336477 0.03317615 0.74028293 -0.03317615]
 [-0.19336477 0.19336477 -0.03317615 -0.03317615 0.74028293]]
```

### K.0.3 $Q_1 \times Q_1$ element

For this element,  $m_v = 4$  and  $m_p = 4$  so  $\mathbb{G}_{el}$  is a  $8 \times 4$  matrix:

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega \quad (\text{K.1})$$

$$= - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\gamma & 0 & \partial_s N_1^\gamma \\ 0 & \partial_s N_1^\gamma & \partial_r N_1^\gamma \\ \partial_r N_2^\gamma & 0 & \partial_s N_2^\gamma \\ 0 & \partial_s N_2^\gamma & \partial_r N_2^\gamma \\ \partial_r N_3^\gamma & 0 & \partial_s N_3^\gamma \\ 0 & \partial_s N_3^\gamma & \partial_r N_3^\gamma \\ \partial_r N_4^\gamma & 0 & \partial_s N_4^\gamma \\ 0 & \partial_s N_4^\gamma & \partial_r N_4^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p & N_2^p & N_3^p & N_4^p \\ N_1^p & N_2^p & N_3^p & N_4^p \\ 0 & 0 & \dots & 0 \end{pmatrix} d\Omega \quad (\text{K.2})$$

$$= - \int_{\Omega_e} \begin{pmatrix} N_1^p \partial_r N_1^\gamma & N_2^p \partial_r N_1^\gamma & N_3^p \partial_r N_1^\gamma & N_4^p \partial_r N_1^\gamma \\ N_1^p \partial_s N_1^\gamma & N_2^p \partial_s N_1^\gamma & N_3^p \partial_s N_1^\gamma & N_4^p \partial_s N_1^\gamma \\ N_1^p \partial_r N_2^\gamma & N_2^p \partial_r N_2^\gamma & N_3^p \partial_r N_2^\gamma & N_4^p \partial_r N_2^\gamma \\ N_1^p \partial_s N_2^\gamma & N_2^p \partial_s N_2^\gamma & N_3^p \partial_s N_2^\gamma & N_4^p \partial_s N_2^\gamma \\ N_1^p \partial_r N_3^\gamma & N_2^p \partial_r N_3^\gamma & N_3^p \partial_r N_3^\gamma & N_4^p \partial_r N_3^\gamma \\ N_1^p \partial_s N_3^\gamma & N_2^p \partial_s N_3^\gamma & N_3^p \partial_s N_3^\gamma & N_4^p \partial_s N_3^\gamma \\ N_1^p \partial_r N_4^\gamma & N_2^p \partial_r N_4^\gamma & N_3^p \partial_r N_4^\gamma & N_4^p \partial_r N_4^\gamma \\ N_1^p \partial_s N_4^\gamma & N_2^p \partial_s N_4^\gamma & N_3^p \partial_s N_4^\gamma & N_4^p \partial_s N_4^\gamma \end{pmatrix} d\Omega \quad (\text{K.3})$$

$$(\text{K.4})$$

We have  $N_i^\gamma = N_i^p$  with  $i = 1, 2, 3, 4$ , so we can drop the superscripts and we can write:

$$\mathbb{G}_{el} = - \int_{\Omega_e} \begin{pmatrix} N_1 \partial_r N_1 & N_2 \partial_r N_1 & N_3 \partial_r N_1 & N_4 \partial_r N_1 \\ N_1 \partial_s N_1 & N_2 \partial_s N_1 & N_3 \partial_s N_1 & N_4 \partial_s N_1 \\ N_1 \partial_r N_2 & N_2 \partial_r N_2 & N_3 \partial_r N_2 & N_4 \partial_r N_2 \\ N_1 \partial_s N_2 & N_2 \partial_s N_2 & N_3 \partial_s N_2 & N_4 \partial_s N_2 \\ N_1 \partial_r N_3 & N_2 \partial_r N_3 & N_3 \partial_r N_3 & N_4 \partial_r N_3 \\ N_1 \partial_s N_3 & N_2 \partial_s N_3 & N_3 \partial_s N_3 & N_4 \partial_s N_3 \\ N_1 \partial_r N_4 & N_2 \partial_r N_4 & N_3 \partial_r N_4 & N_4 \partial_r N_4 \\ N_1 \partial_s N_4 & N_2 \partial_s N_4 & N_3 \partial_s N_4 & N_4 \partial_s N_4 \end{pmatrix} d\Omega \quad (\text{K.5})$$

$$\int_{\Omega_e} N_1 \partial_r N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (-1)(1-s) dr ds = -1/3 \quad (\text{K.6})$$

$$\int_{\Omega_e} N_1 \partial_s N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (1-r)(-1) dr ds = -1/3 \quad (\text{K.7})$$

$$\int_{\Omega_e} N_1 \partial_r N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (+1)(1-s) dr ds = 1/3 \quad (\text{K.8})$$

$$\int_{\Omega_e} N_1 \partial_s N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (1+r)(-1) dr ds = -1/6 \quad (\text{K.9})$$

$$\int_{\Omega_e} N_1 \partial_r N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (+1)(1+s) dr ds = 1/6 \quad (\text{K.10})$$

$$\int_{\Omega_e} N_1 \partial_s N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (1+r)(+1) dr ds = 1/6 \quad (\text{K.11})$$

$$\int_{\Omega_e} N_1 \partial_r N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (-1)(1+s) dr ds = -1/6 \quad (\text{K.12})$$

$$\int_{\Omega_e} N_1 \partial_s N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{4} (1-r)(+1) dr ds = 1/3 \quad (\text{K.13})$$

$$\int_{\Omega_e} N_2 \partial_r N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (-1)(1-s) dr ds = -1/3 \quad (\text{K.14})$$

$$\int_{\Omega_e} N_2 \partial_s N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (1-r)(-1) dr ds = -1/6 \quad (\text{K.15})$$

$$\int_{\Omega_e} N_2 \partial_r N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (+1)(1-s) dr ds = 1/3 \quad (\text{K.16})$$

$$\int_{\Omega_e} N_2 \partial_s N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (1+r)(-1) dr ds = -1/3 \quad (\text{K.17})$$

$$\int_{\Omega_e} N_2 \partial_r N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (+1)(1+s) dr ds = 1/6 \quad (\text{K.18})$$

$$\int_{\Omega_e} N_2 \partial_s N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (1+r)(+1) dr ds = 1/3 \quad (\text{K.19})$$

$$\int_{\Omega_e} N_2 \partial_r N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (-1)(1+s) dr ds = -1/6 \quad (\text{K.20})$$

$$\int_{\Omega_e} N_2 \partial_s N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{4} (1-r)(+1) dr ds = 1/6 \quad (\text{K.21})$$

$$\int_{\Omega_e} N_3 \partial_r N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (-1)(1-s) dr ds = -1/6 \quad (\text{K.22})$$

$$\int_{\Omega_e} N_3 \partial_s N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (1-r)(-1) dr ds = -1/6 \quad (\text{K.23})$$

$$\int_{\Omega_e} N_3 \partial_r N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (+1)(1-s) dr ds = 1/6 \quad (\text{K.24})$$

$$\int_{\Omega_e} N_3 \partial_s N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (1+r)(-1) dr ds = -1/3 \quad (\text{K.25})$$

$$\int_{\Omega_e} N_3 \partial_r N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (+1)(1+s) dr ds = 1/3 \quad (\text{K.26})$$

$$\int_{\Omega_e} N_3 \partial_s N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (1+r)(+1) dr ds = 1/3 \quad (\text{K.27})$$

$$\int_{\Omega_e} N_3 \partial_r N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (-1)(1+s) dr ds = -1/3 \quad (\text{K.28})$$

$$\int_{\Omega_e} N_3 \partial_s N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{4} (1-r)(+1) dr ds = 1/6 \quad (\text{K.29})$$

$$\int_{\Omega_e} N_4 \partial_r N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (-1)(1-s) dr ds = -1/6 \quad (\text{K.30})$$

$$\int_{\Omega_e} N_4 \partial_s N_1 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (1-r)(-1) dr ds = -1/3 \quad (\text{K.31})$$

$$\int_{\Omega_e} N_4 \partial_r N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (+1)(1-s) dr ds = 1/6 \quad (\text{K.32})$$

$$\int_{\Omega_e} N_4 \partial_s N_2 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (1+r)(-1) dr ds = -1/6 \quad (\text{K.33})$$

$$\int_{\Omega_e} N_4 \partial_r N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (+1)(1+s) dr ds = 1/3 \quad (\text{K.34})$$

$$\int_{\Omega_e} N_4 \partial_s N_3 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (1+r)(+1) dr ds = 1/6 \quad (\text{K.35})$$

$$\int_{\Omega_e} N_4 \partial_r N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (-1)(1+s) dr ds = -1/3 \quad (\text{K.36})$$

$$\int_{\Omega_e} N_4 \partial_s N_4 d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{4} (1-r)(+1) dr ds = 1/3 \quad (\text{K.37})$$

Putting it all together:

$$\mathbb{G}_{el} = - \begin{pmatrix} -1/3 & -1/3 & -1/6 & -1/6 \\ -1/3 & -1/6 & -1/6 & -1/3 \\ 1/3 & 1/3 & 1/6 & 1/6 \\ -1/6 & -1/3 & -1/3 & -1/6 \\ 1/6 & 1/6 & 1/3 & 1/3 \\ 1/6 & 1/3 & 1/3 & 1/6 \\ -1/6 & -1/6 & -1/3 & -1/3 \\ 1/3 & 1/6 & 1/6 & 1/3 \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 2 & 2 & 1 & 1 \\ 2 & 1 & 1 & 2 \\ -2 & -2 & -1 & -1 \\ 1 & 2 & 2 & 1 \\ -1 & -1 & -2 & -2 \\ -1 & -2 & -2 & -1 \\ 1 & 1 & 2 & 2 \\ -2 & -1 & -1 & -2 \end{pmatrix} \quad (\text{K.38})$$

I have implemented a 3x3 quadrature integration to numerically compute the matrix in the file `python_codes/Gel/programQ1Q1.py`. The code returns:

```

[[0.33333333 0.33333333 0.16666667 0.16666667]
 [0.33333333 0.16666667 0.16666667 0.33333333]
 [-0.33333333 -0.33333333 -0.16666667 -0.16666667]
 [0.16666667 0.33333333 0.33333333 0.16666667]
 [-0.16666667 -0.16666667 -0.33333333 -0.33333333]
 [-0.16666667 -0.33333333 -0.33333333 -0.16666667]
 [0.16666667 0.16666667 0.33333333 0.33333333]
 [-0.33333333 -0.16666667 -0.16666667 -0.33333333]]

```

which is indeed what we have obtained above.

Similarly to the Q1P0 element, a macroelement made of a single element has zero left over vel dof after b.c. are applied on all sides, so we resort to the following macroelement:

|             |             |      |
|-------------|-------------|------|
| velocity    | pressure    |      |
| 7====8====9 | 7====8====9 |      |
|             |             |      |
| 4====5====6 | 4====5====6 | NV=9 |
|             |             |      |
| 1====2====3 | 1====2====3 | NP=9 |

After assembly we have  $\mathbb{G}$  is a  $ndofV * NV \times ndofP * NP = 18 * 9$  matrix:

$$\mathbb{G} = \frac{1}{6} \begin{pmatrix} 2 & 2 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 2 & 1 & 0 & 2 & 1 & 0 & 0 & 0 & 0 \\ -2 & 0 & 2 & -1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 4 & 1 & 1 & 4 & 1 & 0 & 0 & 0 \\ 0 & -2 & -2 & 0 & -1 & -1 & 0 & 0 & 0 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 0 & 0 \\ 1 & 1 & 0 & 4 & 4 & 0 & 1 & 1 & 0 \\ -2 & -1 & 0 & 0 & 0 & 0 & 2 & 1 & 0 \\ -1 & 0 & 1 & -4 & 0 & 4 & -1 & 0 & 1 \\ -1 & -4 & -1 & 0 & 0 & 0 & 1 & 4 & 1 \\ 0 & -1 & -1 & 0 & -4 & -4 & 0 & -1 & -1 \\ 0 & -1 & -2 & 0 & 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 1 & 1 & 0 & 2 & 2 & 0 \\ 0 & 0 & 0 & -2 & -1 & 0 & -2 & -1 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 & -2 & 0 & 2 \\ 0 & 0 & 0 & -1 & -4 & -1 & -1 & -4 & -1 \\ 0 & 0 & 0 & 0 & -1 & -1 & 0 & -2 & -2 \\ 0 & 0 & 0 & 0 & -1 & -2 & 0 & -1 & -2 \end{pmatrix}$$

and after imposing boundary conditions on nodes 1,2,3,4,6,7,8,9:

$$\mathbb{G} = \frac{1}{6} \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & -4 & 0 & 4 & -1 & 0 & 1 \\ -1 & -4 & -1 & 0 & 0 & 0 & 1 & 4 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

or,

$$\mathbb{G} = \frac{1}{6} \begin{pmatrix} -1 & 0 & 1 & -4 & 0 & 4 & -1 & 0 & 1 \\ -1 & -4 & -1 & 0 & 0 & 0 & 1 & 4 & 1 \end{pmatrix}$$

When passed to *null\_space* as argument it returns the following nullspace:

```
[[1.47619e-01 -6.57142e-01 -7.40148e-18 6.57142e-01 -1.47619e-01 6.66666e-02 1.80952e-
[-1.90476e-01 9.52380e-02 -7.40148e-17 -9.52380e-02 1.90476e-01 6.66666e-01 1.42857e-
[9.57142e-01 1.04761e-01 -7.40148e-18 -1.04761e-01 4.28571e-02 6.66666e-02 -9.52380e-
[9.52380e-02 6.19047e-01 2.34780e-34 3.80952e-01 -9.52380e-02 -2.11471e-18 9.52380e-
[-8.45884e-18 4.22942e-18 1.00000e+00 -4.22942e-18 8.45884e-18 2.96059e-17 6.34413e-
[-9.52380e-02 3.80952e-01 5.35591e-34 6.19047e-01 9.52380e-02 -4.82418e-18 -9.52380e-
[4.28571e-02 -1.04761e-01 7.40148e-18 1.04761e-01 9.57142e-01 -6.66666e-02 9.52380e-
[7.61904e-02 -3.80952e-02 2.96059e-17 3.80952e-02 -7.61904e-02 7.33333e-01 -5.71428e-
[-4.76190e-03 8.57142e-02 7.40148e-18 -8.57142e-02 4.76190e-03 -6.66666e-02 9.61904e-
```

which is very bad: the dimension of the nullspace is 9!

Note that it does not mean that this element is unstable (see Q2Q1) since it is a sufficient but not necessary condition. We could test with larger macroelements (see Q2Q1) and these could prove to have a properly sized nullspace.

See code `python_codes/Gel/macro_element_q1q1.py`



### K.0.4 $Q_1^+ \times Q_1$ element

For the quadrilateral MINI element,  $m_v = 5$  and  $m_p = 4$  so  $\mathbb{G}_{el}$  is a  $10 \times 4$  matrix:

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega \quad (\text{K.39})$$

$$= - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\gamma & 0 & \partial_s N_1^\gamma \\ 0 & \partial_s N_1^\gamma & \partial_r N_1^\gamma \\ \partial_r N_2^\gamma & 0 & \partial_s N_2^\gamma \\ 0 & \partial_s N_2^\gamma & \partial_r N_2^\gamma \\ \partial_r N_3^\gamma & 0 & \partial_s N_3^\gamma \\ 0 & \partial_s N_3^\gamma & \partial_r N_3^\gamma \\ \partial_r N_4^\gamma & 0 & \partial_s N_4^\gamma \\ 0 & \partial_s N_4^\gamma & \partial_r N_4^\gamma \\ \partial_r N_5^\gamma & 0 & \partial_s N_5^\gamma \\ 0 & \partial_s N_5^\gamma & \partial_r N_5^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p & N_2^p & N_3^p & N_4^p \\ N_1^p & N_2^p & N_3^p & N_4^p \\ 0 & 0 & \dots & 0 \end{pmatrix} d\Omega \quad (\text{K.40})$$

$$= - \int_{\Omega_e} \begin{pmatrix} N_1^p \partial_r N_1^\gamma & N_2^p \partial_r N_1^\gamma & N_3^p \partial_r N_1^\gamma & N_4^p \partial_r N_1^\gamma \\ N_1^p \partial_s N_1^\gamma & N_2^p \partial_s N_1^\gamma & N_3^p \partial_s N_1^\gamma & N_4^p \partial_s N_1^\gamma \\ N_1^p \partial_r N_2^\gamma & N_2^p \partial_r N_2^\gamma & N_3^p \partial_r N_2^\gamma & N_4^p \partial_r N_2^\gamma \\ N_1^p \partial_s N_2^\gamma & N_2^p \partial_s N_2^\gamma & N_3^p \partial_s N_2^\gamma & N_4^p \partial_s N_2^\gamma \\ N_1^p \partial_r N_3^\gamma & N_2^p \partial_r N_3^\gamma & N_3^p \partial_r N_3^\gamma & N_4^p \partial_r N_3^\gamma \\ N_1^p \partial_s N_3^\gamma & N_2^p \partial_s N_3^\gamma & N_3^p \partial_s N_3^\gamma & N_4^p \partial_s N_3^\gamma \\ N_1^p \partial_r N_4^\gamma & N_2^p \partial_r N_4^\gamma & N_3^p \partial_r N_4^\gamma & N_4^p \partial_r N_4^\gamma \\ N_1^p \partial_s N_4^\gamma & N_2^p \partial_s N_4^\gamma & N_3^p \partial_s N_4^\gamma & N_4^p \partial_s N_4^\gamma \\ N_1^p \partial_r N_5^\gamma & N_2^p \partial_r N_5^\gamma & N_3^p \partial_r N_5^\gamma & N_4^p \partial_r N_5^\gamma \\ N_1^p \partial_s N_5^\gamma & N_2^p \partial_s N_5^\gamma & N_3^p \partial_s N_5^\gamma & N_4^p \partial_s N_5^\gamma \end{pmatrix} d\Omega \quad (\text{K.41})$$

We have :

$$N_1^\gamma = N_1^p - \frac{1}{4}b(r, s) \quad (\text{K.42})$$

$$N_2^\gamma = N_2^p - \frac{1}{4}b(r, s) \quad (\text{K.43})$$

$$N_3^\gamma = N_3^p - \frac{1}{4}b(r, s) \quad (\text{K.44})$$

$$N_4^\gamma = N_4^p - \frac{1}{4}b(r, s) \quad (\text{K.45})$$

$$N_5^\gamma = b(r, s) \quad (\text{K.46})$$

so that (once again I drop the superscripts)

$$\mathbb{G}_{el} = - \int_{\Omega_e} \begin{pmatrix} N_1 \partial_r N_1 & N_2 \partial_r N_1 & N_3 \partial_r N_1 & N_4 \partial_r N_1 \\ N_1 \partial_s N_1 & N_2 \partial_s N_1 & N_3 \partial_s N_1 & N_4 \partial_s N_1 \\ N_1 \partial_r N_2 & N_2 \partial_r N_2 & N_3 \partial_r N_2 & N_4 \partial_r N_2 \\ N_1 \partial_s N_2 & N_2 \partial_s N_2 & N_3 \partial_s N_2 & N_4 \partial_s N_2 \\ N_1 \partial_r N_3 & N_2 \partial_r N_3 & N_3 \partial_r N_3 & N_4 \partial_r N_3 \\ N_1 \partial_s N_3 & N_2 \partial_s N_3 & N_3 \partial_s N_3 & N_4 \partial_s N_3 \\ N_1 \partial_r N_4 & N_2 \partial_r N_4 & N_3 \partial_r N_4 & N_4 \partial_r N_4 \\ N_1 \partial_s N_4 & N_2 \partial_s N_4 & N_3 \partial_s N_4 & N_4 \partial_s N_4 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} d\Omega + \frac{1}{4} \int_{\Omega_e} \begin{pmatrix} N_1 \partial_r b & N_2 \partial_r b & N_3 \partial_r b & N_4 \partial_r b \\ N_1 \partial_s b & N_2 \partial_s b & N_3 \partial_s b & N_4 \partial_s b \\ N_1 \partial_r b & N_2 \partial_r b & N_3 \partial_r b & N_4 \partial_r b \\ N_1 \partial_s b & N_2 \partial_s b & N_3 \partial_s b & N_4 \partial_s b \\ N_1 \partial_r b & N_2 \partial_r b & N_3 \partial_r b & N_4 \partial_r b \\ N_1 \partial_s b & N_2 \partial_s b & N_3 \partial_s b & N_4 \partial_s b \\ N_1 \partial_r b & N_2 \partial_r b & N_3 \partial_r b & N_4 \partial_r b \\ N_1 \partial_s b & N_2 \partial_s b & N_3 \partial_s b & N_4 \partial_s b \\ -4N_1 \partial_r b & -4N_2 \partial_r b & -4N_3 \partial_r b & -4N_4 \partial_r b \\ -4N_1 \partial_s b & -4N_2 \partial_s b & -4N_3 \partial_s b & -4N_4 \partial_s b \end{pmatrix} d\Omega$$

The matrix which only contains  $N_i$  functions is in fact the  $\mathbb{G}_{el}$  matrix for standard  $Q_1 \times Q_1$  elements as we have seen in the previous section so we need not recompute it.

## Bubble function 1

Let us now assume the bubble is bubble 1:

$$\begin{aligned} b_1(r, s) &= (1-r)(1-s)(1-r^2)(1-s^2) \\ \partial_r b_1(r, s) &= (3r^2 - 2r - 1)(1-s)(1-s^2) \\ \partial_s b_1(r, s) &= (3s^2 - 2s - 1)(1-r)(1-r^2) \end{aligned}$$

$$\begin{aligned} \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_1 \partial_r b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s)(3r^2 - 2r - 1)(1-s)(1-s^2) dr ds = 2/15 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_1 \partial_s b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s)(3s^2 - 2s - 1)(1-r)(1-r^2) dr ds = 2/15 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_2 \partial_r b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s)(3r^2 - 2r - 1)(1-s)(1-s^2) dr ds = -2/15 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_2 \partial_s b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s)(3s^2 - 2s - 1)(1-r)(1-r^2) dr ds = 4/45 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_3 \partial_r b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s)(3r^2 - 2r - 1)(1-s)(1-s^2) dr ds = -4/45 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_3 \partial_s b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s)(3s^2 - 2s - 1)(1-r)(1-r^2) dr ds = -4/45 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_4 \partial_r b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s)(3r^2 - 2r - 1)(1-s)(1-s^2) dr ds = 4/45 \\ \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_4 \partial_s b_1 dr ds &= \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s)(3s^2 - 2s - 1)(1-r)(1-r^2) dr ds = -2/15 \end{aligned}$$

$$\mathbb{G}_{el} = - \begin{pmatrix} -1/3 & -1/3 & -1/6 & -1/6 \\ -1/3 & -1/6 & -1/6 & -1/3 \\ 1/3 & 1/3 & 1/6 & 1/6 \\ -1/6 & -1/3 & -1/3 & -1/6 \\ 1/6 & 1/6 & 1/3 & 1/3 \\ 1/6 & 1/3 & 1/3 & 1/6 \\ -1/6 & -1/6 & -1/3 & -1/3 \\ 1/3 & 1/6 & 1/6 & 1/3 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 2/15 & -2/15 & -4/45 & 4/45 \\ 2/15 & 4/45 & -4/45 & -2/15 \\ 2/15 & -2/15 & -4/45 & 4/45 \\ 2/15 & 4/45 & -4/45 & -2/15 \\ 2/15 & -2/15 & -4/45 & 4/45 \\ 2/15 & 4/45 & -4/45 & -2/15 \\ 2/15 & -2/15 & -4/45 & 4/45 \\ 2/15 & 4/45 & -4/45 & -2/15 \\ -8/15 & 8/15 & 16/45 & -16/45 \\ -8/15 & -16/45 & 16/45 & 8/15 \end{pmatrix} \quad (\text{K.47})$$

I have implemented a 3x3 quadrature integration to numerically compute the matrix in the file `python_codes/Gel/programQ1pQ1.py`. The code returns:

```
[[0.46666667 0.2 0.07777778 0.25555556]
 [0.46666667 0.25555556 0.07777778 0.2]
 [-0.2 -0.46666667 -0.25555556 -0.07777778]
 [0.3 0.42222222 0.24444444 0.03333333]
 [-0.03333333 -0.3 -0.42222222 -0.24444444]
 [-0.03333333 -0.24444444 -0.42222222 -0.3]
 [0.3 0.03333333 0.24444444 0.42222222]
 [-0.2 -0.07777778 -0.25555556 -0.46666667]
 [-0.53333333 0.53333333 0.35555556 -0.35555556]
 [-0.53333333 -0.35555556 0.35555556 0.53333333]]
```

which is indeed what we have obtained above. It can be rewritten

$$\mathbb{G}_{el} = \frac{1}{90} \begin{pmatrix} 42 & 18 & 7 & 23 \\ 42 & 23 & 7 & 18 \\ -18 & -42 & -23 & -7 \\ 27 & 38 & 22 & 3 \\ -3 & -27 & -38 & -22 \\ -3 & -22 & -38 & -27 \\ 27 & 3 & 22 & 38 \\ -18 & -7 & -23 & -42 \\ -48 & 48 & 32 & -32 \\ -48 & -32 & 32 & 48 \end{pmatrix}$$

Let us now build a macroelement of size LxxLy=4x4 made of 2x2 elements. Each element has a Gel like the one above since they are of size 2x2:

| velocity    | pressure    |       |
|-------------|-------------|-------|
| 7====8====9 | 7====8====9 |       |
| 12   13     |             |       |
| 4====5====6 | 4====5====6 | NV=13 |
| 10   11     |             |       |
| 1====2====3 | 1====2====3 | NP=9  |

I am here following the approach by Lamichhane [741] but I am not sure why he did not use a single element macro-element? Probably because when applying bc on all four nodes of a single element, the left over matrix  $\tilde{G}_{el}$  is composed of the last two rows of  $\mathbb{G}_{el}$  and this has a nullspace of dimension 2.

After assembly we have  $\mathbb{G}$  is a  $ndofV * NV \times ndofP * NP = 26 * 9$  matrix:

$$\mathbb{G} = \frac{1}{90} \begin{pmatrix} 42 & 18 & 0 & 23 & 7 & 0 & 0 & 0 & 0 \\ 42 & 23 & 0 & 18 & 7 & 0 & 0 & 0 & 0 \\ -18 & 0 & 18 & -7 & 0 & 7 & 0 & 0 & 0 \\ 27 & 80 & 23 & 3 & 40 & 7 & 0 & 0 & 0 \\ 0 & -18 & -42 & 0 & -7 & -23 & 0 & 0 & 0 \\ 0 & 27 & 38 & 0 & 3 & 22 & 0 & 0 & 0 \\ 27 & 3 & 0 & 80 & 40 & 0 & 23 & 7 & 0 \\ -18 & -7 & 0 & 0 & 0 & 0 & 18 & 7 & 0 \\ -3 & 0 & 3 & -40 & 0 & 40 & -7 & 0 & 7 \\ -3 & -40 & -7 & 0 & 0 & 0 & 3 & 40 & 7 \\ 0 & -3 & -27 & 0 & -40 & -80 & 0 & -7 & -23 \\ 0 & -3 & -22 & 0 & 0 & 0 & 0 & 3 & 22 \\ 0 & 0 & 0 & 27 & 3 & 0 & 38 & 22 & 0 \\ 0 & 0 & 0 & -18 & -7 & 0 & -42 & -23 & 0 \\ 0 & 0 & 0 & -3 & 0 & 3 & -22 & 0 & 22 \\ 0 & 0 & 0 & -3 & -40 & -7 & -27 & -80 & -23 \\ 0 & 0 & 0 & 0 & -3 & -27 & 0 & -22 & -38 \\ 0 & 0 & 0 & 0 & -3 & -22 & 0 & -27 & -38 \\ -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 & 0 \\ -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 & 0 \\ 0 & -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 \\ 0 & -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 \\ 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 & 0 \\ 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 & 0 \\ 0 & 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 \\ 0 & 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 \end{pmatrix}$$

After boundary conditions on Vnodes 1,2,3,4,6,7,8,9, the matrix  $G$  looks like:

$$\mathbb{G} = \frac{1}{90} \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -3 & 0 & 3 & -40 & 0 & 40 & -7 & 0 & 7 \\ -3 & -40 & -7 & 0 & 0 & 0 & 3 & 40 & 7 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 & 0 \\ -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 & 0 \\ 0 & -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 \\ 0 & -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 \\ 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 & 0 \\ 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 & 0 \\ 0 & 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 \\ 0 & 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 \end{pmatrix}$$

or, removing the lines with only zeros, we arrive at a  $10 \times 9$  matrix as in Lamichhane [741]:

$$\tilde{\mathbb{G}} = \frac{1}{90} \begin{pmatrix} -3 & 0 & 3 & -40 & 0 & 40 & -7 & 0 & 7 \\ -3 & -40 & -7 & 0 & 0 & 0 & 3 & 40 & 7 \\ -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 & 0 \\ -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 & 0 \\ 0 & -48 & 48 & 0 & -32 & 32 & 0 & 0 & 0 \\ 0 & -48 & -32 & 0 & 48 & 32 & 0 & 0 & 0 \\ 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 & 0 \\ 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 & 0 \\ 0 & 0 & 0 & 0 & -48 & 48 & 0 & -32 & 32 \\ 0 & 0 & 0 & 0 & -48 & -32 & 0 & 48 & 32 \end{pmatrix}$$

This matrix is then passed as argument to the *null\_space* function which returns a single vector such that  $\ker(\tilde{\mathbb{G}}) = (1, 1, 1, 1, 1, 1, 1, 1, 1)$ .

It must be said that the matrix above contains similar values as the one in [741] as well as the same number of nonzeros. Similarities: 16 times  $\pm 48/90 = \pm 8/15 = 2 \cdot 4/15$  and 16 times  $\pm 32/90 = \pm 16/45 = 2 \cdot 8/45$  as in the paper (aside from scaling factor 2). However the 12 remaining values differ ?

## Bubble function 2

When looking at bubble 2 with  $\beta = 1/4$ , we get

$$\mathbb{G}_{el} = \frac{1}{180} \begin{pmatrix} 79 & 41 & 9 & 51 \\ 79 & 51 & 9 & 41 \\ -41 & -79 & -51 & -9 \\ 49 & 81 & 39 & 11 \\ -11 & -49 & -81 & -39 \\ -11 & -39 & -81 & -49 \\ 49 & 11 & 39 & 81 \\ -41 & -9 & -51 & -79 \\ -76 & 76 & 84 & -84 \\ -76 & -84 & 84 & 76 \end{pmatrix}$$

After assembly we have  $\mathbb{G}$  is a  $ndofV * NV \times ndofP * NP = 26 * 9$  matrix:

$$\mathbb{G} = \frac{1}{180} \begin{pmatrix} 79 & 41 & 0 & 51 & 9 & 0 & 0 & 0 & 0 \\ 79 & 51 & 0 & 41 & 9 & 0 & 0 & 0 & 0 \\ -41 & 0 & 41 & -9 & 0 & 9 & 0 & 0 & 0 \\ 49 & 160 & 51 & 11 & 80 & 9 & 0 & 0 & 0 \\ 0 & -41 & -79 & 0 & -9 & -51 & 0 & 0 & 0 \\ 0 & 49 & 81 & 0 & 11 & 39 & 0 & 0 & 0 \\ 49 & 11 & 0 & 160 & 80 & 0 & 51 & 9 & 0 \\ -41 & -9 & 0 & 0 & 0 & 0 & 41 & 9 & 0 \\ -11 & 0 & 11 & -80 & 0 & 80 & -9 & 0 & 9 \\ -11 & -80 & -9 & 0 & 0 & 0 & 11 & 80 & 9 \\ 0 & -11 & -49 & 0 & -80 & -160 & 0 & -9 & -51 \\ 0 & -11 & -39 & 0 & 0 & 0 & 0 & 11 & 39 \\ 0 & 0 & 0 & 49 & 11 & 0 & 81 & 39 & 0 \\ 0 & 0 & 0 & -41 & -9 & 0 & -79 & -51 & 0 \\ 0 & 0 & 0 & -11 & 0 & 11 & -39 & 0 & 39 \\ 0 & 0 & 0 & -11 & -80 & -9 & -49 & -160 & -51 \\ 0 & 0 & 0 & 0 & -11 & -49 & 0 & -39 & -81 \\ 0 & 0 & 0 & 0 & -11 & -39 & 0 & -49 & -81 \\ -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 & 0 \\ -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 & 0 \\ 0 & -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 \\ 0 & -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 \\ 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 & 0 \\ 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 & 0 \\ 0 & 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 \\ 0 & 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 \end{pmatrix}$$

After boundary conditions are applied:

$$\mathbb{G} = \frac{1}{180} \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -11 & 0 & 11 & -80 & 0 & 80 & -9 & 0 & 9 \\ -11 & -80 & -9 & 0 & 0 & 0 & 11 & 80 & 9 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 & 0 \\ -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 & 0 \\ 0 & -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 \\ 0 & -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 \\ 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 & 0 \\ 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 & 0 \\ 0 & 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 \\ 0 & 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 \end{pmatrix}$$

or,

$$\tilde{\mathbb{G}} = \frac{1}{180} \begin{pmatrix} -11 & 0 & 11 & -80 & 0 & 80 & -9 & 0 & 9 \\ -11 & -80 & -9 & 0 & 0 & 0 & 11 & 80 & 9 \\ -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 & 0 \\ -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 & 0 \\ 0 & -76 & 76 & 0 & -84 & 84 & 0 & 0 & 0 \\ 0 & -76 & -84 & 0 & 76 & 84 & 0 & 0 & 0 \\ 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 & 0 \\ 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 & 0 \\ 0 & 0 & 0 & 0 & -76 & 76 & 0 & -84 & 84 \\ 0 & 0 & 0 & 0 & -76 & -84 & 0 & 76 & 84 \end{pmatrix}$$

We make the same observation as for bubble 1: when this matrix is passed as argument to the *null\_space* function, it returns a single vector such that  $\ker(\tilde{\mathbb{G}}) = (1, 1, 1, 1, 1, 1, 1, 1, 1)$ .

**Special case:**  $\beta = 0$  The bubble is then

$$b(r, s) = (1 - r^2)(1 - s^2)$$

We repeat the same process and arrive at

$$\mathbb{G}_{el} = \frac{1}{18} \begin{pmatrix} 8 & 4 & 1 & 5 \\ 8 & 5 & 1 & 4 \\ -4 & -8 & -5 & -1 \\ 5 & 8 & 4 & 1 \\ -1 & -5 & -8 & -4 \\ -1 & -4 & -8 & -5 \\ 5 & 1 & 4 & 8 \\ -4 & -1 & -5 & -8 \\ -8 & 8 & 8 & -8 \\ -8 & -8 & 8 & 8 \end{pmatrix}$$

and

$$\tilde{\mathbb{G}} = \frac{1}{18} \begin{pmatrix} -1 & 0 & 1 & -8 & 0 & 8 & -1 & 0 & 1 \\ -1 & -8 & -1 & 0 & 0 & 0 & 1 & 8 & 1 \\ -8 & 8 & 0 & -8 & 8 & 0 & 0 & 0 & 0 \\ -8 & -8 & 0 & 8 & 8 & 0 & 0 & 0 & 0 \\ 0 & -8 & 8 & 0 & -8 & 8 & 0 & 0 & 0 \\ 0 & -8 & -8 & 0 & 8 & 8 & 0 & 0 & 0 \\ 0 & 0 & 0 & -8 & 8 & 0 & -8 & 8 & 0 \\ 0 & 0 & 0 & -8 & -8 & 0 & 8 & 8 & 0 \\ 0 & 0 & 0 & 0 & -8 & 8 & 0 & -8 & 8 \\ 0 & 0 & 0 & 0 & -8 & -8 & 0 & 8 & 8 \end{pmatrix}$$

Finally the *null\_space* function returns:

```
[[0.27870965 -0.34974409]
 [0.39102578 0.31160686]
 [0.27870965 -0.34974409]
 [0.39102578 0.31160686]
 [0.27870965 -0.34974409]
 [0.39102578 0.31160686]
 [0.27870965 -0.34974409]
 [0.39102578 0.31160686]
 [0.27870965 -0.34974409]]
```

i.e. the null space has dimension 2, so that the element is then not stable.

See code `python_codes/Gel/macro_element_q1pq1.py`

### K.0.5 $Q_1^+ \times Q_1$ element in 3D

For the quadrilateral MINI element,  $m_v = 9$  and  $m_p = 8$  so  $\mathbb{G}_{el}$  is a  $27 \times 8$  matrix (obtained with  $10^3$  quadrature points, no difference with  $6^3$  points).

$$\mathbb{G}_{el} = - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega = - \int_{\Omega_e} \begin{pmatrix} \partial_r N_1^\gamma & 0 & 0 & \partial_s N_1^\gamma & \partial_t N_1^\gamma & 0 \\ 0 & \partial_s N_1^\gamma & 0 & \partial_r N_1^\gamma & 0 & \partial_t N_1^\gamma \\ 0 & 0 & \partial_t N_1^\gamma & 0 & \partial_r N_1^\gamma & \partial_s N_1^\gamma \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \partial_r N_{27}^\gamma & 0 & 0 & \partial_s N_{27}^\gamma & \partial_t N_{27}^\gamma & 0 \\ 0 & \partial_s N_{27}^\gamma & 0 & \partial_r N_{27}^\gamma & 0 & \partial_t N_{27}^\gamma \\ 0 & 0 & \partial_t N_{27}^\gamma & 0 & \partial_r N_{27}^\gamma & \partial_s N_{27}^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p & N_2^p & \cdots & N_8^p \\ N_1^p & N_2^p & \cdots & N_8^p \\ N_1^p & N_2^p & \cdots & N_8^p \\ 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \end{pmatrix} d\Omega$$



[illegible]

$$\int \int \int N_1^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_1 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_1 \partial_r b dr ds dt \quad (\text{K.48})$$

$$= \frac{8}{27} - \frac{1}{8} \frac{4}{75} = \frac{1}{1350} (400 - 9) = \frac{372}{1350} \quad (\text{K.49})$$

$$\int \int \int N_2^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_2 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_2 \partial_r b dr ds dt \quad (\text{K.50})$$

$$= \frac{4}{27} - \frac{1}{8} \quad (\text{K.51})$$

$$\int \int \int N_3^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_3 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_3 \partial_r b dr ds dt \quad (\text{K.52})$$

$$= \frac{2}{27} - \frac{1}{8} \quad (\text{K.53})$$

$$\int \int \int N_4^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_4 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_4 \partial_r b dr ds dt \quad (\text{K.54})$$

$$= \frac{4}{27} \quad (\text{K.55})$$

$$\int \int \int N_5^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_5 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_5 \partial_r b dr ds dt \quad (\text{K.56})$$

$$= \frac{4}{27} \quad (\text{K.57})$$

$$\int \int \int N_6^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_6 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_6 \partial_r b dr ds dt \quad (\text{K.58})$$

$$= \frac{2}{27} \quad (\text{K.59})$$

$$\int \int \int N_7^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_7 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_7 \partial_r b dr ds dt \quad (\text{K.60})$$

$$= \frac{1}{27} \quad (\text{K.61})$$

$$\int \int \int N_8^p \partial_r N_1^\gamma dr ds dt = \int \int \int N_8 \partial_r N_1 dr ds dt - \frac{1}{8} \int \int \int N_8 \partial_r b dr ds dt \quad (\text{K.62})$$

$$= \frac{2}{27} \quad (\text{K.63})$$

$$(\text{K.64})$$

For bubble function #1:

$$\mathbb{G}_{el} = \frac{1}{1350} \begin{pmatrix} 372 & 228 & 102 & 198 & 198 & 102 & 43 & 107 \\ 372 & 198 & 102 & 228 & 198 & 107 & 43 & 102 \\ 372 & 198 & 107 & 198 & 228 & 102 & 43 & 102 \\ -228 & -372 & -198 & -102 & -102 & -198 & -107 & -43 \\ 222 & 348 & 252 & 78 & 123 & 182 & 118 & 27 \\ 222 & 348 & 182 & 123 & 78 & 252 & 118 & 27 \\ -78 & -222 & -348 & -252 & -27 & -123 & -182 & -118 \\ -78 & -252 & -348 & -222 & -27 & -118 & -182 & -123 \\ 147 & 198 & 332 & 198 & 3 & 102 & 268 & 102 \\ 222 & 78 & 252 & 348 & 123 & 27 & 118 & 182 \\ -228 & -102 & -198 & -372 & -102 & -43 & -107 & -198 \\ 222 & 123 & 182 & 348 & 78 & 27 & 118 & 252 \\ 222 & 78 & 27 & 123 & 348 & 252 & 118 & 182 \\ 222 & 123 & 27 & 78 & 348 & 182 & 118 & 252 \\ -228 & -102 & -43 & -102 & -372 & -198 & -107 & -198 \\ -78 & -222 & -123 & -27 & -252 & -348 & -182 & -118 \\ 147 & 198 & 102 & 3 & 198 & 332 & 268 & 102 \\ -78 & -252 & -118 & -27 & -222 & -348 & -182 & -123 \\ -3 & -147 & -198 & -102 & -102 & -198 & -332 & -268 \\ -3 & -102 & -198 & -147 & -102 & -268 & -332 & -198 \\ -3 & -102 & -268 & -102 & -147 & -198 & -332 & -198 \\ 147 & 3 & 102 & 198 & 198 & 102 & 268 & 332 \\ -78 & -27 & -123 & -222 & -252 & -118 & -182 & -348 \\ -78 & -27 & -118 & -252 & -222 & -123 & -182 & -348 \\ -576 & 576 & 384 & -384 & -384 & 384 & 256 & -256 \\ -576 & -384 & 384 & 576 & -384 & -256 & 256 & 384 \\ -576 & -384 & -256 & -384 & 576 & 384 & 256 & 384 \end{pmatrix}$$

Considering a 2x2x2 macroelement of size 4x4x4. Then NV=3\*3\*3+8=35, NP=3\*3\*3=27 Matrix  $\mathbb{G}$  is 3\*35x27=105\*27

After bc are imposed on all nodes on the boundary, 9 Vnodes are still free (8 bubble nodes and the node in the middle), i.e. 9\*3 dofs = 27 and there are 3x3x3=27 pressure nodes. So  $\tilde{G}$  is 27\*27. no less.

We get

$$\tilde{G} = \frac{1}{1350} ( j f g l j h h g l )$$

I have tried all kinds of bubbles but I usually arrive at a null space of dimension 3 to 5... typically for bubble 1, dim=5, while for bubble=2 dim=3.

### K.0.6 $Q_2 \times Q_1$ element

$$\begin{aligned}
\mathbb{G}_{el} &= - \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N} d\Omega \\
&= - \int_{\Omega_e} \begin{pmatrix} \partial_x N_1^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_1^\gamma & \partial_x N_1^\gamma \\ \partial_x N_2^\gamma & 0 & \partial_y N_2^\gamma \\ 0 & \partial_y N_2^\gamma & \partial_x N_2^\gamma \\ \partial_x N_3^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_3^\gamma & \partial_x N_1^\gamma \\ \partial_x N_4^\gamma & 0 & \partial_y N_2^\gamma \\ 0 & \partial_y N_4^\gamma & \partial_x N_2^\gamma \\ \partial_x N_5^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_5^\gamma & \partial_x N_1^\gamma \\ \partial_x N_6^\gamma & 0 & \partial_y N_2^\gamma \\ 0 & \partial_y N_6^\gamma & \partial_x N_2^\gamma \\ \partial_x N_7^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_7^\gamma & \partial_x N_1^\gamma \\ \partial_x N_8^\gamma & 0 & \partial_y N_2^\gamma \\ 0 & \partial_y N_8^\gamma & \partial_x N_2^\gamma \\ \partial_x N_9^\gamma & 0 & \partial_y N_1^\gamma \\ 0 & \partial_y N_9^\gamma & \partial_x N_1^\gamma \end{pmatrix} \cdot \begin{pmatrix} N_1^p & N_2^p & N_3^p & N_4^p \\ N_1^p & N_2^p & N_3^p & N_4^p \\ 0 & 0 & \dots & 0 \end{pmatrix} d\Omega \\
&= - \int_{\Omega_e} \begin{pmatrix} N_1^p \partial_r N_1^\gamma & N_2^p \partial_r N_1^\gamma & N_3^p \partial_r N_1^\gamma & N_4^p \partial_r N_1^\gamma \\ N_1^p \partial_s N_1^\gamma & N_2^p \partial_s N_1^\gamma & N_3^p \partial_s N_1^\gamma & N_4^p \partial_s N_1^\gamma \\ N_1^p \partial_r N_2^\gamma & N_2^p \partial_r N_2^\gamma & N_3^p \partial_r N_2^\gamma & N_4^p \partial_r N_2^\gamma \\ N_1^p \partial_s N_2^\gamma & N_2^p \partial_s N_2^\gamma & N_3^p \partial_s N_2^\gamma & N_4^p \partial_s N_2^\gamma \\ N_1^p \partial_r N_3^\gamma & N_2^p \partial_r N_3^\gamma & N_3^p \partial_r N_3^\gamma & N_4^p \partial_r N_3^\gamma \\ N_1^p \partial_s N_3^\gamma & N_2^p \partial_s N_3^\gamma & N_3^p \partial_s N_3^\gamma & N_4^p \partial_s N_3^\gamma \\ N_1^p \partial_r N_4^\gamma & N_2^p \partial_r N_4^\gamma & N_3^p \partial_r N_4^\gamma & N_4^p \partial_r N_4^\gamma \\ N_1^p \partial_s N_4^\gamma & N_2^p \partial_s N_4^\gamma & N_3^p \partial_s N_4^\gamma & N_4^p \partial_s N_4^\gamma \\ N_1^p \partial_r N_5^\gamma & N_2^p \partial_r N_5^\gamma & N_3^p \partial_r N_5^\gamma & N_4^p \partial_r N_5^\gamma \\ N_1^p \partial_s N_5^\gamma & N_2^p \partial_s N_5^\gamma & N_3^p \partial_s N_5^\gamma & N_4^p \partial_s N_5^\gamma \\ N_1^p \partial_r N_6^\gamma & N_2^p \partial_r N_6^\gamma & N_3^p \partial_r N_6^\gamma & N_4^p \partial_r N_6^\gamma \\ N_1^p \partial_s N_6^\gamma & N_2^p \partial_s N_6^\gamma & N_3^p \partial_s N_6^\gamma & N_4^p \partial_s N_6^\gamma \\ N_1^p \partial_r N_7^\gamma & N_2^p \partial_r N_7^\gamma & N_3^p \partial_r N_7^\gamma & N_4^p \partial_r N_7^\gamma \\ N_1^p \partial_s N_7^\gamma & N_2^p \partial_s N_7^\gamma & N_3^p \partial_s N_7^\gamma & N_4^p \partial_s N_7^\gamma \\ N_1^p \partial_r N_8^\gamma & N_2^p \partial_r N_8^\gamma & N_3^p \partial_r N_8^\gamma & N_4^p \partial_r N_8^\gamma \\ N_1^p \partial_s N_8^\gamma & N_2^p \partial_s N_8^\gamma & N_3^p \partial_s N_8^\gamma & N_4^p \partial_s N_8^\gamma \\ N_1^p \partial_r N_9^\gamma & N_2^p \partial_r N_9^\gamma & N_3^p \partial_r N_9^\gamma & N_4^p \partial_r N_9^\gamma \\ N_1^p \partial_s N_9^\gamma & N_2^p \partial_s N_9^\gamma & N_3^p \partial_s N_9^\gamma & N_4^p \partial_s N_9^\gamma \end{pmatrix} d\Omega \tag{K.65}
\end{aligned}$$

$$\int_{\Omega_e} N_1^p \partial_r N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{2} (2r-1) \frac{1}{2} s(s-1) dr ds = -5/18 \quad (\text{K.66})$$

$$\int_{\Omega_e} N_2^p \partial_s N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s) \frac{1}{2} r(r-1) \frac{1}{2} (2s-1) dr ds = -5/18 \quad (\text{K.67})$$

$$\int_{\Omega_e} N_3^p \partial_r N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{2} (2r-1) \frac{1}{2} s(s-1) dr ds = -1/18 \quad (\text{K.68})$$

$$\int_{\Omega_e} N_4^p \partial_s N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s) \frac{1}{2} r(r-1) \frac{1}{2} (2s-1) dr ds = 0 \quad (\text{K.69})$$

$$\int_{\Omega_e} N_1^p \partial_r N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{2} (2r-1) \frac{1}{2} s(s-1) dr ds = 0 \quad (\text{K.70})$$

$$\int_{\Omega_e} N_2^p \partial_s N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s) \frac{1}{2} r(r-1) \frac{1}{2} (2s-1) dr ds = 0 \quad (\text{K.71})$$

$$\int_{\Omega_e} N_3^p \partial_r N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{2} (2r-1) \frac{1}{2} s(s-1) dr ds = 0 \quad (\text{K.72})$$

$$\int_{\Omega_e} N_4^p \partial_s N_1^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s) \frac{1}{2} r(r-1) \frac{1}{2} (2s-1) dr ds = -1/18 \quad (\text{K.73})$$

... same procedure for 1,2,3,4,5,6,7...

$$\int_{\Omega_e} N_1^p \partial_r N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s)(-2r)(1-s^2) dr ds = 4/9 \quad (\text{K.74})$$

$$\int_{\Omega_e} N_2^p \partial_s N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1-s)(1-r^2)(-2s) dr ds = 4/9 \quad (\text{K.75})$$

$$\int_{\Omega_e} N_3^p \partial_r N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s)(-2r)(1-s^2) dr ds = -4/9 \quad (\text{K.76})$$

$$\int_{\Omega_e} N_4^p \partial_s N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1-s)(1-r^2)(-2s) dr ds = 4/9 \quad (\text{K.77})$$

$$\int_{\Omega_e} N_1^p \partial_r N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s)(-2r)(1-s^2) dr ds = -4/9 \quad (\text{K.78})$$

$$\int_{\Omega_e} N_2^p \partial_s N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1+r)(1+s)(1-r^2)(-2s) dr ds = -4/9 \quad (\text{K.79})$$

$$\int_{\Omega_e} N_3^p \partial_r N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s)(-2r)(1-s^2) dr ds = 4/9 \quad (\text{K.80})$$

$$\int_{\Omega_e} N_4^p \partial_s N_9^\gamma d\Omega = \int_{-1}^{+1} \int_{-1}^{+1} \frac{1}{4} (1-r)(1+s)(1-r^2)(-2s) dr ds = -4/9 \quad (\text{K.81})$$

We obtain

$$-\int_{\Omega_e} \begin{pmatrix} N_1^p \partial_r N_9^\gamma & N_2^p \partial_r N_9^\gamma & N_3^p \partial_r N_9^\gamma & N_4^p \partial_r N_9^\gamma \\ N_1^p \partial_s N_9^\gamma & N_2^p \partial_s N_9^\gamma & N_3^p \partial_s N_9^\gamma & N_4^p \partial_s N_9^\gamma \end{pmatrix} d\Omega = -\frac{4}{9} \begin{pmatrix} 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \end{pmatrix}$$

which is identical to Eq.(3.53) of Elman *et al.* [371].

I have implemented a 3x3 quadrature integration to numerically compute the matrix in the file `python_codes/Gel/programQ2Q1.py`:

```
[[2.77777778e-01 5.55555556e-02 1.73472348e-18 3.68628739e-18]
 [2.77777778e-01 -1.29020059e-17 -3.46944695e-18 5.55555556e-02]
 [-5.55555556e-02 -2.77777778e-01 -2.77555756e-17 -3.03576608e-18]
 [-2.19008839e-17 2.77777778e-01 5.55555556e-02 -4.33680869e-18]
 [-3.46944695e-18 -1.38777878e-17 -2.77777778e-01 -5.55555556e-02]
 [2.60208521e-18 -5.55555556e-02 -2.77777778e-01 6.93889390e-18]
 [4.01154804e-18 4.33680869e-18 5.55555556e-02 2.77777778e-01]
 [-5.55555556e-02 6.07153217e-18 2.08166817e-17 -2.77777778e-01]
 [-2.22222222e-01 2.22222222e-01 0.00000000e+00 -8.67361738e-19]
 [5.55555556e-01 5.55555556e-01 1.11111111e-01 1.11111111e-01]
 [-1.11111111e-01 -5.55555556e-01 -5.55555556e-01 -1.11111111e-01]
 [-8.67361738e-19 -2.22222222e-01 2.22222222e-01 6.93889390e-18]
 [-8.67361738e-18 -6.93889390e-18 2.22222222e-01 -2.22222222e-01]
 [-1.11111111e-01 -1.11111111e-01 -5.55555556e-01 -5.55555556e-01]
 [5.55555556e-01 1.11111111e-01 1.11111111e-01 5.55555556e-01]
 [-2.22222222e-01 -5.20417043e-18 6.93889390e-18 2.22222222e-01]
 [-4.44444444e-01 4.44444444e-01 4.44444444e-01 -4.44444444e-01]
 [-4.44444444e-01 -4.44444444e-01 4.44444444e-01 4.44444444e-01]]
```

or,

$$\mathbb{G}_{el} = \frac{1}{18} \begin{pmatrix} 5 & 1 & 0 & 0 \\ 5 & 0 & 0 & 1 \\ -1 & -5 & 0 & 0 \\ 0 & 5 & 1 & 0 \\ 0 & 0 & -5 & -1 \\ 0 & -1 & -5 & 0 \\ 0 & 0 & 1 & 5 \\ -1 & 0 & 0 & -5 \\ -4 & 4 & 0 & 0 \\ 10 & 10 & 2 & 2 \\ -2 & -10 & -10 & -2 \\ 0 & -4 & 4 & 0 \\ 0 & 0 & 4 & -4 \\ -2 & -2 & -10 & -10 \\ 10 & 2 & 2 & 10 \\ -4 & 0 & 0 & 4 \\ -8 & 8 & 8 & -8 \\ -8 & -8 & 8 & 8 \end{pmatrix}$$

Reading Elman [371] we see that this element is stable but patches of even and odd elements actually are needed to establish the stability of the element.

# Appendix L

## Computational Geophysics GEO4-1427 - projects

compgeo.tex

### L.0.1 Convection in a box \*

This exercise builds on your existing 2D advection-diffusion code. Scale up the benchmark described in Section 12.2.9 so that it runs in a 1000x1000 km domain with Earth-like parameters and velocities (the maximum velocity is denoted by  $\mathbf{v}_{conv}$  and will be varied). Start with an initial zero temperature field and Earth-like boundary conditions on the top and bottom, e.g.  $T = 20$  at the top and  $T = 1000$  at the bottom. Set  $k = 3$ ,  $C_p = 1250$  and  $\rho = 3000$ .

Run the code until steady state is reached. Implement an algorithm which computes the average temperature

$$\langle T \rangle = \frac{1}{L_x L_y} \iint T(x, y) dx dy$$

in the domain and plot it as a function of time. Also compute the root mean square velocity in the domain:

$$\mathbf{v}_{rms} = \sqrt{\frac{1}{L_x L_y} \iint (u^2 + v^2) dx dy}$$

Plot the steady state  $\langle T \rangle$  and  $\mathbf{v}_{rms}$  as a function of the resolution  $h$ . Plot the temperature on the  $x = L_x/2$  line for different values of  $\mathbf{v}_{conv}$ . When possible, make a link with the Mantle Dynamics practical.

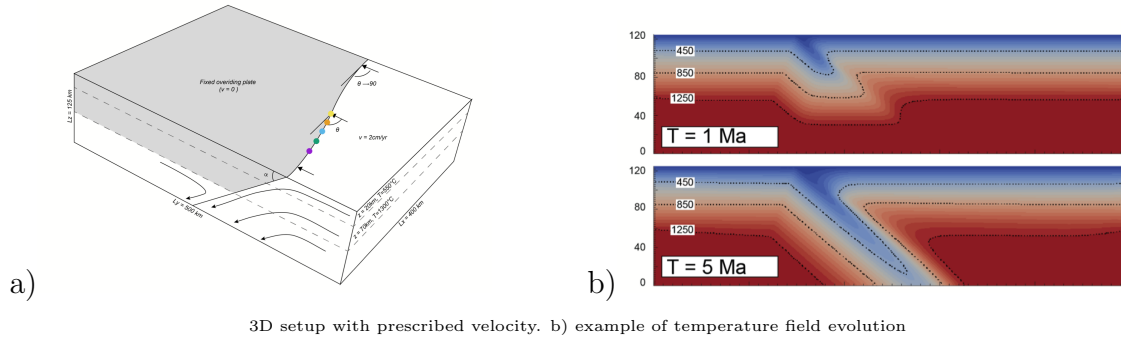
Bonus: Compute and plot the heat flux  $\vec{q} = -\vec{\nabla}T$  in the center of the elements.

### L.0.2 Corner flow subduction \*

- In this experiment the velocity is prescribed in geometrically simple subducting and overriding plates, while the velocity in the mantle is computed by means of an analytical solution coined 'corner flow velocity'. Details are to be found in G.K. Batchelor, *An introduction to fluid dynamics*. Cambridge University Press, 1967.
- Write a function which prescribes the velocity in a lithospheric sized domain.
- Use this velocity to drive the system in time (choose the appropriate values for the coefficients in the heat transport equation)
- prescribe a constant temperature value at the top, and fix the temperature on both sides, but only in the plates (along lengths  $l_1$  and  $l_2$ . Choose an appropriate plate temperature model.

- Run the model over millions of years with different velocities.
- Measure the depth of the isotherm  $800^\circ$  as a function of time (bonus).
- Is steady state ever reached ?

If all goes well, you should be able to recover similar results:



what are l1, l2 ? rephrase !

### L.0.3 From 2D to 3D \*\*

Rewrite your 2D FEM advection-diffusion code so that it now runs on a cube. You will need to create a new connectivity array, compute new elemental matrices, etc ... Center the cube on the origin of the axis system.

Compute the steady-state solution of a problem without heat advection. The domain is a unit cube,  $k = \rho = C_p = 1$ . A temperature  $T_{max} = 1$  is prescribed at the bottom and  $T_{min} = 0$  at the top.

Same problem when these boundary conditions are now prescribed on the faces.

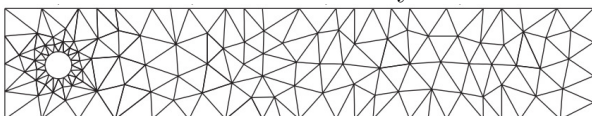
Prescribe an temperature field such that it is 1 everywhere in the domain but 2 inside a sphere of radius 0.1 centered at  $(0.66, 0.66, 0.66)$ . This time no diffusion takes place but we wish to advect the field using the velocity  $\vec{v} = (y, x, 0)$ . What is the highest resolution that is achievable on your computer?

### L.0.4 Triangular linear elements \*/\*\*

Redo the 2D advection-diffusion exercises with triangular elements. You will need to make a new icon array, and recompute the mass matrix and other matrices. The triangular elements are constructed by splitting square elements along the diagonal. See Section 5.3.7 for the basis functions and their derivatives. See Appendix E.0.5 for the calculations of the matrices.

### L.0.5 Triangular linear elements \*\*\*

Same exercise as above, with an additional task: run the benchmark presented in Section 12.2.11. For this you will need to generate a mesh such that nodes are placed on the perimeter of the cylinder and there is no node inside the cylinder:



You can build it by hand, or you can use an external mesher library (see Delaunay triangulation inside scipy). Vary the heat conduction coefficient to show the effect of diffusion on the obtained steady state temperature field.



## L.0.6 Diffusion of topography \*\*\*\*


In a 2D plane assign each node an initial topography  $h(x, y, t = 0)$  given by

$$h(x, y, t = 0) = h_0 \sin(\pi x / L_x) + \xi(x, y) \delta h$$

where  $L_x$  and  $L_y$  are the dimensions of the domain,  $h_0$  is the height of the orogen,  $\xi(x, y)$  is a random perturbation in  $[-1, 1]$  and  $\delta h$  is the amplitude of the perturbation.

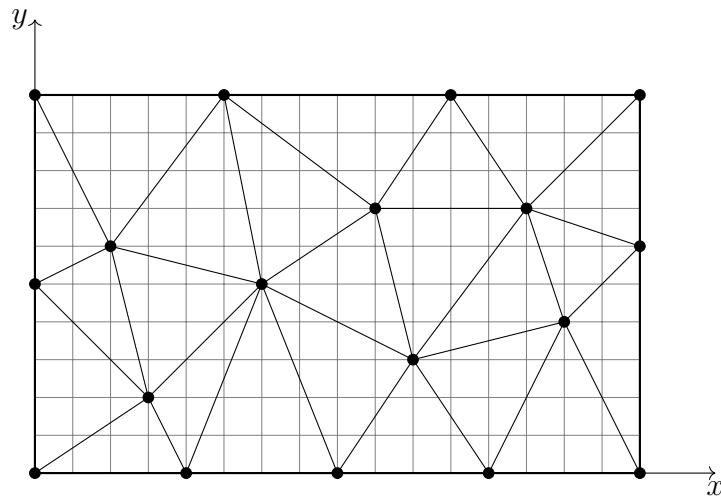
We wish to 'erode' the topography by means of a (nonlinear) diffusion law as in section 2.1.1 of Burov & Cloetingh [183].

1. What are the physical parameters needed to carry out this experiment? What are the appropriate boundary conditions? What is the steady state? What are the relevant time scales? How should we choose the time step?
2. Write a code which solves the linear diffusion equation until steady state is reached. Explore the effect of  $\delta h$ . Compute the slope  $\vec{\nabla} h$  inside each element and plot its time evolution.
3. Implement the nonlinear diffusion law and run the model once again.
4. If a source term is added to the diffusion equation it is in fact a vertical velocity ( $\partial h / \partial t$  has the dimensions of a velocity). Add a source term which generates uplift in a symmetric and asymmetric manner.

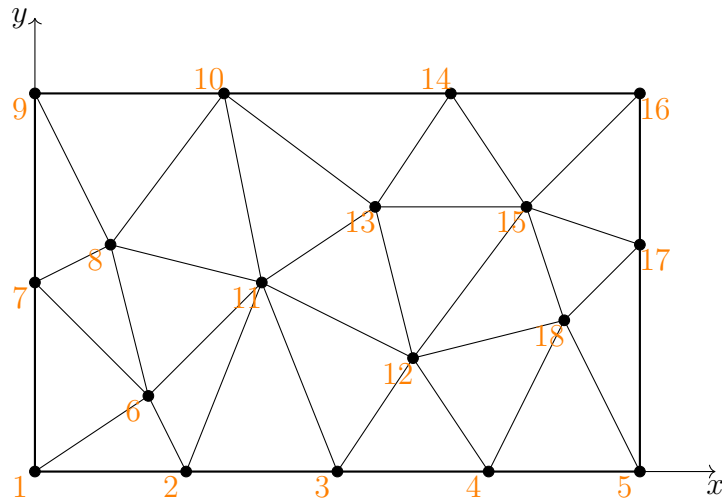
 Relevant Literature[1262] [1209], also check Appendix H.

## L.0.7 An example of a hand-built triangular mesh

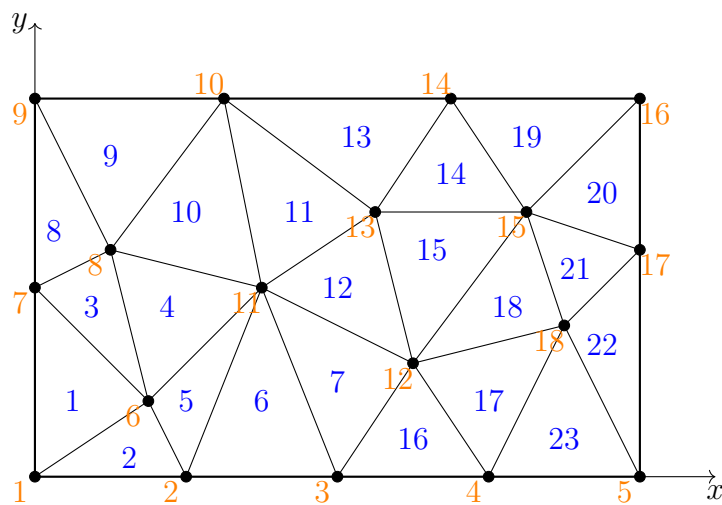
We start from a 8x5 domain which is tessellated as follows:



We can first label the nodes:



and then label the elements:



We can finally build the connectivity array by hand:

```

icon(1,1)=1
icon(2,1)=6
icon(3,1)=7
icon(1,2)=1
icon(2,2)=2
icon(3,2)=6
icon(1,3)=7
icon(2,3)=6
icon(3,3)=8
...
icon(1,12)=11
icon(2,12)=12
icon(3,12)=13
...
icon(1,19)=14
icon(2,19)=15
icon(3,19)=16

```

The labelling of nodes and elements above is done by a human so it starts at 1. When implementing this in python, you know what to do ...

## L.0.8 How to visualise data on a triangular mesh with Paraview

If arrays **x**,**y** contain the coordinates of the nodes, your connectivity array is called **icon**, and your mesh consists of **nel** elements and comprises **nnp** nodes, you can use the following code to generate a vtu file to be opened with Paraview. You also need a temperature array **T**.

Code at [https](https://raw.githubusercontent.com/cedrict/fieldstone/master/images/compgeo/mesh_visu.py):

[//raw.githubusercontent.com/cedrict/fieldstone/master/images/compgeo/mesh\\_visu.py](https://raw.githubusercontent.com/cedrict/fieldstone/master/images/compgeo/mesh_visu.py)

```
vtufile=open('mesh.vtu ', "w")
vtufile.write("<VTKFile type='UnstructuredGrid' version='0.1' byte_order='BigEndian'>
_\n")
vtufile.write("<UnstructuredGrid>\n")
vtufile.write("<Piece NumberOfPoints='_%5d' NumberOfCells='_%5d'>\n" %(nnp, nel))
vtufile.write("<Points>\n")
vtufile.write("<DataArray type='Float32' NumberOfComponents='3' Format='ascii'>\n")
for i in range(0, nnp):
 vtufile.write("%10e_%10e_%10e\n" %(x[i], y[i], 0.))
vtufile.write("</DataArray>\n")
vtufile.write("</Points>\n")
#vtufile.write("<CellData Scalars='scalars'>\n")
#vtufile.write("<DataArray type='Float32' Name='area' Format='ascii'> \n")
#for iel in range(0, nel):
vtufile.write("%10e\n" %(area[iel]))
#vtufile.write("</DataArray>\n")
#vtufile.write("</CellData>\n")
vtufile.write("<PointData Scalars='scalars'>\n")
vtufile.write("<DataArray type='Float32' Name='T' Format='ascii'>\n")
for i in range(0, nnp):
 vtufile.write("%10e\n" %T[i])
vtufile.write("</DataArray>\n")
vtufile.write("</PointData>\n")
vtufile.write("<Cells>\n")
vtufile.write("<DataArray type='Int32' Name='connectivity' Format='ascii'>\n")
for iel in range(0, nel):
 vtufile.write("%d_%d_%d\n" %(icon[0, iel], icon[1, iel], icon[2, iel]))
vtufile.write("</DataArray>\n")
vtufile.write("<DataArray type='Int32' Name='offsets' Format='ascii'>\n")
for iel in range(0, nel):
 vtufile.write("%d\n" %((iel+1)*3))
vtufile.write("</DataArray>\n")
vtufile.write("<DataArray type='Int32' Name='types' Format='ascii'>\n")
for iel in range(0, nel):
 vtufile.write("%d\n" %5)
vtufile.write("</DataArray>\n")
vtufile.write("</Cells>\n")
vtufile.write("</Piece>\n")
vtufile.write("</UnstructuredGrid>\n")
vtufile.write("</VTKFile>\n")
vtufile.close()
```

# Appendix M

## Using prisms in forward gravity modelling

*This appendix was written by Sverre Hassing as part of his Bachelor thesis. Although the final formula are definitely correct, the derivations below may still contain a typo.*

The derivations for prisms have been published in the early 50's [822]. However, to the best of our knowledge the full derivation has not been carried out in English in full detail. The derivations are based on those of Mader (1951) [822] and of Nagy *et al.* (2000) [922, 921]. Mader provided the derivations in some detail, while Nagy *et al.* interpreted the results in a more modern style.

### M.0.1 Basic formulas

The derivations for prisms are a lot more complicated than that for the point masses. We start with the following two integral equations which are integral part of the derivations:

$$\int \frac{x^2 dx}{x^2 + z^2} = x - z \arctan \frac{x}{z} \quad (\text{M.1})$$

$$\int \frac{dx}{\sqrt{x^2 + y^2 + z^2}} = \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) \quad (\text{M.2})$$

Another equation that will come back multiple times is of the form:

$$\int \frac{du}{(v^2 + w^2)\sqrt{u^2 + v^2 + w^2}} \quad (\text{M.3})$$

This can be solved with a trigonometric substitution, where  $u = \sqrt{v^2 + w^2} \tan \phi$ . This means that  $du = \frac{\sqrt{v^2 + w^2}}{\cos^2 \phi} d\phi$ .

$$\begin{aligned} \int \frac{du}{(v^2 + w^2)\sqrt{u^2 + v^2 + w^2}} &= \int \frac{\sqrt{v^2 + w^2}}{\cos^2 \phi} \frac{1}{(v^2 + w^2) \tan^2 \phi + v^2} \frac{d\phi}{\sqrt{v^2 + w^2 + (v^2 + w^2) \tan^2 \phi}} \\ &= \int \frac{\sqrt{v^2 + w^2}}{\cos^2 \phi} \frac{1}{v^2(\tan^2 \phi + 1) + w^2 \tan^2 \phi} \frac{d\phi}{\sqrt{(v^2 + w^2)(\tan^2 \phi + 1)}} \\ &= \int \frac{\sqrt{v^2 + w^2}}{\cos^2 \phi} \frac{1}{\frac{v^2 + w^2 \sin^2 \phi}{\cos^2 \phi} \frac{\sqrt{v^2 + w^2}}{\cos \phi}} \frac{d\phi}{\sqrt{(v^2 + w^2)(\tan^2 \phi + 1)}} \\ &= \int \frac{\cos \phi d\phi}{v^2 + w^2 \sin^2 \phi} \end{aligned} \quad (\text{M.4})$$

A second substitution is needed where  $t = \frac{w}{v} \sin \phi$  and  $dt = \frac{v}{w} \cos \phi d\phi$ :

$$\begin{aligned}
\int \frac{\cos \phi d\phi}{v^2 + w^2 \sin^2 \phi} &= \int \frac{v dt}{w(v^2 + v^2 t^2)} \\
&= \frac{1}{vw} \int \frac{dt}{1 + t^2} \\
&= \frac{1}{vw} \arctan t \\
&= \frac{1}{vw} \arctan \frac{w \sin \phi}{v}
\end{aligned} \tag{M.5}$$

Now the  $\sin \phi$  needs to be converted back to  $u, v, w$ . If it is known that  $\tan \phi = \frac{u}{\sqrt{v^2 + w^2}}$ , then it follows that  $\sin \phi = \frac{u}{\sqrt{u^2 + v^2 + w^2}}$ . Eq.(M.5) then becomes

$$\frac{1}{vw} \arctan \frac{w \sin \phi}{v} = \frac{1}{vw} \arctan \frac{uw}{v\sqrt{u^2 + v^2 + w^2}} \tag{M.6}$$

and finally

$$\boxed{\int \frac{du}{(v^2 + w^2)\sqrt{u^2 + v^2 + w^2}} = \frac{1}{vw} \arctan \frac{uw}{v\sqrt{u^2 + v^2 + w^2}}} \tag{M.7}$$

## M.0.2 The gravitational potential

Each prism is assumed to have constant density  $\rho$ . The gravitational potential is integrated over the whole volume of the prism:

$$U(P) = -\mathcal{G}\rho \underbrace{\int_{x_1}^{x_2} \int_{y_1}^{y_2} \int_{z_1}^{z_2} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}}}_I \tag{M.8}$$

In what follows we work out the exact form for the triple integral term. Elementary Eq. (M.2) can be applied to the integral for  $dx$  in Eq. (M.8).

$$\begin{aligned}
I &= \iiint \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \iint \left( \int \frac{dx}{\sqrt{x^2 + y^2 + z^2}} \right) dydz \\
&= \iint \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) dydz
\end{aligned} \tag{M.9}$$

We further proceed with the integration with respect to  $y$ . We define

$$\left. \begin{aligned} f &= \int \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) dz \\ f' &= \frac{y}{(x + \sqrt{x^2 + y^2 + z^2})\sqrt{x^2 + y^2 + z^2}} \end{aligned} \right| \begin{aligned} g' &= dy \\ g &= y \end{aligned}$$

and using  $\int fg' = fg - \int fg'$  we have

$$I = \underbrace{y \int \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) dz}_A - \underbrace{\iint \frac{y^2 dz}{\left( x + \sqrt{x^2 + y^2 + z^2} \right) \sqrt{x^2 + y^2 + z^2}}}_{B} dy \tag{M.10}$$

The calculation of  $I$  is then split into two large integrals denoted  $A$  and  $B$ , calculated in the following subsections. Note that we have not made use of the integral bounds yet.

## The calculation of $A$

The first step in calculating  $A$  is to carry out a similar partial integration as seen before.

$$\begin{aligned}
 & \left. \begin{array}{l} f = \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) \\ f' = \frac{\partial f}{\partial z} = \frac{z}{(x + \sqrt{x^2 + y^2 + z^2}) \sqrt{x^2 + y^2 + z^2}} \end{array} \right| \begin{array}{l} g' = dz \\ g = z \end{array} \\
 A &= y \left( z \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) - \int \frac{z^2 dz}{\left( x + \sqrt{x^2 + y^2 + z^2} \right) \sqrt{x^2 + y^2 + z^2}} \right) \\
 &= \underbrace{yz \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right)}_{A_0} - y \underbrace{\int \frac{z^2 dz}{\left( x + \sqrt{x^2 + y^2 + z^2} \right) \sqrt{x^2 + y^2 + z^2}}}_{A_1} \quad (M.11)
 \end{aligned}$$

We now focus on the  $A_1$  integral. We first multiply the numerator and denominator by  $-x + \sqrt{x^2 + y^2 + z^2}$ . The last step uses Eqs. (M.2), (M.7) and (M.1) respectively for each term.

$$\begin{aligned}
 A_1 &= \int \frac{z^2 dz}{\left( x + \sqrt{x^2 + y^2 + z^2} \right) \sqrt{x^2 + y^2 + z^2}} \frac{-x + \sqrt{x^2 + y^2 + z^2}}{-x + \sqrt{x^2 + y^2 + z^2}} \\
 &= \int \frac{(-xz^2 + z^2 \sqrt{x^2 + y^2 + z^2}) dz}{(x^2 + y^2 + z^2 - x^2) \sqrt{x^2 + y^2 + z^2}} \\
 &= \int \frac{-xz^2 dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} + \int \frac{z^2 \sqrt{x^2 + y^2 + z^2} dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} \\
 &= \int \frac{-x(z^2 + y^2 - y^2) dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} + \int \frac{z^2 dz}{y^2 + z^2} \\
 &= \int \frac{-x dz}{\sqrt{x^2 + y^2 + z^2}} + \int \frac{xy^2 dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} + \int \frac{z^2 dz}{y^2 + z^2} \\
 &= -x \int \frac{dz}{\sqrt{x^2 + y^2 + z^2}} + xy^2 \int \frac{dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} + \int \frac{z^2 dz}{y^2 + z^2} \\
 &= -x \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) + y \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} + z - y \arctan \frac{z}{y} \quad (M.12)
 \end{aligned}$$

This can be combined to get the final expression for  $A$ :

$$A = y \left( z \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) + x \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) - y \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} - z + y \arctan \frac{z}{y} \right) \quad (M.13)$$

The last two terms can be left out because they will cancel out when computing the integration boundaries from  $x_1$  to  $x_2$ , because these terms do not contain the variable  $x$ . Finally we arrive at the following expression for  $A$ :

$$A = yz \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) + xy \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) - y^2 \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} \quad (M.14)$$

## The calculation of $B$

The inner integral can be simplified similarly to how  $A_1$  was simplified in Eq. (M.12), by multiplying both numerator and denominator with  $-x + \sqrt{x^2 + y^2 + z^2}$ . The last step uses Eqs. (M.7) and (M.1):

$$\begin{aligned}
B &= - \int y^2 \int \frac{dz}{\left(x + \sqrt{x^2 + y^2 + z^2}\right) \sqrt{x^2 + y^2 + z^2}} \frac{-x + \sqrt{x^2 + y^2 + z^2}}{-x + \sqrt{x^2 + y^2 + z^2}} dy \\
&= - \int y^2 \int \frac{-x + \sqrt{x^2 + y^2 + z^2}}{(x^2 + y^2 + z^2 - x^2) \sqrt{x^2 + y^2 + z^2}} dz dy \\
&= - \int y^2 \left( - \int \frac{x dz}{(y^2 + z^2) \sqrt{x^2 + y^2 + z^2}} + \int \frac{dz}{y^2 + z^2} \right) dy \\
&= - \int y^2 \left( - \frac{1}{y} \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} + \frac{1}{y} \arctan \frac{z}{y} \right) dy \\
&= \int y \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} dy
\end{aligned} \tag{M.15}$$

Again the second term can be left out, because it does not contain the variable  $x$ . The next step is to apply a partial integration to  $B$ .

$$\begin{array}{c|c}
f = \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} & g' = y \\
\hline
f' = -xz \frac{\frac{1}{y^2 \sqrt{x^2 + y^2 + z^2}} + \frac{1}{(x^2 + y^2 + z^2)^{\frac{3}{2}}}}{\frac{x^2 z^2}{y^2 (x^2 + y^2 + z^2)} + 1} & g = \frac{y^2}{2}
\end{array}$$

$$B = \frac{y^2}{2} \arctan \frac{xz}{y \sqrt{x^2 + y^2 + z^2}} + \underbrace{\frac{xz}{2} \int y^2 \frac{\frac{1}{y^2 \sqrt{x^2 + y^2 + z^2}} + \frac{1}{(x^2 + y^2 + z^2)^{\frac{3}{2}}}}{\frac{x^2 z^2}{y^2 (x^2 + y^2 + z^2)} + 1} dy}_{B_1} \tag{M.16}$$

Let us finish by calculating the integral  $B_1$ :

$$\begin{aligned}
B_1 &= \frac{xz}{2} \int y^2 \frac{\frac{1}{y^2 \sqrt{x^2+y^2+z^2}} + \frac{1}{(x^2+y^2+z^2)^{3/2}}}{\frac{x^2 z^2}{y^2(x^2+y^2+z^2)} + 1} dy \\
&= \frac{xz}{2} \int y^2 \frac{\frac{x^2+y^2+z^2}{y^2(x^2+y^2+z^2)^{3/2}} + \frac{y^2}{y^2(x^2+y^2+z^2)^{3/2}}}{\frac{x^2 z^2 + y^2(x^2+y^2+z^2)}{y^2(x^2+y^2+z^2)}} dy \\
&= \frac{xz}{2} \int y^2 \frac{\frac{x^2+2y^2+z^2}{y^2(x^2+y^2+z^2)^{3/2}}}{\frac{x^2 z^2 + y^2(x^2+y^2+z^2)}{y^2(x^2+y^2+z^2)}} dy \\
&= \frac{xz}{2} \int y^2 \frac{x^2 + 2y^2 + z^2}{\sqrt{x^2 + y^2 + z^2}(x^2 z^2 + y^2(x^2 + y^2 + z^2))} dy \\
&= \frac{xz}{2} \int y^2 \frac{x^2 + 2y^2 + z^2}{\sqrt{x^2 + y^2 + z^2}(x^2 + y^2)(z^2 + y^2)} dy \\
&= \frac{xz}{2} \left( \int \frac{2dy}{\sqrt{x^2 + y^2 + z^2}} + \int \frac{-(x^2 + z^2)y^2 - 2x^2 z^2}{(x^2 + y^2)(y^2 + z^2)\sqrt{x^2 + y^2 + z^2}} dy \right) \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{xz}{2} \int \frac{(x^2 + z^2)y^2 + 2x^2 z^2}{(x^2 + y^2)(y^2 + z^2)\sqrt{x^2 + y^2 + z^2}} dy \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{xz}{2} \int \frac{x^2 y^2 + y^2 z^2 + 2x^2 z^2}{(x^2 + y^2)(y^2 + z^2)\sqrt{x^2 + y^2 + z^2}} dy \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{xz}{2} \int \frac{x^2(y^2 + z^2) + z^2(x^2 + y^2)}{(x^2 + y^2)(y^2 + z^2)\sqrt{x^2 + y^2 + z^2}} dy \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{xz}{2} \int \frac{x^2}{(x^2 + y^2)\sqrt{x^2 + y^2 + z^2}} dy - \frac{xz}{2} \int \frac{z^2}{(y^2 + z^2)\sqrt{x^2 + y^2 + z^2}} dy \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{xz}{2} \frac{x^2 \arctan \frac{yz}{x\sqrt{x^2+y^2+z^2}}}{xz} - \frac{xz}{2} \frac{z^2 \arctan \frac{xy}{z\sqrt{x^2+y^2+z^2}}}{xz} \\
&= xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) - \frac{x^2}{2} \arctan \frac{yz}{x\sqrt{x^2+y^2+z^2}} - \frac{z^2}{2} \arctan \frac{xy}{z\sqrt{x^2+y^2+z^2}} \quad (M.17)
\end{aligned}$$

This can be combined to get the full expression for  $B$ :

$$B = xz \ln \left( \sqrt{x^2 + y^2 + z^2} + y \right) - \frac{x^2}{2} \arctan \frac{zy}{x\sqrt{x^2 + y^2 + z^2}} + \frac{y^2}{2} \arctan \frac{xz}{y\sqrt{x^2 + y^2 + z^2}} - \frac{z^2}{2} \arctan \frac{xy}{x\sqrt{z^2 + y^2 + z^2}}$$

### Combining $A$ and $B$

Now  $A$  and  $B$  can be combined to get the expression of  $I$

$$\begin{aligned}
I &= A + B \\
&= yz \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) + xy \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) - y^2 \arctan \frac{xz}{y\sqrt{x^2 + y^2 + z^2}} \\
&+ xz \ln \left( \sqrt{x^2 + y^2 + z^2} + y \right) - \frac{x^2}{2} \arctan \frac{zy}{x\sqrt{x^2 + y^2 + z^2}} + \frac{y^2}{2} \arctan \frac{xz}{y\sqrt{x^2 + y^2 + z^2}} - \frac{z^2}{2} \arctan \frac{xy}{x\sqrt{z^2 + y^2 + z^2}} \\
&= yz \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) + xy \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) + xz \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) \\
&- \frac{x^2}{2} \arctan \frac{zy}{x\sqrt{x^2 + y^2 + z^2}} - \frac{y^2}{2} \arctan \frac{xz}{y\sqrt{x^2 + y^2 + z^2}} - \frac{z^2}{2} \arctan \frac{xy}{x\sqrt{z^2 + y^2 + z^2}}
\end{aligned}$$



The boundaries for the volume from Eq. (??) need to be applied to the result of the integration. The boundary conditions are computed by plugging the upper value into the equation and subtracting the equation with the lower value plugged in. When the upper and lower values are respectively  $x_2$  and  $x_1$  for some function  $f(x)$ , this is  $f(x_2) - f(x_1)$ . This can be represented more efficiently with a summation over the subscript. Something needs to be added to still keep the subtraction in there. This can be done by adding a factor of  $-1^i$ , where  $i$  is the summation index. This will be positive when  $i$  is even and negative when  $i$  is odd. The new way of showing the result would be  $\sum_{i=1}^2 -1^i f(x_i)$ . This is especially useful when there are three different integration boundaries to resolve.  $r$  will be used instead of  $\sqrt{x^2 + y^2 + z^2}$ .

$$\begin{aligned}
I &= \left\| \left\| yz \ln(x+r) + xy \ln(z+r) + xz \ln(y+r) - \frac{x^2}{2} \arctan \frac{zy}{xr} - \frac{y^2}{2} \arctan \frac{xz}{yr} - \frac{z^2}{2} \arctan \frac{xy}{xr} \right\|_{x_1}^{x_2} \right\|_{y_1}^{y_2} \Big|_{z_1}^{z_2} \\
&= \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( y_j z_k \ln(x_i + r_{ijk}) + x_i y_j \ln(z_k + r_{ijk}) + x_i z_k \ln(y_j + r_{ijk}) \right. \\
&\quad \left. - \frac{x_i^2}{2} \arctan \frac{z_k y_j}{x_i r_{ijk}} - \frac{y_j^2}{2} \arctan \frac{x_i z_k}{y_j r_{ijk}} - \frac{z_k^2}{2} \arctan \frac{x_i y_j}{x_i r_{ijk}} \right) \quad (M.19)
\end{aligned}$$

There is probably a mistake in eq above and below, last term, most likely should contain zk in denominator?

Finally,

$$\begin{aligned}
U(\vec{r}) &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( y_j z_k \ln(x_i + r_{ijk}) + x_i y_j \ln(z_k + r_{ijk}) + x_i z_k \ln(y_j + r_{ijk}) \right. \\
&\quad \left. - \frac{x_i^2}{2} \arctan \frac{z_k y_j}{x_i r_{ijk}} - \frac{y_j^2}{2} \arctan \frac{x_i z_k}{y_j r_{ijk}} - \frac{z_k^2}{2} \arctan \frac{x_i y_j}{x_i r_{ijk}} \right) \quad (M.20)
\end{aligned}$$

### M.0.3 The gravity vector $\vec{g}$

In 3D Cartesian coordinates the gravity vector is expressed as

$$\vec{g} = -\vec{\nabla}U = \begin{pmatrix} -\frac{\partial U}{\partial x} \\ -\frac{\partial U}{\partial y} \\ -\frac{\partial U}{\partial z} \end{pmatrix} \quad (M.21)$$

The easiest way to calculate this is by including the partial derivatives in the original integral (M.8).

$$\begin{aligned}
I_x(\vec{r}) &= \iiint \frac{\partial}{\partial x} \frac{dx dy dz}{\sqrt{x^2 + y^2 + z^2}} \\
&= - \iiint \frac{x dx dy dz}{(\sqrt{x^2 + y^2 + z^2})^3} \\
&= \iint \frac{dy dz}{\sqrt{x^2 + y^2 + z^2}} \quad (M.22)
\end{aligned}$$

The integral (M.2) can be used, followed by the calculation of  $A$  as seen in Section M.0.2 without the multiplication with  $y$ :

$$\begin{aligned} I_x(\vec{r}) &= \int \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right) dz \\ &= z \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right) + y \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right) - x \arctan \frac{yz}{x\sqrt{x^2 + y^2 + z^2}} \end{aligned} \quad (\text{M.23})$$

The integration boundaries can be applied. Multiplication with  $\mathcal{G}$  and  $\rho$  is the final step in deriving the element of the gravity vector component ( $g_x$ ).

$$g_x = \mathcal{G}\rho \sum_{i,j,k=1}^2 (-1)^{i+j+k} \left( z_k \ln \left( y_j + \sqrt{x_i^2 + y_j^2 + z_k^2} \right) + y_j \ln \left( z_j + \sqrt{x_i^2 + y_j^2 + z_k^2} \right) - x_i \arctan \frac{y_j z_k}{x_i \sqrt{x_i^2 + y_j^2 + z_k^2}} \right)$$

The same can be done for the  $y$ - and  $z$ -components and in the end we obtain

$$\begin{aligned} g_x &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( z_k \ln (y_j + r_{ijk}) + y_j \ln (z_j + r_{ijk}) - x_i \arctan \frac{y_j z_k}{x_i r_{ijk}} \right) \\ g_y &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( z_k \ln (x_i + r_{ijk}) + x_i \ln (z_j + r_{ijk}) - y_j \arctan \frac{x_i z_k}{y_j r_{ijk}} \right) \\ g_z &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( x_i \ln (y_j + r_{ijk}) + y_j \ln (x_i + r_{ijk}) - z_k \arctan \frac{x_i y_j}{z_k r_{ijk}} \right) \end{aligned}$$

These equations can be found in various other papers such as Eq. (6) in Heck and Seitz, 2007 [556], Eqs. (8,11,12) in Nagy *et al.*, 2000 [922] (note that there is a mistake there that is later fixed in [921]), appendix A in Couder-Castaneda *et al.*, 2015 [279] and the derivation between (14) and (15) in Mader, 1951 [822].

## M.0.4 The gravity gradient tensor

The different elements of the gravity gradient tensor can be determined by partially differentiating each component of the gravity vector with respect to each space coordinate. As will be shown later,  $\mathbf{T}$  should be a symmetric matrix and its trace should equal zero.

### The diagonal terms

$$T_{xx} = \frac{\partial}{\partial x} g_x = -\frac{\partial^2}{\partial x^2} U(\vec{r}) = \mathcal{G}\rho \frac{\partial^2}{\partial x^2} (-I(\vec{r})) \quad (\text{M.24})$$

$$\begin{aligned} I_{xx}(\vec{r}) &= \iiint \frac{\partial^2}{\partial x^2} \frac{dx dy dz}{\sqrt{x^2 + y^2 + z^2}} \\ &= \iint \frac{\partial}{\partial x} \frac{dy dz}{\sqrt{x^2 + y^2 + z^2}} \\ &= - \iint \frac{x dy dz}{\sqrt{x^2 + y^2 + z^2}^3} \end{aligned} \quad (\text{M.25})$$

A trigonometric substitution is applied to solve this integral. This uses  $y = \sqrt{x^2 + z^2} \tan \phi$  and  $dy = \frac{\sqrt{x^2 + z^2}}{\cos^2 \phi} d\phi$ .

$$\begin{aligned}
I_{xx} &= -x \iint \frac{\sqrt{x^2 + y^2}}{\cos^2 \phi} \frac{d\phi dz}{\sqrt{(x^2 + z^2) \tan^2 \phi + x^2 + z^2}^3} \\
&= -x \iint \frac{\sqrt{x^2 + y^2}}{\cos^2 \phi} \frac{d\phi dz}{\sqrt{(x^2 + z^2)(\tan^2 \phi + 1)}^3} \\
&= -x \iint \frac{\sqrt{x^2 + y^2}}{\cos^2 \phi} \frac{d\phi dz}{\left( \frac{\sqrt{x^2 + z^2}}{\cos \phi} \right)^3} \\
&= -x \int \frac{1}{x^2 + z^2} \int \cos \phi d\phi dz \\
&= -x \int \frac{1}{x^2 + z^2} \sin \phi dz
\end{aligned} \tag{M.26}$$

Now the substitution needs to be undone. If  $\tan \phi = \frac{y}{\sqrt{x^2 + z^2}}$ , then  $\sin \phi = \frac{y}{\sqrt{x^2 + y^2 + z^2}}$  and then

$$I_{xx} = -xy \int \frac{dz}{(x^2 + z^2) \sqrt{x^2 + y^2 + z^2}} \tag{M.27}$$

This can be solved by applying equation (M.7).

$$\begin{aligned}
I_{xx} &= -\frac{xy}{xy} \arctan \frac{yz}{x \sqrt{x^2 + y^2 + z^2}} \\
&= -\arctan \frac{yz}{x \sqrt{x^2 + y^2 + z^2}}
\end{aligned}$$

The tensor element  $T_{xx}$  is then formulated as follows (the other elements of the diagonal are found by cyclic permutation of  $x$ ,  $y$  and  $z$ ):

$$\begin{aligned}
T_{xx} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( -\arctan \frac{y_j z_k}{x_i r_{ijk}} \right) \\
T_{yy} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( -\arctan \frac{x_i z_k}{y_j r_{ijk}} \right) \\
T_{zz} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} \left( -\arctan \frac{x_i y_j}{z_k r_{ijk}} \right)
\end{aligned}$$

### The off-diagonal terms of the tensor

The other elements are easier to calculate, because the partial derivatives cancel out the integrals:

$$\begin{aligned}
I_{xy} &= \iiint \frac{\partial^2}{\partial x \partial y} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \iint \frac{\partial}{\partial y} \frac{dydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \int \frac{dz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \ln \left( z + \sqrt{x^2 + y^2 + z^2} \right)
\end{aligned} \tag{M.28}$$

$$\begin{aligned}
I_{xz} &= \iiint \frac{\partial^2}{\partial x \partial z} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \iint \frac{\partial}{\partial z} \frac{dydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \int \frac{dy}{\sqrt{x^2 + y^2 + z^2}} \\
&= \ln \left( y + \sqrt{x^2 + y^2 + z^2} \right)
\end{aligned} \tag{M.29}$$

$$\begin{aligned}
I_{yz} &= \iiint \frac{\partial^2}{\partial y \partial z} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \iint \frac{\partial}{\partial z} \frac{dxdz}{\sqrt{x^2 + y^2 + z^2}} \\
&= \int \frac{dx}{\sqrt{x^2 + y^2 + z^2}} \\
&= \ln \left( x + \sqrt{x^2 + y^2 + z^2} \right)
\end{aligned} \tag{M.30}$$

From these calculations it should be obvious why  $\mathbf{T}$  is a symmetric tensor. When applying the second partial derivatives, their order does not matter:

$$I_{xy} = \iiint \frac{\partial^2}{\partial x \partial y} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} = \iiint \frac{\partial^2}{\partial y \partial x} \frac{dxdydz}{\sqrt{x^2 + y^2 + z^2}} = I_{yx} \tag{M.31}$$

The tensor elements following from this are:

$$\begin{aligned}
T_{xy} = T_{yx} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} (\ln(z_k + r_{ijk})) \\
T_{xz} = T_{zx} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} (\ln(y_j + r_{ijk})) \\
T_{yz} = T_{zy} &= \mathcal{G}\rho \sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 (-1)^{i+j+k} (\ln(x_i + r_{ijk}))
\end{aligned}$$

## M.0.5 Revisiting Poisson's equation

The gravitational potential Poisson equation is  $\nabla^2 U = 4\pi\mathcal{G}\rho$ . This can and should be verified for the derived equations for prisms. Inside the prism, the density has an assigned constant value. Outside of the prism, the density is zero, so the result is  $\nabla^2 U = 0$ . These cases will be treated separately.

### Outside the prism

$\nabla^2 U = 0$  can be written  $\frac{\partial^2 U}{\partial x^2} + \frac{\partial^2 U}{\partial y^2} + \frac{\partial^2 U}{\partial z^2} = 0$ . which is the trace of  $\mathbf{T}$ . We first add the terms and then the boundary conditions are applied.

We will need the formula to add arctangents together:

$$\arctan a + \arctan b = \arctan \frac{a+b}{1-ab} \quad (\text{M.32})$$

We start by adding the terms  $T_{xx}$  and  $T_{yy}$  together:

$$\begin{aligned} T_{xx} + T_{yy} &= (-1)^{i+j+k} \arctan \frac{yz}{xr} + (-1)^{i+j+k} \arctan \frac{xz}{yr} \\ &= (-1)^{i+j+k} \arctan \frac{\frac{yz}{xr} + \frac{xz}{yr}}{1 - \frac{yz}{xr} \frac{xz}{yr}} \\ &= (-1)^{i+j+k} \arctan \frac{\frac{y^2 z}{xyr} + \frac{x^2 z}{xyr}}{1 - \frac{xyz^2}{xyr^2}} \\ &= (-1)^{i+j+k} \arctan \frac{\frac{z(x^2+y^2)}{xyr}}{\frac{xy(r^2-z^2)}{xyr^2}} \\ &= (-1)^{i+j+k} \arctan \frac{xyzr^2(x^2+y^2)}{x^2y^2r(r^2-z^2)} \\ &= (-1)^{i+j+k} \arctan \frac{zr(x^2+y^2)}{xy(x^2+y^2+z^2-z^2)} \\ &= (-1)^{i+j+k} \arctan \frac{zr}{xy} \end{aligned} \quad (\text{M.33})$$

By considering a right triangle with sides 1 and  $x$ , it easy to prove that:

$$\arctan x + \arctan \frac{1}{x} = \frac{\pi}{2} \quad (\text{M.34})$$

This can be used to transform the arctan to one that is similar to  $T_{zz}$ .

$$T_{xx} + T_{yy} = (-1)^{i+j+k} \arctan \frac{zr}{xy} = (-1)^{i+j+k} \left( \frac{\pi}{2} - \arctan \frac{xy}{zr} \right) \quad (\text{M.35})$$

The last step is to add the term  $T_{zz}$ :

$$\nabla^2 U = T_{xx} + T_{yy} + T_{zz} = (-1)^{i+j+k} \left( \frac{\pi}{2} - \arctan \frac{xy}{zr} + \arctan \frac{xy}{zr} \right) = (-1)^{i+j+k} \frac{\pi}{2} \quad (\text{M.36})$$

The end result is a single value. When the boundary conditions are applied this single value will be subtracted from itself resulting in zero, so  $\nabla^2 U(\vec{r}) = 0$ .

### Inside the prism

We can simply put the observation point at the centre of the prism. The coordinates of the prism are now such that  $-x_1 = x_2$ ,  $-y_1 = y_2$  and  $-z_1 = z_2$ . All eight terms for these conditions results give  $\frac{\pi}{2}$ , so the result is:

$$\nabla^2 U = \mathcal{G}\rho 8 \frac{\pi}{2} = 4\pi\mathcal{G}\rho \quad (\text{M.37})$$

## M.0.6 Better numerical stability

Heck and Seitz (2007) [556] modify the standard formulae for the prism to get a better numerical stability in the logarithms. This is done by dividing the inside of the logs by an extra factor:

$$\ln(z_k + r_{ijk}) \rightarrow \ln \frac{z_k + r_{ijk}}{\sqrt{x_i^2 + y_j^2}}$$

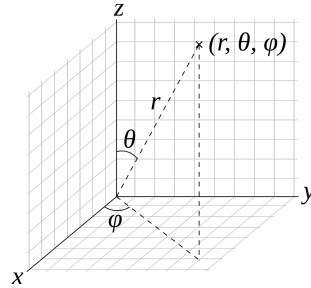
This extra factor disappears when applying the boundary conditions in the  $z$  direction, so that the results remain identical:

$$\begin{aligned} \left| \ln \frac{z_k + r_{ijk}}{\sqrt{x_i^2 + y_j^2}} \right|_{z_1}^{z_2} &= \left| \ln(z_k + r_{ijk}) - \ln \sqrt{x_i^2 + y_j^2} \right|_{z_1}^{z_2} \\ &= \ln(z_2 + r_{ijk}) - \ln \sqrt{x_i^2 + y_j^2} - \ln(z_1 + r_{ijk}) + \ln \sqrt{x_i^2 + y_j^2} \\ &= \ln(z_2 + r_{ijk}) - \ln(z_1 + r_{ijk}) \end{aligned} \tag{M.38}$$

# Appendix N

## Solutions to exercises of GEO3-1313

### N.0.1 Problem 1



$$\begin{aligned}
 I &= \frac{1}{3}(I_x + I_y + I_z) \\
 &= \frac{1}{3} \iiint_V \rho(r)(y^2 + z^2)dV + \frac{1}{3} \iiint_V \rho(r)(x^2 + z^2)dV + \frac{1}{3} \iiint_V \rho(r)(x^2 + y^2)dV \\
 &= \frac{1}{3} \iiint_V \rho(r)(2x^2 + 2y^2 + 2z^2)dV \\
 &= \frac{2}{3} \iiint_V \rho(r)r^2 dV \\
 &= \frac{2}{3} \iiint_V \rho(r)r^2 r^2 \sin \theta dr d\theta d\phi \\
 &= \frac{2}{3} \int_0^R \int_0^\pi \int_0^{2\pi} \rho(r)r^2 r^2 \sin \theta dr d\theta d\phi \\
 &= \frac{2}{3} \left( \int_0^\pi \sin \theta d\theta \right) \left( \int_0^{2\pi} d\phi \right) \int_0^R \rho(r)r^4 dr \\
 &= \frac{8\pi}{3} \int_0^R \rho(r)r^4 dr
 \end{aligned} \tag{N.1}$$

Alternative solution:  $I$  is evaluated for the special case where the rotation axis is the  $z$ -axis, where  $d = r \sin \theta$ . Substitution in  $I = \int_V \rho d^2 dV$  yields

$$\begin{aligned}
 I &= \iiint_V \rho(r)(r^2 \sin^2 \theta)dV \\
 &= \int_0^R \int_0^\pi \int_0^{2\pi} \rho(r)(r^2 \sin^2 \theta)r^2 \sin \theta dr d\theta d\phi \\
 &= \left( \int_0^\pi \sin^3 \theta d\theta \right) \left( \int_0^{2\pi} d\phi \right) \int_0^R \rho(r)r^4 dr
 \end{aligned}$$

By writing  $\sin^3 \theta = \sin \theta (1 - \cos^2 \theta)$  we can operate a change of variables  $u = \cos \theta$ , and we arrive at  $\int \sin^3 \theta d\theta = 4/3$  (skipping a few steps here) and we obtain the expected result.



## N.0.2 Problem 2

For a uniform sphere, we have  $\rho(r) = \rho_0$ . Then

$$\begin{aligned} I &= \frac{8\pi}{3} \int_0^R \rho_0 r^4 dr \\ &= \frac{8\pi}{3} \rho_0 \int_0^R r^4 dr \\ &= \frac{8\pi}{15} \rho_0 R^5 \int_0^R r^4 dr \end{aligned}$$

The mass of the sphere is

$$\begin{aligned} M &= \iiint_V \rho(r) dV \\ &= \rho_0 \int_0^R \int_0^\pi \int_0^{2\pi} r^2 \sin \theta dr d\theta d\phi \\ &= \rho_0 \frac{4}{3} \pi R^3 \end{aligned} \tag{N.2}$$

In the end we get

$$I = \frac{2}{5} M R^2$$

When all the mass is concentrated in the center, then  $\rho(r) = \rho_0 \delta(r)$  where  $\delta$  is the Dirac delta function<sup>1</sup>. Then

$$\begin{aligned} I &= \frac{8\pi}{3} \int_0^R \rho_0 \delta(r) r^4 dr \\ &= \frac{8\pi}{3} \rho_0 \int_0^R \delta(r) r^4 dr \\ &= 0 \end{aligned} \tag{N.3}$$

When all the mass is concentrated in a shell of zero thickness of radius  $R$ , then  $\rho(r) = \rho_0 \delta(r - R)$ , so

$$\begin{aligned} I &= \frac{8\pi}{3} \int_0^R \rho_0 \delta(r - R) r^4 dr \\ &= \frac{8\pi}{3} \rho_0 \int_0^R \delta(r - R) r^4 dr \\ &= \frac{8\pi}{3} \rho_0 R^4 \end{aligned} \tag{N.4}$$

Conversely, its mass is

$$\begin{aligned} M &= \iiint_V \rho(r) dV \\ &= \rho_0 \int_0^R \int_0^\pi \int_0^{2\pi} \delta(r - R) r^2 \sin \theta dr d\theta d\phi \\ &= \rho_0 4\pi R^2 \end{aligned} \tag{N.5}$$

and then

$$I = \frac{2}{3} M R^2$$

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Dirac\\_delta\\_function](https://en.wikipedia.org/wiki/Dirac_delta_function)

### N.0.3 Problem 3

$$\begin{aligned}
 \langle \rho \rangle &= \frac{1}{V} \iiint_V \rho dV \\
 &= \frac{1}{\frac{4}{3}\pi R^3} \iiint_V \rho(r) r^2 \sin \theta dr d\theta d\phi \\
 &= \frac{1}{\frac{4}{3}\pi R^3} 4\pi \int_0^R \rho(r) r^2 dr \\
 &= \frac{3}{R^3} \int_0^R \rho(r) r^2 dr
 \end{aligned} \tag{N.6}$$

We then turn to the moment of inertia:

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr \tag{N.7}$$

$$= fMR^2 \tag{N.8}$$

where

$$M = \iiint_V \rho(r) dV = V \underbrace{\frac{1}{V} \iiint_V \rho(r) dV}_{\langle \rho \rangle} = \frac{4}{3}\pi R^3 \langle \rho \rangle$$

We then insert this expression of  $M$  in Eq. (N.8):

$$I = f \frac{4}{3}\pi R^3 \langle \rho \rangle R^2 = f \frac{4}{3}\pi R^5 \langle \rho \rangle$$

Equating this to Eq. N.7 yields

$$\frac{8\pi}{3} \int_0^R \rho(r) r^4 dr = f \frac{4}{3}\pi R^5 \langle \rho \rangle$$

or,

$$f \langle \rho \rangle R^5 = 2 \int_0^R \rho(r) r^4 dr \tag{N.9}$$

We now make use of the expression for the density:

$$\rho(r) = \begin{cases} \rho_c & 0 \leq r \leq R_c \\ \rho_m & R_c \leq r \leq R \end{cases}$$

Then Eq. (N.6) yields

$$\begin{aligned}
 \langle \rho \rangle &= \frac{3}{R^3} \int_0^R \rho(r) r^2 dr \\
 &= \frac{3}{R^3} \left( \int_0^{R_c} \rho_c r^2 dr + \int_{R_c}^R \rho_m r^2 dr \right) \\
 &= \frac{3}{R^3} \left( \frac{R_c^3}{3} \rho_c + \frac{1}{3} (R^3 - R_c^3) \rho_m \right) \\
 &= \frac{R_c^3}{R^3} \rho_c + \left( 1 - \frac{R_c^3}{R^3} \right) \rho_m \\
 &= \left( \frac{R_c}{R} \right)^3 \rho_c + \left[ 1 - \left( \frac{R_c}{R} \right)^3 \right] \rho_m
 \end{aligned} \tag{N.10}$$

We know  $\rho_m$ , but not  $\rho_c$ , so we write:

$$\rho_c = \rho_m \left[ 1 + \left( \frac{R}{R_c} \right)^3 \left( \frac{\langle \rho \rangle}{\rho_m} - 1 \right) \right] \quad (\text{N.11})$$

We now turn to Eq. (N.9). Since

$$\begin{aligned} 2 \int_0^R \rho(r) r^4 dr &= 2 \int_0^{R_c} \rho_c r^4 dr + 2 \int_{R_c}^R \rho_m r^4 dr \\ &= \frac{2}{5} [R_c^5 \rho_c + (R^5 - R_c^5) \rho_m] \end{aligned}$$

then

$$f \langle \rho \rangle R^5 = \frac{2}{5} [R_c^5 \rho_c + (R^5 - R_c^5) \rho_m]$$

or,

$$\begin{aligned} \frac{5f \langle \rho \rangle R^5}{2\rho_m} &= R_c^5 \frac{\rho_c}{\rho_m} + (R^5 - R_c^5) \\ \frac{5f \langle \rho \rangle R^5}{2\rho_m} &= R_c^5 \left( \frac{\rho_c}{\rho_m} - 1 \right) + R^5 \end{aligned}$$

Now dividing by  $R^5$  on each side:

$$\frac{5f \langle \rho \rangle}{2\rho_m} = \left( \frac{R_c}{R} \right)^5 \left( \frac{\rho_c}{\rho_m} - 1 \right) + 1$$

Using Eq. (N.11), we can write

$$\frac{\rho_c}{\rho_m} - 1 = \left( \frac{R}{R_c} \right)^3 \left( \frac{\langle \rho \rangle}{\rho_m} - 1 \right)$$

so finally

$$\frac{5f \langle \rho \rangle}{2\rho_m} - 1 = \left( \frac{R_c}{R} \right)^5 \left( \frac{R}{R_c} \right)^3 \left( \frac{\langle \rho \rangle}{\rho_m} - 1 \right)$$

or,

$$\frac{5f \langle \rho \rangle}{2\rho_m} - 1 = \left( \frac{R_c}{R} \right)^2 \left( \frac{\langle \rho \rangle}{\rho_m} - 1 \right)$$

and then

$$\left( \frac{R_c}{R} \right) = \left( \frac{\frac{5f \langle \rho \rangle}{2\rho_m} - 1}{\left( \frac{\langle \rho \rangle}{\rho_m} - 1 \right)} \right)^{1/2}$$

#### N.0.4 Problem 4

Assume a uniform mantle  $\rho_m$  and core  $\rho_c$ . For the total mass we have

$$M = \int_V \rho dV = 4\pi \int_0^R \rho(r) r^2 dr = \frac{4\pi}{3} R_c^3 \rho_c + \frac{4\pi}{3} (R^3 - R_c^3) \rho_m$$

For the moment of inertia we have,

$$I = \frac{8\pi}{3} \int_0^R \rho(r) r^4 dr = \frac{8\pi}{15} R_c^5 \rho_c + \frac{8\pi}{15} (R^5 - R_c^5) \rho_m$$

In matrix form the above equations become:

$$\begin{pmatrix} \frac{4\pi}{3} R_c^3 & \frac{4\pi}{3} (R^3 - R_c^3) \\ \frac{8\pi}{15} R_c^5 & \frac{8\pi}{15} (R^5 - R_c^5) \end{pmatrix} \cdot \begin{pmatrix} \rho_c \\ \rho_m \end{pmatrix} = \begin{pmatrix} M \\ I \end{pmatrix}$$

We use Cramers rule

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \quad \Rightarrow \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \frac{1}{\Delta} \begin{pmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

In our case the determinant of the matrix is

$$\Delta = \frac{32\pi^2}{45} [R_c^3 (R^5 - R_c^3) - R_c^5 (R^3 - R_c^3)]$$

#### N.0.5 Problem 5

skipped

#### N.0.6 Problem 6

see lecture notes

#### N.0.7 Problem 7

$$g = \frac{GM}{R^2} = \frac{6.67e-11 * 5.97e24}{6371000^2} \simeq 9.8$$

#### N.0.8 Problem 8

$$p_{CMB} = \int \rho g dz \simeq \rho_0 g_0 (R_{Earth} - R_{CMB}) \simeq 127.5 GPa$$

### N.0.9 Problem 9

$$U(\vec{r}) = -\frac{\mathcal{G}m_1}{|\vec{r}_1 - \vec{r}|} = -\frac{\mathcal{G}m_1}{\sqrt{(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2}}$$

$$\vec{\nabla}U = \begin{pmatrix} \partial_x U \\ \partial_y U \\ \partial_z U \end{pmatrix}$$

We have

$$\begin{aligned} \partial_x U &= -\mathcal{G}m_1 \frac{\partial}{\partial x} \frac{1}{\sqrt{(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2}} \\ &= -\mathcal{G}m_1 \cdot -\frac{1}{2} \frac{-2(x_1 - x)}{[(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2]^{3/2}} \\ &= \mathcal{G}m_1 \frac{1}{[(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2]} \frac{-(x_1 - x)}{\sqrt{(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2}} \\ &= \mathcal{G}m_1 \frac{1}{|\vec{r}_1 - \vec{r}|^2} \frac{-(x_1 - x)}{\sqrt{(x_1 - x)^2 + (y_1 - y)^2 + (z_1 - z)^2}} \end{aligned}$$

We repeat this operation for  $\partial_y$  and  $\partial_z$  and finally:

$$-\vec{\nabla}U = \begin{pmatrix} \partial_x U \\ \partial_y U \\ \partial_z U \end{pmatrix} = \mathcal{G}m_1 \frac{1}{|\vec{r}_1 - \vec{r}|^2} \begin{pmatrix} \frac{(x_1 - x)}{|\vec{r}_1 - \vec{r}|} \\ \frac{(y_1 - y)}{|\vec{r}_1 - \vec{r}|} \\ \frac{(z_1 - z)}{|\vec{r}_1 - \vec{r}|} \end{pmatrix} = \frac{\mathcal{G}m_1}{|\vec{r}_1 - \vec{r}|^2} \vec{e}_{\vec{r}_1 \vec{r}} = \vec{g}$$

### N.0.10 Problem 10

$$\int_V \Delta U dV = \int_V 4\pi\mathcal{G}\rho dV = 4\pi\mathcal{G} \int_V M\delta(\vec{r} - \vec{r}_0) dV = 4\pi\mathcal{G}M$$

$$\int_V \Delta U dV = \int_V \vec{\nabla}^2 U dV = \int_V \vec{\nabla} \cdot \vec{\nabla} U dV = \int_\Gamma \vec{\nabla} U \cdot \vec{n} dS = \int_\Gamma (-\vec{g}) \cdot \vec{n} dS = \int_\Gamma g dS = 4\pi r^2 g$$

Note that we have  $\vec{g}$  which is pointing towards the center and therefore has the opposite direction to  $\vec{n}$  which is normal to the shell surface so that  $(-\vec{g}) \cdot \vec{n} = g$ . Also the integral on  $\Gamma$  is at constant radius so  $g(r)$  can be taken out of the integral.

Finally we obtain

$$g = \frac{\mathcal{G}M}{r^2}$$

### N.0.11 Problem 11

We start from  $\vec{\nabla}^2 U = 4\pi\mathcal{G}\rho$ . We then have

$$[\vec{\nabla}^2][U] = [\mathcal{G}][\rho]$$

so

$$[U] = [\mathcal{G}][\rho]/[\vec{\nabla}^2] = M^{-1}L^3T^{-2} \cdot ML^{-3} \cdot L^2 = L^2T^{-2} = \underbrace{ML^2T^{-2}}_{\text{energy}}/M$$

(see lecture notes on physical dimensions)

### N.0.12 Problem 12

Escape velocity is speed at which kinetic energy is equal to gravitational potential energy, i.e.

$$\frac{1}{2}mv^2 = mgH$$

so  $v = \sqrt{2gH}$  and since  $g = \mathcal{G}M/R^2$  then the escape velocity at the surface (i.e.  $H = R$ ) is given by

$$v = \sqrt{2\mathcal{G}M/R}$$

$$v_{earth} \simeq 11.2 \text{ km/s}$$

$$v_{moon} \simeq 2.4 \text{ km/s}$$

See [https://www.newworldencyclopedia.org/entry/Escape\\_velocity](https://www.newworldencyclopedia.org/entry/Escape_velocity)

### N.0.13 Problem 13

$$\begin{aligned} E &= - \int_V \rho U dV \\ &= - \int_V \rho_0 U(r) r^2 \sin \theta dr d\theta d\phi \\ &= -4\pi \rho_0 \int_0^R r^2 U(r) dr \\ &= -4\pi \rho_0 \int_0^R r^2 \left[ \frac{2\pi}{3} \mathcal{G} \rho_0 r^2 - \frac{3}{2} \frac{\mathcal{G}M}{R} \right] dr \\ &= -\frac{8\pi^2 \rho_0^2}{3} \mathcal{G} \int_0^R r^4 dr + 6\rho_0 \pi \frac{\mathcal{G}M}{R} \int_0^R r^2 dr \\ &= \dots \\ &= \frac{8\pi}{5} \rho_0 M R^2 \mathcal{G} \end{aligned} \tag{N.12}$$

### N.0.14 Problem 14

We start from the Laplace operator in spherical coordinates:

$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left( \sin \theta \frac{\partial}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2}{\partial \phi^2}$$

Because of the symmetry of the problem, the solution is expected to only depend on  $r$ , and not on  $\theta$  nor  $\phi$ , so that  $\partial_\theta \rightarrow 0$  and  $\partial_\phi \rightarrow 0$  in the equation above. We end up with:

$$\Delta = \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial}{\partial r} \right)$$

Inside the planet, the density is not zero, so we need to solve a Poisson equation

$$\Delta U = \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi \mathcal{G} \rho_0$$

Outside the planet the density is zero and we need to solve a Laplace equation:

$$\Delta U = \frac{1}{r^2} \frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 0$$

We start with the Poisson equation which we rewrite as follows:

$$\frac{\partial}{\partial r} \left( r^2 \frac{\partial U}{\partial r} \right) = 4\pi \mathcal{G} \rho_0 r^2$$

We integrate once and obtain

$$r^2 \frac{\partial U}{\partial r} = \frac{4\pi}{3} \mathcal{G} \rho_0 r^3 + A$$

where  $A$  is a constant to be specified later. We divide by  $r^2$  and obtain

$$\frac{\partial U}{\partial r} = \frac{4\pi}{3} \mathcal{G} \rho_0 r + \frac{A}{r^2}$$

The radial component of the gradient operator is simply  $\partial_r$  so that the equation above is (save a minus sign)  $g_r$ :

$$g_r(r) = -\frac{4\pi}{3} \mathcal{G} \rho_0 r - \frac{A}{r^2}$$

When  $r \rightarrow 0$  the gravity acceleration must remain finite so we need to set  $A = 0$ . Then

$$\boxed{g_r(r)|_{inside} = -\frac{4\pi}{3} \mathcal{G} \rho_0 r}$$

We integrate once more and obtain

$$U(r)|_{inside} = \frac{2\pi}{3} \mathcal{G} \rho_0 r^2 + B$$

where  $B$  is a constant.

We now turn to the Laplace equation outside the planet which yields

$$r^2 \frac{\partial U}{\partial r} = C$$

where  $C$  is a constant. Then it follows that

$$\frac{\partial U}{\partial r} = \frac{C}{r^2}$$

or,

$$U(r) = -\frac{C}{r} + D$$

When  $r \rightarrow \infty$  the potential tends to zero, so that  $D = 0$ . Then

$$\boxed{U(r)_{outside} = -\frac{C}{r}}$$

and from  $g_r(r) = -\partial_r U$  we get

$$g_r(r)|_{outside} = -\frac{C}{r^2}$$

We have solved the Poisson and Laplace equations but remain two constants to be specified. In order to do so we invoke the continuity of the gravity acceleration and potential at the surface of the planet:

$$\begin{aligned} g_r(r = R)|_{inside} &= g_r(r = R)|_{outside} \\ U(r = R)|_{inside} &= U(r = R)|_{outside} \end{aligned}$$

The first continuity condition yields

$$-\frac{C}{R^2} = -\frac{4\pi}{3}\mathcal{G}\rho_0 R$$

i.e.,  $C = M\mathcal{G}$ . The second continuity condition then yields

$$-\frac{C}{R} = -\frac{M\mathcal{G}}{R} = \frac{2\pi}{3}\mathcal{G}\rho_0 R^2 + B$$

i.e.  $B = -\frac{3}{2}\frac{M\mathcal{G}}{R}$ .

If all the mass is concentrated at the origin then by definition

$$g_r(r) = \frac{M\mathcal{G}}{r^2} \quad U(r) = -\frac{M\mathcal{G}}{R}$$

Finally

$$P(r) = -\int_r^R \rho(r')g(r')dr' = \int_r^R \rho_0 \frac{4\pi}{3}\mathcal{G}\rho_0 r' dr' = \frac{2\pi}{3}\mathcal{G}\rho_0^2(R^2 - r^2)$$

### N.0.15 Problem 15

We start from (10.36), i.e

$$U(r) = -\int_r^\infty \frac{Gm(r')}{r'^2} dr'$$

Outside the sphere,  $r > R$  and the mass at any location  $r' > R$  is simply  $M$ . Then

$$U(r) = -\int_r^\infty \frac{GM}{r'^2} dr' = -\frac{GM}{r}$$

When  $r < R$  we can split the integral in two:

$$U(r) = -\int_r^R \frac{Gm(r')}{r'^2} dr' - \int_R^\infty \frac{Gm(r')}{r'^2} dr'$$

We have just computed the second term so we focus on the first one. In this integral the mass  $m(r') = \frac{4}{3}\pi r'^3 \rho_0$  so that

$$U(r) = -\int_r^R \frac{G}{r'^2} \frac{4}{3}\pi r'^3 \rho_0 dr' - \frac{GM}{r} = -\frac{1}{2} \frac{4}{3}\pi \mathcal{G}(R^2 - r^2) - \frac{GM}{r} = \frac{2\pi}{3}\rho_0 \mathcal{G}r^2 - \frac{3}{2} \frac{GM}{r}$$



# Appendix O

## A quick guide to Paraview and gnuplot

### O.0.1 Paraview

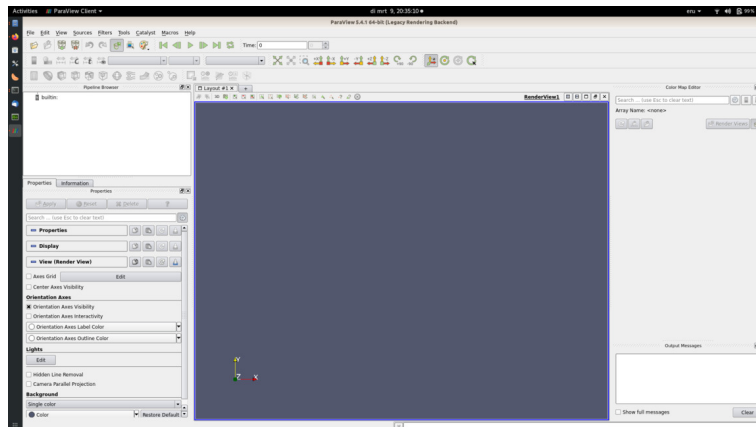
#### Installation procedure

If you have Ubuntu, type the following in a terminal and follow the instructions.

```
sudo apt-get install paraview
```

Upon completion, type 'paraview' followed by Enter in the terminal and your screen should look similar to Screen Capture 1.

If you run Windows or MacOS<sup>1</sup>, go to [www.paraview.org](http://www.paraview.org). Click on 'download'. The website automatically detects your OS<sup>2</sup>. Download the latest version(.exe for Windows, .dmg for Apple), and install it on your computer. Find the icon on your computer, double click on it and your screen should look similar to Screen Capture 1.



Screen Capture 1.

#### Opening a file

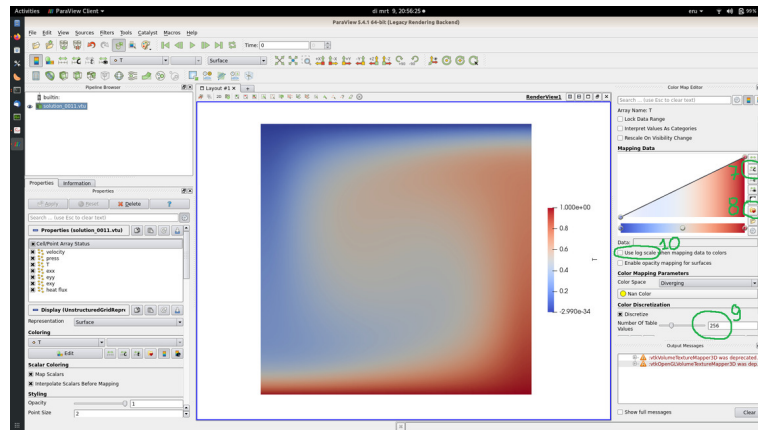
Press Ctrl+O on your keyboard or click File>Open> and the following window should appear after you press the green button Apply:

---

<sup>1</sup>You could also use Home Brew <https://formulae.brew.sh/cask/paraview>

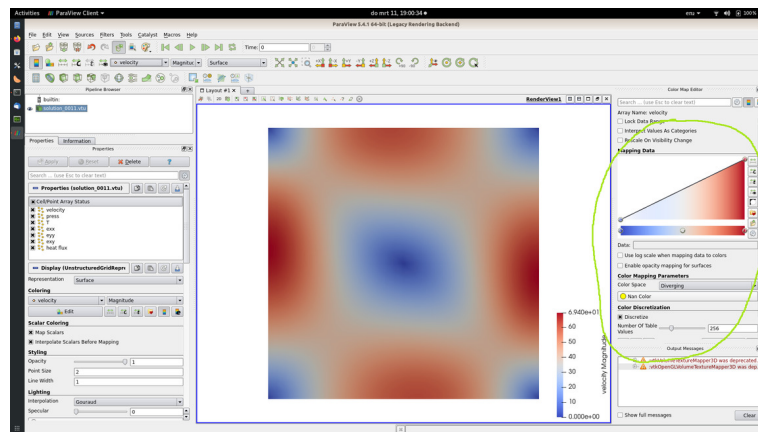
<sup>2</sup>[https://en.wikipedia.org/wiki/Operating\\_system](https://en.wikipedia.org/wiki/Operating_system)





- 7: change the range of the variable;
- 8: change the colour scale;
- 9: change the number of colours inside the scale;
- 10: switch on logarithmic scale.

When it comes to choosing colours, please see: Crameri, Shephard, and Heron [286] and Zelst, Crameri, Pusok, Glerum, Dannberg, and Thieulot [1403].

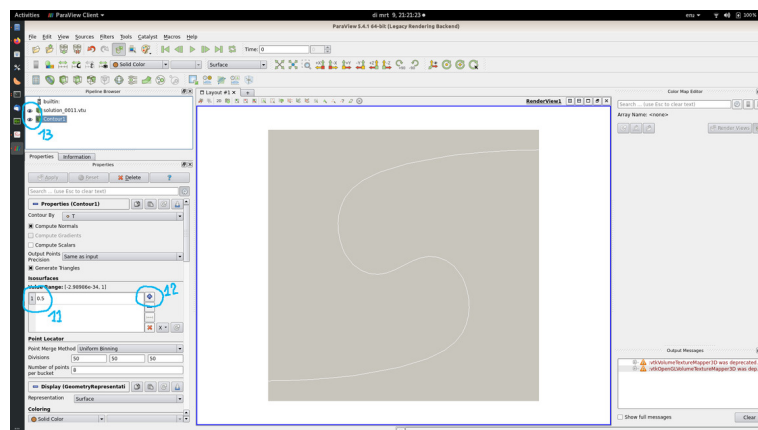


If you find that the circled area is missing on your screen, go to View and click on Color Map Editor.

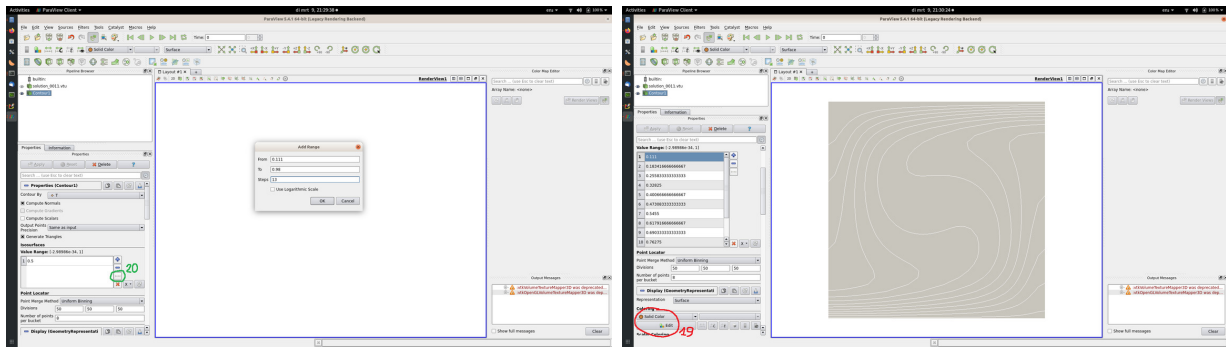
add how to add color scales

## Isocontours

Having clicked on the icon numbered 3 in the panels above, your screen should look like this:

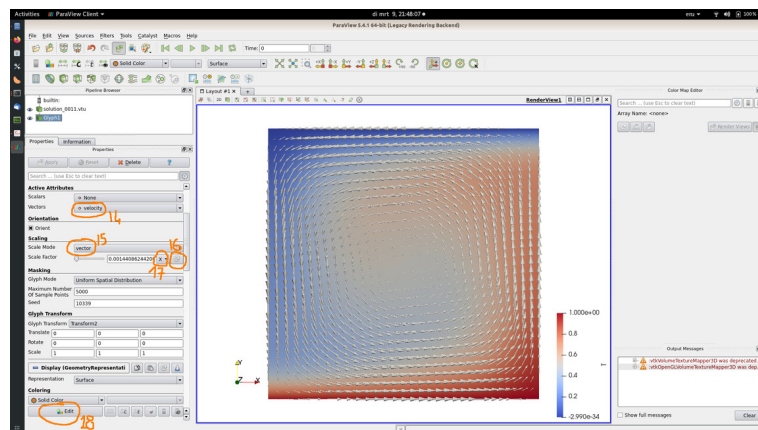


- 11: value of the isocontour;
- 12: add or remove an isocontour;
- 13: toggle the background grey square and the isocontours on/off by clicking on their respective eye.



Left: If you want automatically generated isocontours, remove the existing one and click on 20. A small window opens: fill the min/max/number values and click OK. Right: Having obtained these isocontours you can change the colour of the lines by clicking on 19.

## Vector field arrows



In order to obtain such arrows, make sure that you go through points 14 and 15.

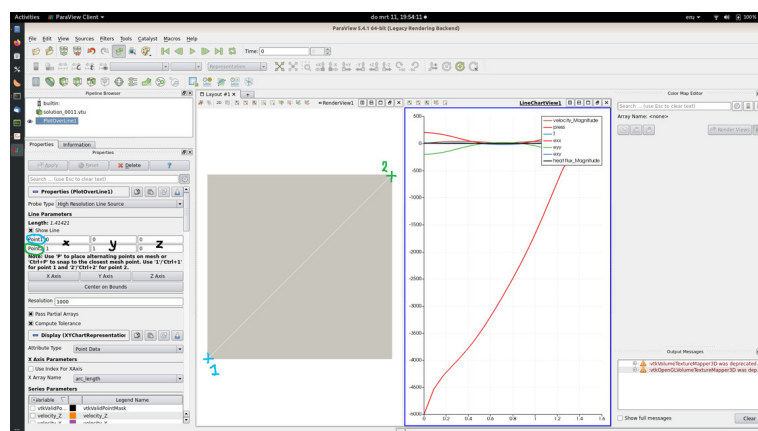
Then click on the icon 16. In order to change the scale of the arrows change the value in 17.

## Exporting to png

File>Save Screenshot. Click OK on the first panel. Enter the name of the file you have chosen and click OK.

## Exporting line data

Filters > Data Analysis > Plot Over Line.

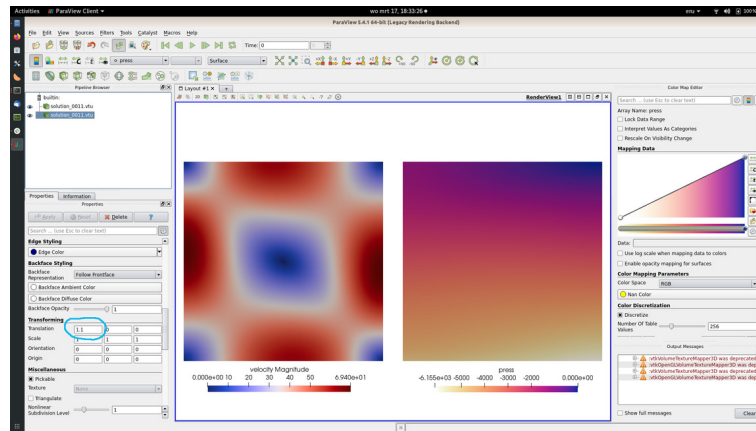


You can change the coordinates of the beginning and the end of the line.

## Getting rid of 'Apply'

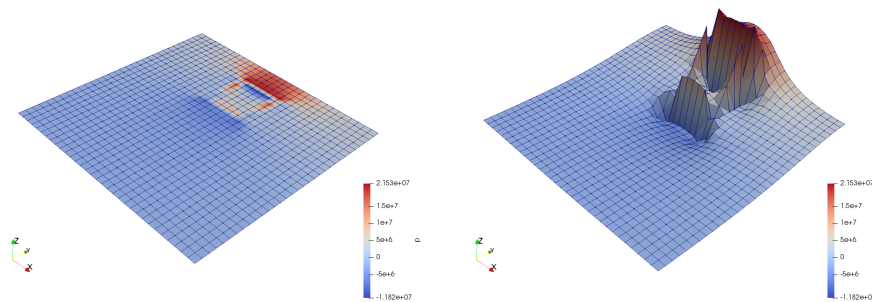
It can be annoying to have to press Apply all the time so if you wish to bypass it, got to Edit > Settings, and tick the 'Auto Apply' box in the window that appears.

## Multiple vtu files at once



You can load multiple vtu files or the same one multiple times and move each where you want it.

## Warp by scalar



Filters  $\downarrow$  Alphabetical  $\downarrow$  Warp By Scalar

## O.0.2 gnuplot

gnuplot is a famous and widely used command-line program that can generate two- and three-dimensional plots of functions, data, and data fits. It dates back to 1986 and runs on all operating systems (Linux, Unix, Microsoft Windows, macOS).

<http://www.gnuplot.info/>

<http://www.gnuplotting.org/>

<http://lowrank.net/gnuplot/index-e.html>

The gray boxes indicate that its content takes place in the terminal.

## Installing gnuplot

If you are using Ubuntu, you can install gnuplot as follows:

```
> sudo apt-get install gnuplot
```

## Interactive use

In the following pages I explain how to use the gnuplot program from the terminal. Having done so, in the terminal simply type

```
> gnuplot
```

The following should then appear:

```
G N U P L O T
Version 5.2 patchlevel 2 last modified 2017-11-01

Copyright (C) 1986-1993, 1998, 2004, 2007-2017
Thomas Williams, Colin Kelley and many others

gnuplot home: http://www.gnuplot.info
faq, bugs, etc: type "help FAQ"
immediate help: type "help" (plot window: hit 'h')

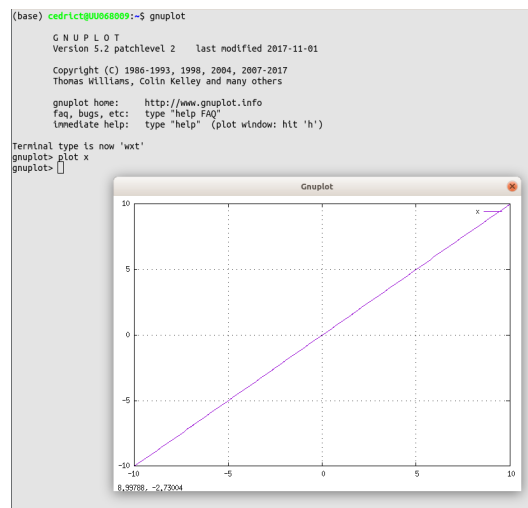
Terminal type is now 'wxt'

gnuplot>
```

The prompt means that gnuplot is expecting instructions. We start by making sure that the terminal type is such that a window appears in this interactive mode. We test this by plotting a simple function,  $f(x) = x$ :

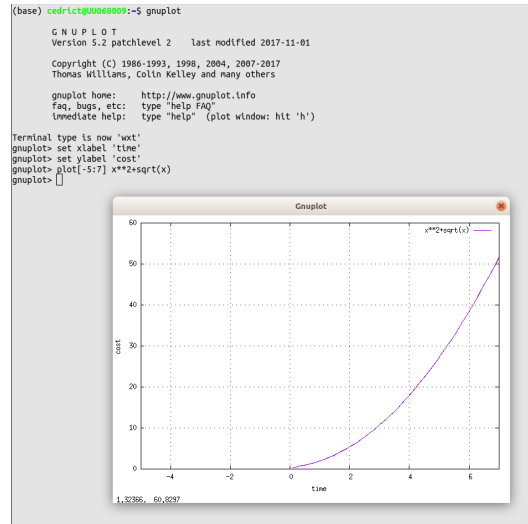
```
gnuplot> plot x
```

You should then obtain something similar:



You can specify the  $x$  range, change the function to  $x^2 + \sqrt{x}$  and label the axes as follows:

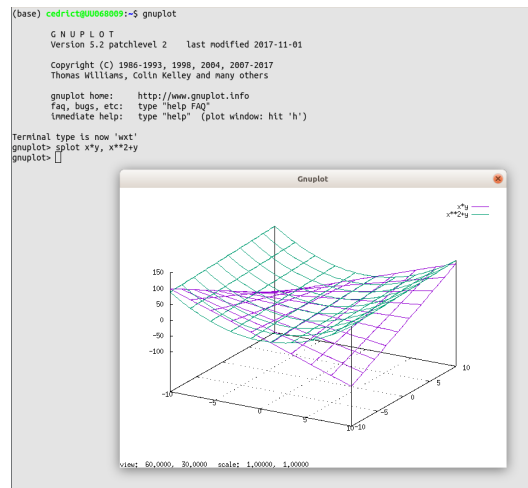
```
gnuplot> set xlabel 'time'
gnuplot> set ylabel 'cost'
gnuplot> plot[-5:7] x**2+sqrt(x)
```



We can also plot functions of both  $x$  and  $y$  as follows:

```
gnuplot> splot x*y, x**2+y
```

and we get



Another nice feature in the interactive is the fact that you can use the left button of the mouse to rotate the plot!

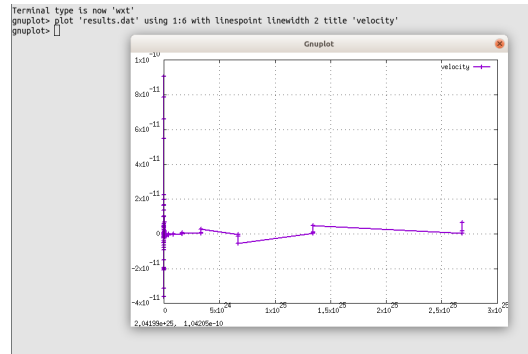
Finally, let us assume that there is a file *results.dat* in the folder and that it contains results from experimental measurements or numerical values organised in columns as follows:

```
1e17 8 0 256000 384000 4.91094e-12 -0.00533647 -714769 0
1e17 32 0 256000 384000 1.96438e-11 -0.0213459 -2.85908e+06 0
1e17 128 0 256000 384000 7.8575e-11 -0.0853835 -1.14363e+07 0
2e17 8 0 256000 384000 3.43871e-12 -0.00533555 -714753 0
2e17 32 0 256000 384000 1.37548e-11 -0.0213422 -2.85901e+06 0
2e17 128 0 256000 384000 5.50193e-11 -0.0853688 -1.1436e+07 0
4e17 8 0 256000 384000 4.13458e-12 -0.00533372 -714720 0
...
67108864e17 128 0 256000 384000 -5.28841e-12 -0.0212269 -1.27701e+07 0
134217728e17 8 0 256000 384000 2.93622e-13 -0.00132619 -798163 0
134217728e17 32 0 256000 384000 1.17449e-12 -0.00530475 -3.19265e+06 0
134217728e17 128 0 256000 384000 4.69795e-12 -0.021219 -1.27706e+07 0
```

```
268435456e17 8 0 256000 384000 4.03077e-13 -0.00132594 -798181 0
268435456e17 32 0 256000 384000 1.61231e-12 -0.00530376 -3.19272e+06 0
268435456e17 128 0 256000 384000 6.44923e-12 -0.0212151 -1.27709e+07 0
```

In this case we may wish to plot the 6th column as a function of the 1st one:

```
plot 'results.dat' using 1:6 with linespoint linewidth 2 title 'velocity'
```



Since typing these instructions time and time again is a bit tedious gnuplot allows the user to use short versions of these commands:

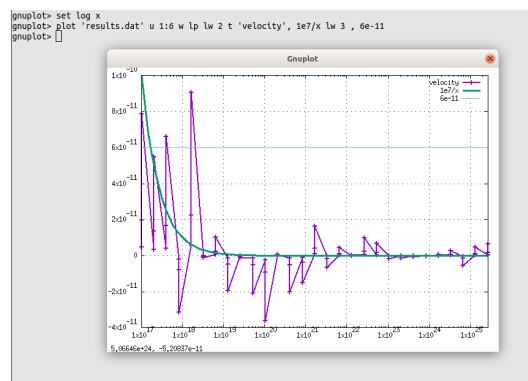
```
gnuplot> plot 'results.dat' u 1:6 w lp lw 2 t 'velocity'
```

We see that the range of  $x$  values spans many order of magnitudes so we wish to use a logarithmic scale on the  $x$ -axis.

```
gnuplot> set log x
```

Also, I can combine data with analytical function:

```
gnuplot> set log x
gnuplot> plot 'results.dat' u 1:6 w lp lw 2 t 'velocity', 1e7/x lw 3 , 6e-11
```



Finally, we may wish to export the plot to a file, say a pdf file. We must then re-assign the terminal, give a name to the file and re-plot:

```
gnuplot> set term pdf
gnuplot> set output 'results.pdf'
gnuplot> plot 'results.dat' u 1:6 w lp lw 2 t 'velocity', 1e7/x lw 3 , 6e-11
```

You can exit the session by typing



```
gnuplot> exit
```

You should find *results.pdf* in your folder next to *results.dat*.

## Scripting gnuplot

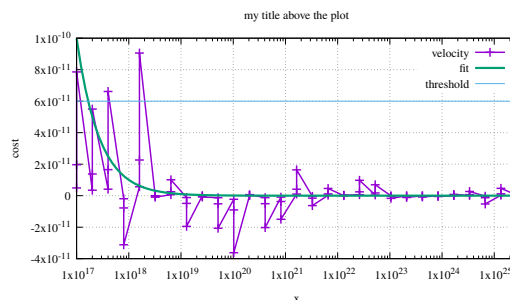
Although the interactive approach is very useful its workflow is not practical if one wishes (for instance) to produce the same plot with different data, or to communicate a figure to another scientist.

We will therefore now turn to scripting. The idea is simple: write all the gnuplot commands in a text file, say *script.gnuplot* and pass this script as argument to gnuplot:

```
> gnuplot script.gnuplot
```

This file contains the following:

```
set term pdf font "Times,12pt"
set output 'results.pdf'
set grid
set xlabel 'x'
set ylabel 'cost'
set log x
set title 'my title above the plot'
plot 'results.dat' u 1:6 w lp lw 2 t 'velocity', 1e7/x lw 3 t 'fit' , 6e-11 t 'threshold'
```



Note that I have added a title to the plot as well.

## Greek letters

In order to display Greek letters the `/Symbol` command:

```
set xlabel '{/Symbol d}/{/Symbol r}'
```

| Alphabet | Symbol  | Alphabet | Symbol  | Alphabet | Symbol                 | Alphabet | Symbol             |
|----------|---------|----------|---------|----------|------------------------|----------|--------------------|
| A        | Alpha   | N        | Nu      | a        | alpha ( $\alpha$ )     | n        | nu $\nu$           |
| B        | Beta    | O        | Omicron | b        | beta ( $\beta$ )       | o        | omicron            |
| C        | Chi     | P        | Pi      | c        | chi ( $\chi$ )         | p        | pi $\pi$           |
| D        | Delta   | Q        | Theta   | d        | delta ( $\delta$ )     | q        | theta $\theta$     |
| E        | Epsilon | R        | Rho     | e        | epsilon ( $\epsilon$ ) | r        | rho $\rho$         |
| F        | Phi     | S        | Sigma   | f        | phi ( $\phi$ )         | s        | sigma $\sigma$     |
| G        | Gamma   | T        | Tau     | g        | gamma ( $\gamma$ )     | t        | tau $\tau$         |
| H        | Eta     | U        | Upsilon | h        | eta ( $\eta$ )         | u        | upsilon $\upsilon$ |
| I        | iota    | W        | Omega   | i        | iota ( $\iota$ )       | w        | omega $\omega$     |
| K        | Kappa   | X        | Xi      | k        | kappa ( $\kappa$ )     | x        | xi $\xi$           |
| L        | Lambda  | Y        | Psi     | l        | lambda ( $\lambda$ )   | y        | psi $\psi$         |
| M        | Mu      | Z        | Zeta    | m        | mu ( $\mu$ )           | z        | zeta $\zeta$       |

## piecewise function

You can define piecewise functions as follows:

$$f(x) = x < a \quad ? \quad 1 \quad : \quad 1/0$$
$$g(x) = x \geq a \quad ? \quad 1 \quad : \quad 1/0$$

and then use these functions like any function, e.g.:

```
plot[-10:12] f(x),g(x)
```

## linetype and dashtype

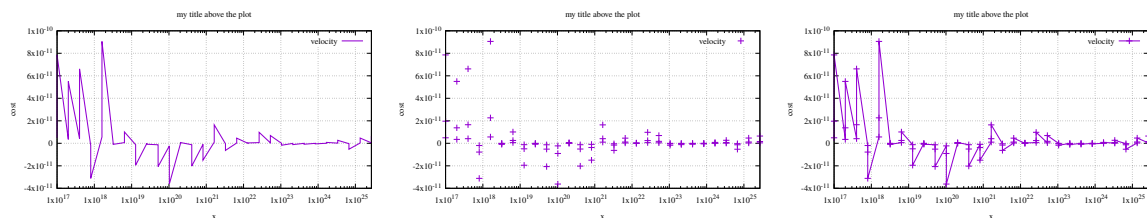
There are essentially three ways of plotting data:

```
plot 'results.dat' u 1:6 w l
```

```
plot 'results.dat' u 1:6 w p
```

```
plot 'results.dat' u 1:6 w lp
```

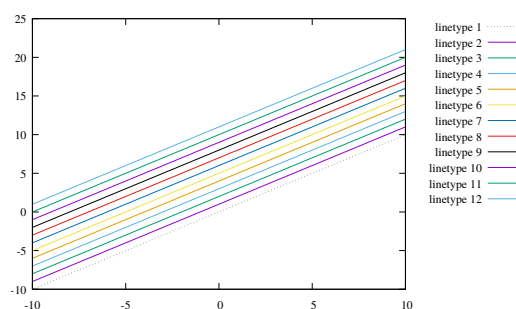
corresponding to (from left to right):



The following script

```
set output 'linetypes.pdf'
plot[] [] \
x+0 w l lt 0 t 'linetype 1', \
x+1 w l lt 1 t 'linetype 2', \
x+2 w l lt 2 t 'linetype 3', \
x+3 w l lt 3 t 'linetype 4', \
x+4 w l lt 4 t 'linetype 5', \
x+5 w l lt 5 t 'linetype 6', \
x+6 w l lt 6 t 'linetype 7', \
x+7 w l lt 7 t 'linetype 8', \
x+8 w l lt 8 t 'linetype 9', \
x+9 w l lt 9 t 'linetype 10', \
x+10 w l lt 10 t 'linetype 11', \
x+11 w l lt 11 t 'linetype 12'
```

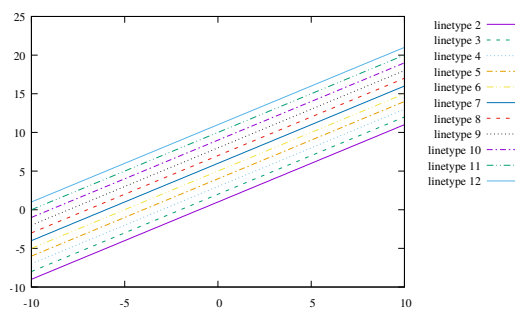
generates the following plot:



We see that the colours repeat from linetype 10. Fortunately we can also combine linetypes with dashtypes. The following script

```
set output 'dashtypes.pdf'
plot[] []\
x+1 w l lt 1 dt 1 t 'linetype 2',\
x+2 w l lt 2 dt 2 t 'linetype 3',\
x+3 w l lt 3 dt 3 t 'linetype 4',\
x+4 w l lt 4 dt 4 t 'linetype 5',\
x+5 w l lt 5 dt 5 t 'linetype 6',\
x+6 w l lt 6 dt 6 t 'linetype 7',\
x+7 w l lt 7 dt 7 t 'linetype 8',\
x+8 w l lt 8 dt 8 t 'linetype 9',\
x+9 w l lt 9 dt 9 t 'linetype 10',\
x+10 w l lt 10 dt 10 t 'linetype 11',\
x+11 w l lt 11 dt 11 t 'linetype 12'
```

generates the following plot

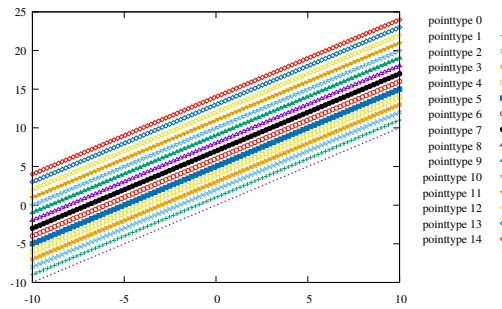


and we see that there are only 5 different dash types.

Finally, we turn to point types. The following script

```
set output 'pointtypes.pdf'
plot[] []\
x+0 w p pt 0 ps .5 t 'linetype 1',\
x+1 w p pt 1 ps .5 t 'linetype 2',\
x+2 w p pt 2 ps .5 t 'linetype 3',\
x+3 w p pt 3 ps .5 t 'linetype 4',\
x+4 w p pt 4 ps .5 t 'linetype 5',\
x+5 w p pt 5 ps .5 t 'linetype 6',\
x+6 w p pt 6 ps .5 t 'linetype 7',\
x+7 w p pt 7 ps .5 t 'linetype 8',\
x+8 w p pt 8 ps .5 t 'linetype 9',\
x+10 w p pt 10 ps .5 t 'pointtype 10',\
x+11 w p pt 11 ps .5 t 'pointtype 11',\
x+12 w p pt 12 ps .5 t 'pointtype 12',\
x+13 w p pt 12 ps .5 t 'pointtype 13',\
x+14 w p pt 12 ps .5 t 'pointtype 14'
```

generates the following plot



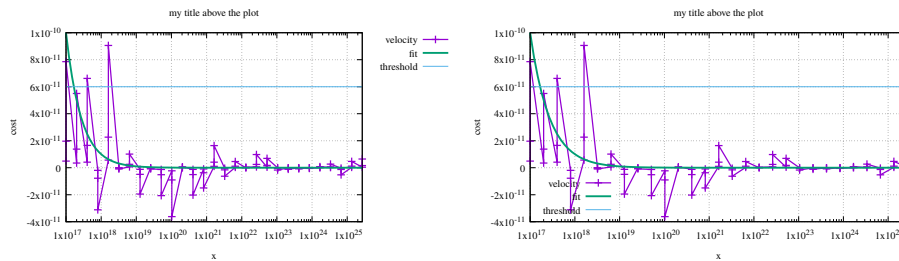
Note that I have used the `ps` command ('point size') to make the points 50% smaller than normal.

## Moving the 'key'

The default is inside top right, but it can be changed, e.g.:

```
set key outside
set key bottom left
```

corresponding to (from left to right):

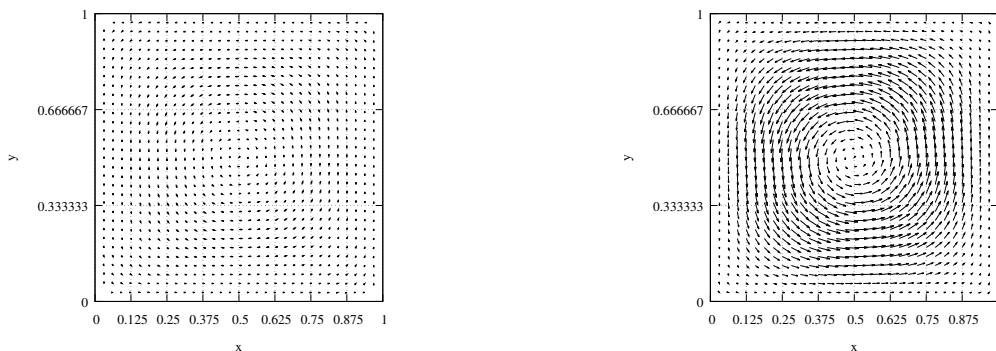


## Plotting arrows

Let us now turn to the *velocity.dat* file which consists of four columns:  $x$ ,  $y$ ,  $v_x$  and  $v_y$ .

```
set output 'velocity_1.pdf'
set xlabel 'x'
set ylabel 'y'
set xtics 0.125
set ytics 0.333333333333
set grid
set size square
plot[0:1][0:1]\
'velocity.dat' u 1:2:3:4 w vectors lt -1 notitle
```

Note that I have required the plot to be square, that the ticks on the  $x$ -axis should be spaced 0.125 while the ticks on the  $y$ -axis should be spaced 0.333. We obtain the left plot a):



The arrows are too small, so we scale each vector component by a factor 4. All we need to do is as follows:

```
plot[0:1][0:1]\
'veLOCITY.dat' u 1:2:($3*4):($4*4) w vectors lt -1 notitle
```

Note the dollar sign which means that gnuplot takes the value in column 3 or 4 and multiplies it by 4. The resulting figure is shown in b).

### **Powers of 10**

```
set format y "10^{%L}"
```

### **Least square fit**

Assuming we have a file containing data, e.g. `data.ascii`, that we want to fit by means of a linear relationship over the range  $[-1, +1]$ :

```
f(x)=a*x+b
fit [-1:1] f(x) 'data.ascii' u 1:2 via a,b
```

This should return a few lines in the console indicating whether convergence was reached and then also the  $a$  and  $b$  values. In order to plot the line, simply do:

```
plot[] 'data.ascii' u 1:2, f(x)
```

### **coloring areas**

```
set style rect fc lt -1 fs solid 0.1 noborder
set obj rect from 0, graph 0 to 15, graph 1
```

### **vertical line**

To draw a vertical line from the bottom to the top of the graph at  $x=3$ , use:

```
set arrow from 3, graph 0 to 3, graph 1 nohead
```

### **Show list of all available colors**

In an interactive gnuplot session type:

```
> show colors
```

# Appendix P

## A few L<sup>A</sup>T<sub>E</sub>X features



### newcommand

This features allows to define new commands which can then be used throughout the document. In *manual.tex* you will find

```
\newcommand{\K}{{\mathbb{K}}}
\newcommand{\Ranb}{{\mathsf{Ra}}}
\newcommand{\nineteeneightysix}{{\color{violet}\bf 1986}}
```

which correspond to  $\mathbb{K}$ ,  $\mathsf{Ra}$  and **1986**.

### fullpage package

How to extend margins for the whole document (such as this one):

```
\usepackage[cm]{fullpage}
```

### siunitx package

```
\usepackage{siunitx}
\DeclareSIUnit\year{yr}
```

### tikz

```
\usetikzlibrary{arrows, arrows.meta}
```

## include verbatim material inside a line

```
\verb"git pull upstream master" and then eat an ice cream.
```

results in git pull upstream master and then eat an ice cream.

## bibliography

```
opla
 \footfullcite
```

## how to refer to a latex document from another document?

This is the first latex file names `manual.tex`. It contains sections and equations, all labelled.

```
\documentclass[a4paper]{book}
\begin{document}
\chapter{opla1} \label{ch1}
\section{meuh} \label{ss1}
\subsection{popo}
opla \ref{ch1} opla \ref{ss1}
\begin{equation}
\alpha= \beta \label{eq:one}
\end{equation}
\end{document}
```

This second file will load the `manual.aux` with the 'xr' package and then all labels of that file are now available in this file:

```
\documentclass[a4paper]{article}
\usepackage{amsmath}

\usepackage{xr}
\externaldocument[MMM-]{manual}

\begin{document}

opla

the introduction to volume1 (\ref{MMM-ss1})
see \eqref{MMM-eq:one}

\end{document}
```

# Appendix Q

## Linux how to

When encountering the terminal for the first time one soon realises that some basic commands are needed to navigate the folders, edit, copy or delete files, etc ...

Please also have a look at this website: <https://ryanstutorials.net/linuxtutorial/>

For Windows users, you must abandon your preconceived (and totally arbitrary) ideas about the 'C:' drive. Before hard drives even existed, the computer would have one or two floppy drives and would reserve the A: and B: drive letters for them. Nowadays these devices are usually not installed in the computer anymore, but the labelling remains<sup>1</sup>.

The root of the file system is simply /. Your 'home' is most likely in /home/your-family-name/. Inside this you will find Documents, Downloads, etc ...

The basic set of commands is in fact not very large (Please also watch <https://youtu.be/6bMYzzrycV0>):

- **man**: Documentation in Linux is mostly available in the form of *man pages*. They are usually written in the style of a reference manual and can be daunting at first. They are usually: the name, a compact formulation of the syntax (that can be scary for more complex programs), a description about what the software actually is and does, examples (if you are lucky) and explanations of all the options mentioned above. Typical use:

```
> man name-of-command-I-want-to-learn-about
```

- **cd**: it stands for 'change directory' (i.e. change folder). Typical use:

```
> cd results
```

where **results** should be an existing folder. You can check this with:

- **ls**: it stands for 'list'. This commands lists all files and folders Typical use:

```
> ls
```

or

```
> ls -la
```

if you wish to have one item per line, or

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Drive\\_letter\\_assignment](https://en.wikipedia.org/wiki/Drive_letter_assignment)



```
> ls -l
```

if you also want to see all files/folders starting with '.' My preference goes to

```
> ls -lhG
```

Use the 'man' command to know what this does!

- **mv**: it stands for 'move' but it has in fact a hidden functionality: rename. If you type

```
> mv garfield.txt odie.txt
```

then the file *garfield.txt* has been renamed *odie.txt* <sup>2</sup>.

You can move a file in a different folder as follows:

```
> mv garfield.txt ../myfolder/
```

This moves the file to a folder that is one level up. This will work only if the folder exists. If you need to make a new folder, then use:

- **mkdir**: it stands for 'make directory'. Typical usage

```
> mkdir res_123
```

creates a folder named *res\_123*.

- **rm**: it stands for 'remove'. Before we go any further: this command is dangerous. Unlike its counterpart based on selecting a file with a mouse and deleting it, **rm** does not send the file to the Trash folder (or Windows Recycle Bin). It simply deletes it forever. No turning back! Typical usage

```
> rm myoldfile.txt
```

deletes the file. Note that by default, it does not remove directories. In order to remove a folder, one needs to type

```
> rm -r myoldfolder
```

Unless you are an experienced user, never use **rm** in conjunction with **\*** and/or in a recursive way. If things go wrong, you will delete entire portions of your hard drive at best, or will destroy your operating system at worse.

- **pwd**: it stands for 'print working directory'. If you are unsure of where the current prompt of the terminal is, simply

```
> pwd
```

and it will return the full path, from the root to where the prompt is.

- **du**: it stands for 'disk usage'. In this case always tag the **-h** option to it:

---

<sup>2</sup><https://en.wikipedia.org/wiki/Odie> Obviously one cannot transform Garfield into Odie.

```
> du -h
```

It will list the size of all folders in the folder the prompt is in. If you wish to know the size of an object simply do

```
> du -h file-or-folder
```

The `ls -lh` command would have told you as much, but for all files inside the folder.

- **grep**: searches for patterns in each file. Typical usage

```
> grep linear *.tex
```

This searches all occurrences of the word 'linear' in all .tex files. If you wish to look for this word in all files inside subfolders:

```
> grep -r linear .
```

- **more**: allows to visualise the content of an [ascii](https://en.wikipedia.org/wiki/ASCII)<sup>3</sup> file inside the terminal without using a text editor. In other words, you can look into the file but not change its content.

```
> more interesting-file
```

- **top/htop**: allows to visualise which process is running on the computer and how much CPU and memory it takes.
- **wget**: it is a computer program that retrieves content from web servers. It is part of the GNU Project. Its name derives from "World Wide Web" and "get." It supports downloading via HTTP, HTTPS, and FTP.

```
> wget address.of.file.online.on.server
```

See example of use in Section [11.12.5](#).

- **ssh**: This command is necessary to connect to a remote computer from the terminal. By default most Linux computers run a client and server ssh program so that one can connect to any other computer if its address is known (as well as the username and password). Here is how I connect to my shrek desktop computer:

```
> ssh shrek.geo.uu.nl -Y -l thieulot
```

The -Y option ensures that I have X11 support. If successful, the prompt of the terminal points to the default folder on the remote machine and I can control the remote computer via the command line. Type exit to break the connection.

This is used to log in on remote servers like clusters where large calculations take place.

- **sftp**: allows to get or put files on a remote computer via a secured connection. By default, SFTP uses the SSH protocol to authenticate and establish a secure connection. Because of this, the same authentication methods are available that are present in SSH.

---

<sup>3</sup><https://en.wikipedia.org/wiki/ASCII>

```
> sftp garfield@remote.univ.nl
```

This will prompt for the password of the remote machine associated to user garfield. If successful the prompt in the terminal then points to the default folder on the remote computer. Once connected simply type exit to break the connection. All shell commands are available: pwd, ls, cd, etc ... you can list the contents of the current directory on the local machine using `lls`. If we want to download files from our remote host, you can do so using the get command:

```
> get remotefile
```

If you wish to download an entire folder, simply use

```
> get -r remotefolder
```

Conversely, you can transfer files from your local machine onto the remote one:

```
> put remotefile
```

- `scp`: to do

In what follows I list a few 'tricks' which I find useful or just can never remember:

- How to convert files in batch

```
> convert '*.png' converted_%04d.jpg
```

- How to Remove unwanted empty lines in a file with vi(m)  
Use either of the following commands to delete all empty lines:

```
:g/~/d
:v/./d
```

If you want to delete all lines that are empty or that contain only whitespace characters (spaces, tabs), use either of:

```
:g/^\s*/d
:v/^\S/d
```

- How to find LAPACK

```
BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.7.1
LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.7.1
```

- How to apt-get MUMPS

```
> sudo apt-get install libmumps-seq-dev
```

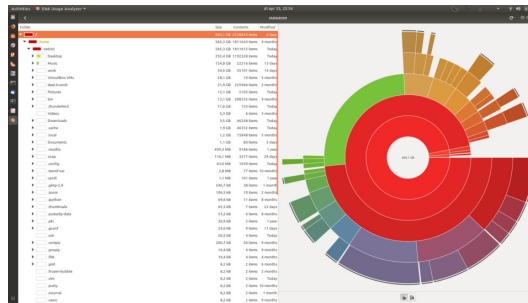
- How to remove a big file wrongly committed

```
> git filter-branch --tree-filter 'rm -rf path/to/your/file' HEAD
> git push
```

- How to check your disk usage (ubuntu)

Disk Usage Analyzer (aka Baobab) is a graphical, menu-driven viewer that you can use to view and monitor your disk usage and folder structure. It is part of every GNOME desktop.

```
> baobab
```



- How to convert images into another format in the command line

```
> convert file.png file.jpg
```

You can also resize on the fly:

```
> convert -resize 70% file.png file.jpg
```

- How to compress files: tar is used to pack files into an archive without zipping them. This enables easy attaching to email, or to sending across the internet without the computer having to start a new transfer for every file. Also it ensures everything belong to a certain package stays together. Optionally it can be used to zip files. Let us have an example of three files; *a.txt*, *b.py* and *c.dat* These can be packed together with

```
> tar -cf allFiles.tar *
```

*c* stands for *compress* and *f* for *file*, as can be seen in the manpage of tar. It can similarly be extracted using the command:

```
> tar -xf allFiles.tar *
```

to recover three original files.

Additionally, it is possible to zip them as well, by adding a *z* to the options. The default extension then becomes *tgz*:

```
> tar -czf allFiles.tgz *
> tar -xzf allFiles.tgz *
```

# Appendix R

## on using Fortran

### R.0.1 Full matrix multiplications in fortran

In fortran there is the intrinsic function *matmul*. However, it turns out that it is not always the fastest option to carry out (full) matrix multiplications.

This code is designed to test this:

```
program test
implicit none
! The order of the square matrices is 2000.
integer(kind=4)::n=1000
! Calculate the matrix multiplications:
! i) c:=a*b in a triple do-loop.
! ii) d:=a*b by matmul(a,b).
! iii) e:=a*b by dgemm in INTEL MKL.
real(kind=8),allocatable::a(:,:),b(:,:),c(:,:),d(:,:),e(:,:)
real(kind=8)::alpha,beta
integer(kind=4)::i,j,k,lda,ldb,lde
real(kind=8)::start,finish

allocate(a(n,n),b(n,n),c(n,n),d(n,n),e(n,n))
alpha=1.0;beta=1.0
lda=n;ldb=n;lde=n

! Generate the matrices, a and b, randomly.
call cpu_time(start)
call random_seed()
do j=1, n
do i=1, n
call random_number(a(i,j))
call random_number(b(i,j))
enddo
enddo
call cpu_time(finish)
write(unit=6,fmt=100) "The_generation_of_two_matrices_takes_",finish-start,"_seconds.
"

! i) c:=a*b in a triple do-loop.
call cpu_time(start)
c=0.0D0
do j=1, n
do i=1, n
do k=1, n
c(i,j)=c(i,j)+a(i,k)*b(k,j)
enddo
enddo
```

```

enddo
call cpu_time(finish)
write(unit=6,fmt=100) "A triple do-loop takes ", finish-start, " seconds."

! ii) d:=a*b by matmul(a,b).
call cpu_time(start)
d=0.0D0
d=matmul(a,b)
call cpu_time(finish)
write(unit=6,fmt=100) "A matmul(a,b) function takes ", finish-start, " seconds."

! iii) e:=a*b by dgemm in INTEL MKL.
call cpu_time(start)
e=0.0D0
call dgemm("N","N",n,n,n,alpha,a,lda,b,ldb,beta,e,lde)
call cpu_time(finish)
write(unit=6,fmt=100) "A DGEMM subroutine takes ", finish-start, " seconds."

deallocate(a,b,c,d,e)

stop
100 format(A,F8.3,A)
end program test

```

It is compiled as follows:

```
> gfortran -O3 prog.f90 -lblas
```

For  $100 \times 100$  matrices:

```

The generation of two matrices takes 0.004 seconds.
A triple do-loop takes 0.009 seconds.
A matmul(a,b) function takes 0.001 seconds.
A DGEMM subroutine takes 0.000 seconds.

```

For  $1000 \times 1000$  matrices:

```

The generation of two matrices takes 0.123 seconds.
A triple do-loop takes 1.527 seconds.
A matmul(a,b) function takes 0.080 seconds.
A DGEMM subroutine takes 0.054 seconds.

```

For  $1000 \times 2000$  matrices:

```

The generation of two matrices takes 0.392 seconds.
A triple do-loop takes 33.785 seconds.
A matmul(a,b) function takes 0.725 seconds.
A DGEMM subroutine takes 0.455 seconds.

```

## R.0.2 A simple example of an Interface

```

program kwadraat

implicit none

integer, parameter :: IntegerRoot = 6
real, parameter :: RealRoot = 4.5

```

```

Interface Square
 function RealSquare(root)
 real :: root
 real :: RealSquare
 end function

 function IntegerSquare(root)
 integer :: root
 integer :: IntegerSquare
 end function
end interface

write(*,*) "Integer_square: ", Square(IntegerRoot)
write(*,*) "Real_square: ", Square(realRoot)

end program

function IntegerSquare(root)
 implicit none
 integer :: root, IntegerSquare
 IntegerSquare = root**2
end function

function RealSquare(root)
 implicit none
 real :: root, RealSquare
 RealSquare = root*root
end function

```

# Appendix S

## mineral parameters

### S.0.1 Olivine

| ref.        | Wet<br><i>A</i>                                                   | Q (kJ/mol)   | <i>V</i> | <i>n</i>      | Dry<br><i>A</i>                                                     | Q (kJ/mol)   |
|-------------|-------------------------------------------------------------------|--------------|----------|---------------|---------------------------------------------------------------------|--------------|
| [505]       | $3.91 \cdot 10^3 MPa^{-n}/s$<br>$= 3.91 \cdot 10^{-15} Pa^{-n}/s$ | 430          | 0        | 3             | $2.42 \times 10^5 MPa^{-n}/s$<br>$= 2.42 \times 10^{-16} Pa^{-n}/s$ | 540          |
| [295]       | $3.91 \cdot 10^{-15} Pa^{-n}/s$                                   | 430          | 0        | 3             |                                                                     |              |
| [613]       |                                                                   |              |          |               | $2.4 \times 10^{-16} Pa^{-n}/s$                                     | 540          |
| [614]       |                                                                   |              | 0        |               | $2.42 \times 10^{-15} Pa^{-n}/s$                                    | 540          |
| [673]       | $3.906 \cdot 10^{-15} Pa^{-n}/s$                                  | 430          | 10-20    | 3             | $2.4169 \cdot 10^{-16} Pa^{-n}/s$                                   | 540          |
| [632]       |                                                                   |              | 0        |               | $1.43 \times 10^{-15} Pa^{-n}/s$                                    | 65           |
| [482]       | $4.89 \cdot 10^{-15} Pa^{-n}/s$                                   | 515          |          | 3.5           |                                                                     |              |
| [1023]      | $4.89 \cdot 10^{-15} Pa^{-n}/s$                                   | 515          |          | 3.5           | $4.85 \times 10^{-17} Pa^{-n}/s$                                    | 535          |
| [577]       | $4.89 \cdot 10^{-15} Pa^{-n}/s$                                   | 515          |          | 3.5           | $4.85 \times 10^{-17} Pa^{-n}/s$                                    | 535          |
| [672]       |                                                                   |              |          |               | $10^{6.1} MPa^{-n}/s$<br>$= 1.26 \cdot 10^{-12} Pa^{-n}/s$          | $510 \pm 30$ |
| [1041]      | $2 \times 10^3 MPa^{-n}/s$                                        | $471 \pm 31$ |          | $4 \pm 0.1$   | $2.5 \times 10^4 MPa^{-n}/s$                                        | $532 \pm 52$ |
|             | $= 2 \cdot 10^{-21} Pa^{-n}/s$                                    |              |          |               | $= 2.5 \cdot 10^{-17} Pa^{-n}/s$                                    |              |
| [1245]      | $5.33 \cdot 10^{-19}$                                             | 480          | 11       | 3.5           |                                                                     |              |
| [1245]      | $1.5 \cdot 10^{-18}$                                              | 335          | 4        | 1             |                                                                     |              |
| [685]       |                                                                   |              |          |               | $10^{5.04} MPa^{-n}/s$<br>$= 1.1 \cdot 10^{-16} Pa^{-n}/s$          | 530          |
| [782]       |                                                                   |              |          |               |                                                                     | 470          |
| [578]       | $3.58 \cdot 10^{-16} Pa^{-n}/s$                                   | $480 \pm 40$ | 11       | $3.5 \pm 0.3$ | $1.1 \cdot 10^{-16} Pa^{-n}/s$                                      | $530 \pm 4$  |
| [578]       | $8 \cdot 10^{-9} Pa^{-n}/s$                                       | $335 \pm 75$ | 4        | 1             | $1.2 \cdot 10^{-8} Pa^{-n}/s$                                       | $375 \pm 50$ |
| ELEFANT     |                                                                   |              |          |               |                                                                     |              |
| wetolivine1 | $3.9063 \cdot 10^{-15} Pa^{-n}/s$                                 | 430          | 15       | 3             |                                                                     |              |
| dryolivine1 |                                                                   |              |          |               | $2.4169 \cdot 10^{-16} Pa^{-n}/s$                                   | 540          |
| wetolivine2 | $4.89 \cdot 10^{-15} Pa^{-n}/s$                                   | 515          |          | 3.5           |                                                                     |              |
| dryolivine2 |                                                                   |              |          |               | $4.85 \times 10^{-17} Pa^{-n}/s$                                    | 535          |



## S.0.2 Quartz

| ref.                  | Wet<br><i>A</i>                                                                 | Q (kJ/mol)   | <i>V</i> ( ) | <i>n</i>    | Dry<br><i>A</i> | Q (kJ/mol) | <i>V</i> ( ) | <i>n</i> | comment |
|-----------------------|---------------------------------------------------------------------------------|--------------|--------------|-------------|-----------------|------------|--------------|----------|---------|
| [505]                 | $3.2 \times 10^{-4} MPa^{-n}/s$<br>$= 5.072 \cdot 10^{-18} Pa^{-n}/s$           | 154          |              | 2.3         |                 |            |              |          |         |
| [1245]                | $8.57 \cdot 10^{-28} Pa^{-n}/s$<br>$\rightarrow 1.1 \cdot 10^{-28} Pa^{-n}/s$   | 223          | 0            | 4           |                 |            |              |          |         |
| [632]                 | $8.574 \times 10^{-28} Pa^{-s}/s$<br>$\rightarrow 1.1 \cdot 10^{-28} Pa^{-n}/s$ | 26.8         | 0            | 4           |                 |            |              |          |         |
| [613, 614, 295, 482]  | $1.10 \times 10^{-28} Pa^{-s}/s$                                                | 223          | 0            | 4           |                 |            |              |          |         |
| [61]                  | $2.91 \times 10^{-3}$                                                           | 151          |              | 1.8         |                 |            |              |          |         |
| [465]                 | $1.8 \cdot 10^{-8 \pm 2} MPa^{-n}/s$<br>$= 1.8 \cdot 10^{-32 \pm 2} Pa^{-n}/s$  | $137 \pm 34$ | 0            | $4 \pm 0.9$ |                 |            |              |          |         |
|                       | $1.1 \cdot 10^{-4 \pm 2} MPa^{-n}/s$<br>$= 1.1 \cdot 10^{-28 \pm 2} Pa^{-n}/s$  | $223 \pm 56$ | 0            | $4 \pm 0.9$ |                 |            |              |          |         |
| ELEFANT<br>wetquartz1 | $1.1 \cdot 10^{-28 \pm 2} Pa^{-n}/s$                                            | $223 \pm 56$ | 0            | $4 \pm 0.9$ |                 |            |              |          |         |

note that Buiter in [1245] says that 8.57 value is already scaled. Same for [632]. There are indeed, because for  $n = 4$  the multiplicative factor is approx. 7.794, and  $8.574/7.794=1.1$  as in [465].

## S.0.3 Plagioclase

| ref.                       | Wet<br><i>A</i> | Q (kJ/mol) | <i>V</i> ( ) | <i>n</i> | Dry<br><i>A</i>                                                                                              | Q (kJ/mol)        | <i>V</i> ( ) | <i>n</i>          | comment |
|----------------------------|-----------------|------------|--------------|----------|--------------------------------------------------------------------------------------------------------------|-------------------|--------------|-------------------|---------|
| [1040]                     |                 |            |              |          | $3.2 \times 10^{-4} MPa^{-n}/s$<br>$= 3.2 \times 10^{-23.2} Pa^{-n}/s$<br>$= 2.02 \times 10^{-23} Pa^{-n}/s$ | 238<br>238<br>238 |              | 3.2<br>3.2<br>3.2 |         |
| ELEFANT<br>dryplagioclase1 |                 |            |              |          | $2.02 \times 10^{-23} Pa^{-n}/s$                                                                             | 238               |              | 3.2               |         |

## S.0.4 Peridotite

| ref.                  | Wet<br><i>A</i>   | Q (kJ/mol) | <i>V</i> ( ) | <i>n</i> | Dry<br><i>A</i> ( $Pa^{-n}$ ) | Q (kJ/mol) | <i>V</i> ( ) | <i>n</i> | comment |
|-----------------------|-------------------|------------|--------------|----------|-------------------------------|------------|--------------|----------|---------|
| [1040]                | $2.0 \times 10^3$ | 471        |              | 4        | $2.5 \times 10^4$             | 532        |              | 3.5      |         |
| ELEFANT<br>peridotite | $2.0 \times 10^3$ | 471        |              | 4        | $2.5 \times 10^4$             | 532        |              | 3.5      |         |

## S.0.5 Diabase

| ref.       | Wet<br><i>A</i> | Q (kJ/mol) | <i>V</i> ( ) | <i>n</i> | Dry<br><i>A</i> ( $Pa^{-n}$ ) | Q (J/mol)    | <i>V</i> ( ) | <i>n</i>      | comment         |
|------------|-----------------|------------|--------------|----------|-------------------------------|--------------|--------------|---------------|-----------------|
| [295, 482] |                 |            |              |          | $5.04 \times 10^{-28}$        | 485          |              | 4.7           | refers to [820] |
| [820]      |                 |            |              |          |                               | $485 \pm 30$ |              | $4.7 \pm 0.6$ |                 |
| ELEFANT    |                 |            |              |          |                               |              |              |               |                 |

S.0.6 Gabbro

| ref.    | $A$                             | $Q$ (kJ/mol) | $V$ ( ) | $n$ |  | comment          |
|---------|---------------------------------|--------------|---------|-----|--|------------------|
| [1245]  | $1.12 \cdot 10^{-10} Pa^{-n}/s$ | 497          | 0       | 3.4 |  | refers to [1358] |
| ELEFANT |                                 |              |         |     |  |                  |

Looking in [1358] , can't find the number !?!

S.0.7 Serpentine

| ref.  | $A$                             | $Q$ (kJ/mol) | $V$ ( )   | $n$ |  | comment |
|-------|---------------------------------|--------------|-----------|-----|--|---------|
| [570] | $4.47 \cdot 10^{-38} Pa^{-n}/s$ | 8.9          | $3.2cm^3$ | 3.8 |  |         |

# Appendix T

## Invariants

Remember:  $\mathcal{I}_{1,2,3}$  are moment invariants while  $\mathcal{K}_{1,2,3}$  are principal invariants.

### Second invariants

Remembering that the deviatoric stress tensor  $\boldsymbol{\tau}$  is symmetric:

$$\begin{aligned}
 \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \boldsymbol{\tau} : \boldsymbol{\tau} \\
 &= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2 + 2\tau_{xy}^2 + 2\tau_{xz}^2 + 2\tau_{yz}^2) \\
 \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \sum_{ij} \tau_{ij} \tau_{ji} \\
 &= \frac{1}{2} \sum_{ij} \tau_{ij} \tau_{ij} \quad (\boldsymbol{\tau} \text{ is symm}) \\
 &= \frac{1}{2} \boldsymbol{\tau} : \boldsymbol{\tau} \\
 \mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau}] \\
 &= \frac{1}{2} \text{tr} \left[ \begin{pmatrix} \tau_{xx}^2 + \tau_{xy}\tau_{yx} + \tau_{xz}\tau_{zx} & \cdot & \cdot \\ \cdot & \tau_{yx}\tau_{xy} + \tau_{yy}^2 + \tau_{yz}\tau_{zy} & \cdot \\ \cdot & \cdot & \tau_{zx}\tau_{xz} + \tau_{zy}\tau_{yz} + \tau_{zz}^2 \end{pmatrix} \right] \\
 &= \frac{1}{2} \text{tr} \left[ \begin{pmatrix} \tau_{xx}^2 + \tau_{xy}^2 + \tau_{xz}^2 & \cdot & \cdot \\ \cdot & \tau_{xy}^2 + \tau_{yy}^2 + \tau_{yz}^2 & \cdot \\ \cdot & \cdot & \tau_{xz}^2 + \tau_{yz}^2 + \tau_{zz}^2 \end{pmatrix} \right] \quad (\boldsymbol{\tau} \text{ is symm}) \\
 &= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2 + 2\tau_{xy}^2 + 2\tau_{xz}^2 + 2\tau_{yz}^2) \tag{T.1}
 \end{aligned}$$

$$\begin{aligned}
\mathcal{K}_2(\boldsymbol{\sigma}) &= \frac{1}{2}[\text{tr}(\boldsymbol{\sigma})^2 - \text{tr}(\boldsymbol{\sigma}^2)] \\
&= \frac{1}{2}[(\sigma_{xx} + \sigma_{yy} + \sigma_{zz})^2 - (\sigma_{xx}^2 + \sigma_{xy}\sigma_{yx} + \sigma_{xz}\sigma_{zx} + \sigma_{yy}^2 + \sigma_{xy}\sigma_{yx} + \sigma_{yz}\sigma_{zy} + \sigma_{zz}^2 + \sigma_{xz}\sigma_{zx} + \sigma_{yz}\sigma_{zy})] \\
&= \frac{1}{2}[(\sigma_{xx} + \sigma_{yy} + \sigma_{zz})^2 - (\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xy}\sigma_{yx} + 2\sigma_{xz}\sigma_{zx} + 2\sigma_{yz}\sigma_{zy})] \\
&= \frac{1}{2}[\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xx}\sigma_{yy} + 2\sigma_{xx}\sigma_{zz} + 2\sigma_{yy}\sigma_{zz} - (\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xy}\sigma_{yx} + 2\sigma_{xz}\sigma_{zx} + 2\sigma_{yz}\sigma_{zy})] \\
&= \frac{1}{2}[2\sigma_{xx}\sigma_{yy} + 2\sigma_{xx}\sigma_{zz} + 2\sigma_{yy}\sigma_{zz} - (2\sigma_{xy}\sigma_{yx} + 2\sigma_{xz}\sigma_{zx} + 2\sigma_{yz}\sigma_{zy})] \\
&= \sigma_{xx}\sigma_{yy} + \sigma_{xx}\sigma_{zz} + \sigma_{yy}\sigma_{zz} - \sigma_{xy}\sigma_{yx} - \sigma_{xz}\sigma_{zx} - \sigma_{yz}\sigma_{zy} \\
&= \sigma_{xx}\sigma_{yy} + \sigma_{yy}\sigma_{zz} + \sigma_{xx}\sigma_{zz} - \sigma_{xy}^2 - \sigma_{xz}^2 - \sigma_{yz}^2
\end{aligned}$$

Let us now express the second invariant of the deviatoric stress tensor as a function of the invariants of the full stress tensor (just to be sure I have carried this out twice in what follows):

$$\begin{aligned}
\mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \sum_{ij} \tau_{ij} \tau_{ji} \\
&= \frac{1}{2} \sum_{ij} \left( \sigma_{ij} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) \left( \sigma_{ij} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) \\
&= \frac{1}{2} \sum_{ij} \left[ \sigma_{ij} \sigma_{ij} + \sigma_{ij} \left( -\frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) + \sigma_{ij} \left( -\frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) + \left( -\frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) \left( -\frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) \right] \\
&= \frac{1}{2} \sum_{ij} \left[ \sigma_{ij} \sigma_{ij} - \frac{2}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \sigma_{ij} \delta_{ij} + \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 \delta_{ij} \right] \\
&= \frac{1}{2} \sum_{ij} \underbrace{\sigma_{ij} \sigma_{ij}}_{\mathcal{I}_2(\boldsymbol{\sigma})} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \underbrace{\sum_{ij} \sigma_{ij} \delta_{ij}}_{\mathcal{I}_1(\boldsymbol{\sigma})} + \frac{1}{18} \mathcal{I}_1(\boldsymbol{\sigma})^2 \underbrace{\sum_{ij} \delta_{ij}}_3 \\
&= \mathcal{I}_2(\boldsymbol{\sigma}) - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 \\
&= -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma})
\end{aligned}$$

$$\begin{aligned}
\mathcal{I}_2(\boldsymbol{\tau}) &= \frac{1}{2} \boldsymbol{\tau} : \boldsymbol{\tau} \\
&= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2) + \tau_{xy}^2 + \tau_{xz}^2 + \tau_{yz}^2 \\
&= \frac{1}{2} \left( (\sigma_{xx} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}))^2 + (\sigma_{yy} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}))^2 + (\sigma_{zz} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}))^2 \right) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= \frac{1}{2} \left( \sigma_{xx}^2 - \frac{2}{3} \sigma_{xx} \mathcal{I}_1(\boldsymbol{\sigma}) + \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \sigma_{yy}^2 - \frac{2}{3} \sigma_{yy} \mathcal{I}_1(\boldsymbol{\sigma}) + \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \sigma_{zz}^2 - \frac{2}{3} \sigma_{zz} \mathcal{I}_1(\boldsymbol{\sigma}) + \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 \right) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= \frac{1}{2} \left( \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{2}{3} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \mathcal{I}_1(\boldsymbol{\sigma}) + \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma})^2 \right) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= \frac{1}{2} \left( \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{2}{3} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma})^2 \right) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= \frac{1}{2} \left( \sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma})^2 \right) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \frac{1}{2} (\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\
&= -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma})
\end{aligned}$$

So there is no doubt:

$$\mathcal{I}_2(\boldsymbol{\tau}) = -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma})$$

Note that this relationship is often found in a very confusing form where moment invariants  $\mathcal{K}_{1,2,3}$  are used instead of principal invariants  $\mathcal{I}_{1,2,3}$  (although the letter  $I$  is used!). We have established that  $\mathcal{I}_2(\boldsymbol{\sigma}) = \frac{1}{2} \mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma})$  with  $\mathcal{K}_1(\boldsymbol{\sigma}) = \mathcal{I}_1(\boldsymbol{\sigma})$  so that

$$\mathcal{I}_2(\boldsymbol{\tau}) = -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma}) = -\frac{1}{6} \mathcal{K}_1(\boldsymbol{\sigma})^2 + \frac{1}{2} \mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma}) = \frac{1}{3} \mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma})$$

that is:

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma})$$

If one replaces  $\mathcal{K}$ 's by  $\mathcal{I}$ 's then one finds the formula in the literature <sup>1</sup> <sup>2</sup>.

---

<sup>1</sup><https://www.pantelisliolios.com/deviatoric-stress-and-invariants/>

<sup>2</sup>[https://en.wikipedia.org/wiki/Cauchy\\_stress\\_tensor](https://en.wikipedia.org/wiki/Cauchy_stress_tensor)

### third invariants

Let us now look at the third invariant:

$$\begin{aligned}\mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3} \sum_{ijk} \tau_{ij} \tau_{jk} \tau_{ki} \\ &= \frac{1}{3} (\tau_{xx}^3 + \tau_{yy}^3 + \tau_{zz}^3) + \tau_{xx}(\tau_{xy}^2 + \tau_{xz}^2) + \tau_{yy}(\tau_{xy}^2 + \tau_{yz}^2) + \tau_{zz}(\tau_{xz}^2 + \tau_{yz}^2) + 2\tau_{xy}\tau_{xz}\tau_{yz}\end{aligned}$$

$$\begin{aligned}I_3(\boldsymbol{\tau}) &= \det(\boldsymbol{\tau}) \\ &= \tau_{xx}(\tau_{yy}\tau_{zz} - \tau_{yz}^2) - \tau_{yx}(\tau_{xy}\tau_{zz} - \tau_{zy}\tau_{xz}) + \tau_{zx}(\tau_{xy}\tau_{yz} - \tau_{yy}\tau_{xz}) \\ &= \tau_{xx}\tau_{yy}\tau_{zz} - \tau_{xx}\tau_{yz}^2 - \tau_{zz}\tau_{xy}^2 + \tau_{xy}\tau_{yz}\tau_{yz} + \tau_{xy}\tau_{yz}\tau_{yz} - \tau_{yy}\tau_{xz}^2 \\ &= \tau_{xx}\tau_{yy}\tau_{zz} - \tau_{xx}\tau_{yz}^2 - \tau_{zz}\tau_{xy}^2 + 2\tau_{xy}\tau_{yz}\tau_{yz} - \tau_{yy}\tau_{xz}^2 \\ &= \tau_{xx}\tau_{yy}\tau_{zz} - (-\tau_{yy} - \tau_{zz})\tau_{yz}^2 - (-\tau_{xx} - \tau_{yy})\tau_{xy}^2 + 2\tau_{xy}\tau_{yz}\tau_{yz} - (-\tau_{xx} - \tau_{zz})\tau_{xz}^2 \\ &= \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{xx}(\tau_{xy}^2 + \tau_{xz}^2) + \tau_{yy}(\tau_{xy}^2 + \tau_{yz}^2) + \tau_{zz}(\tau_{xz}^2 + \tau_{yz}^2) + 2\tau_{xy}\tau_{xz}\tau_{yz}\end{aligned}$$

The first term is still different than  $\frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^3 + \tau_{zz}^3)$ ... or is it? Let us have a go using the fact that  $\boldsymbol{\tau}$  is deviatoric:

$$\begin{aligned}\tau_{xx}\tau_{yy}\tau_{zz} &= \tau_{xx}(-\tau_{xx} - \tau_{zz})(-\tau_{xx} - \tau_{yy}) \\ &= \tau_{xx}(\tau_{xx}^2 + \tau_{xx}\tau_{yy} + \tau_{xx}\tau_{zz} + \tau_{yy}\tau_{zz}) \\ &= \tau_{xx}^3 + \tau_{xx}^2\tau_{yy} + \tau_{xx}^2\tau_{zz} + \tau_{xx}\tau_{yy}\tau_{zz} \\ \tau_{xx}\tau_{yy}\tau_{zz} &= (-\tau_{yy} - \tau_{zz})\tau_{yy}(-\tau_{xx} - \tau_{yy}) \\ &= \tau_{yy}(\tau_{xx}\tau_{yy} + \tau_{yy}^2 + \tau_{xx}\tau_{zz} + \tau_{yy}\tau_{zz}) \\ &= \tau_{xx}\tau_{yy}^2 + \tau_{yy}^3 + \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{yy}^2\tau_{zz} \\ \tau_{xx}\tau_{yy}\tau_{zz} &= (-\tau_{yy} - \tau_{zz})(-\tau_{xx} - \tau_{zz})\tau_{zz} \\ &= \tau_{zz}(\tau_{xx}\tau_{yy} + \tau_{yy}\tau_{zz} + \tau_{xx}\tau_{zz} + \tau_{zz}^2) \\ &= \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{yy}\tau_{zz}^2 + \tau_{xx}\tau_{zz}^2 + \tau_{zz}^3 \\ \Rightarrow \tau_{xx}\tau_{yy}\tau_{zz} &= \frac{1}{3}(\tau_{xx}\tau_{yy}\tau_{zz} + \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{xx}\tau_{yy}\tau_{zz}) \\ &= \frac{1}{3}(\tau_{xx}^3 + \tau_{xx}^2\tau_{yy} + \tau_{xx}^2\tau_{zz} + \tau_{xx}\tau_{yy}\tau_{zz} \\ &\quad + \tau_{xx}\tau_{yy}^2 + \tau_{yy}^3 + \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{yy}^2\tau_{zz} \\ &\quad + \tau_{xx}\tau_{yy}\tau_{zz} + \tau_{yy}\tau_{zz}^2 + \tau_{xx}\tau_{zz}^2 + \tau_{zz}^3) \\ &= \frac{1}{3}[\tau_{xx}^3 + \tau_{yy}^3 + \tau_{zz}^3 + \tau_{xx}\tau_{yy}(\underbrace{\tau_{xx} + \tau_{yy} + \tau_{zz}}_{=0}) + \tau_{xx}\tau_{zz}(\underbrace{\tau_{xx} + \tau_{yy} + \tau_{zz}}_{=0}) + \tau_{yy}\tau_{zz}(\underbrace{\tau_{xx} + \tau_{yy} + \tau_{zz}}_{=0})] \\ &= \frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^3 + \tau_{zz}^3)\end{aligned}$$

Let us now turn to  $I_3(\boldsymbol{\tau}) = \frac{1}{3}\text{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau} \cdot \boldsymbol{\tau}]$ . Assuming the tensor  $\boldsymbol{\tau}$  to be symmetric then

$$\boldsymbol{\tau} = \begin{pmatrix} a & d & e \\ d & b & f \\ e & f & c \end{pmatrix}$$

then, thanks to <https://www.wolframalpha.com/> I find that

$$\boldsymbol{\tau} \cdot \boldsymbol{\tau} \cdot \boldsymbol{\tau} = \begin{pmatrix} a(a^2 + d^2 + e^2) + d(ad + bd + ef) + e(ae + ce + df) & d(a^2 + d^2 + e^2) + b(ad + bd + ef) + f(ae + ce + df) & e(a^2 + d^2 + e^2) + f(ad + bd + ef) + c(ae + ce + df) \\ a(ad + bd + ef) + d(b^2 + d^2 + f^2) + e(bf + cf + de) & d(ad + bd + ef) + b(b^2 + d^2 + f^2) + f(bf + cf + de) & e(ad + bd + ef) + f(b^2 + d^2 + f^2) + c(bf + cf + de) \\ a(ae + ce + df) + d(bf + cf + de) + e(c^2 + e^2 + f^2) & d(ae + ce + df) + b(bf + cf + de) + f(c^2 + e^2 + f^2) & e(ae + ce + df) + f(bf + cf + de) + c(c^2 + e^2 + f^2) \end{pmatrix}$$

and then

$$\begin{aligned}\frac{1}{3}\mathrm{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau} \cdot \boldsymbol{\tau}] &= \frac{1}{3}(a^3 + b^3 + c^3) + c(e^2 + f^2) + a(d^2 + e^2) + 2def + b(d^2 + f^2) \\ &= \frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^3 + \tau_{zz}^3) + \tau_{xx}(\tau_{xy}^2 + \tau_{xz}^2) + \tau_{yy}(\tau_{xy}^2 + \tau_{yz}^2) + \tau_{zz}(\tau_{xz}^2 + \tau_{yz}^2) + 2\tau_{xy}\tau_{xz}\tau_{yz}\end{aligned}$$



Let us now express the third invariant of the deviatoric stress tensor as a function of the invariants of the full stress tensor (just to be sure I have carried this out twice in what follows):

$$\begin{aligned}
\mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3} \sum_{ijk} \tau_{ij} \tau_{jk} \tau_{ki} \\
&= \frac{1}{3} \sum_{ijk} \left( \sigma_{ij} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \right) \tau_{jk} \tau_{ki} \\
&= \frac{1}{3} \sum_{ijk} \left[ \sigma_{ij} \tau_{jk} \tau_{ki} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \tau_{jk} \tau_{ki} \right] \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \tau_{jk} \tau_{ki} - \frac{1}{9} \sum_{ijk} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ij} \tau_{jk} \tau_{ki} \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \tau_{jk} \tau_{ki} - \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \sum_{ik} \tau_{ik} \tau_{ki} \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \tau_{jk} \tau_{ki} - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \underbrace{\frac{1}{2} \sum_{ik} \tau_{ik} \tau_{ki}}_{\mathcal{I}_2(\boldsymbol{\tau})} \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \tau_{jk} \tau_{ki} - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \left( \sigma_{jk} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{jk} \right) \left( \sigma_{ki} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ki} \right) - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \frac{1}{3} \sum_{ijk} \left( \sigma_{ij} \sigma_{jk} \sigma_{ki} - \sigma_{ij} \sigma_{jk} \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ki} - \sigma_{ij} \sigma_{ki} \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{jk} + \sigma_{ij} \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 \delta_{jk} \delta_{ki} \right) - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \frac{1}{3} \sum_{ijk} \sigma_{ij} \sigma_{jk} \sigma_{ki} - \frac{1}{3} \sum_{ijk} \sigma_{ij} \sigma_{jk} \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{ki} - \frac{1}{3} \sum_{ijk} \sigma_{ij} \sigma_{ki} \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \delta_{jk} + \frac{1}{3} \sum_{ijk} \sigma_{ij} \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma})^2 \delta_{jk} \delta_{ki} - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \underbrace{\frac{1}{3} \sum_{ijk} \sigma_{ij} \sigma_{jk} \sigma_{ki}}_{\mathcal{I}_3(\boldsymbol{\sigma})} - \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \sum_{ij} \sigma_{ij} \sigma_{ji} - \frac{1}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \sum_{ij} \sigma_{ij} \sigma_{ji} + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^2 \underbrace{\sum_{ijk} \sigma_{ij} \delta_{jk} \delta_{ki}}_{\mathcal{I}_1(\boldsymbol{\sigma})} - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \mathcal{I}_3(\boldsymbol{\sigma}) - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \underbrace{\frac{1}{2} \sum_{ij} \sigma_{ij} \sigma_{ji}}_{\mathcal{I}_2(\boldsymbol{\sigma})} - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \underbrace{\frac{1}{2} \sum_{ij} \sigma_{ij} \sigma_{ji}}_{\mathcal{I}_2(\boldsymbol{\sigma})} + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau}) \\
&= \mathcal{I}_3(\boldsymbol{\sigma}) - \frac{4}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\sigma}) + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\tau})
\end{aligned}$$

Then we use  $\mathcal{I}_2(\boldsymbol{\tau}) = -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma})$  so

$$\begin{aligned}
\mathcal{I}_3(\boldsymbol{\tau}) &= \mathcal{I}_3(\boldsymbol{\sigma}) - \frac{4}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2 + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \left( -\frac{1}{6} \mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma}) \right) \\
&= \mathcal{I}_3(\boldsymbol{\sigma}) - \frac{4}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2 + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 + \frac{1}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\sigma}) \\
&= \frac{2}{27} \mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{9} \mathcal{I}_1(\boldsymbol{\sigma}) \mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma})
\end{aligned}$$

We start this time from

$$\mathcal{I}_3(\boldsymbol{\tau}) = \frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^2 + \tau_{zz}^3) + \tau_{xx}(\tau_{xy}^2 + \tau_{xz}^2) + \tau_{yy}(\tau_{xy}^2 + \tau_{yz}^2) + \tau_{zz}(\tau_{xz}^2 + \tau_{yz}^2) + 2\tau_{xy}\tau_{yz}\tau_{yz}$$

We have

$$\begin{aligned}\tau_{xx}^3 &= \left(\sigma_{xx} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\right)^3 \\ &= \sigma_{xx}^3 - 3\sigma_{xx}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{xx}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ \tau_{yy}^3 &= \left(\sigma_{yy} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\right)^3 \\ &= \sigma_{yy}^3 - 3\sigma_{yy}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{yy}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ \tau_{zz}^3 &= \left(\sigma_{zz} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\right)^3 \\ &= \sigma_{zz}^3 - 3\sigma_{zz}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{zz}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3\end{aligned}$$

Then

$$\begin{aligned}\tau_{xx}^3 + \tau_{yy}^2 + \tau_{zz}^3 &= \sigma_{xx}^3 - 3\sigma_{xx}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{xx}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &+ \sigma_{yy}^3 - 3\sigma_{yy}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{yy}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &+ \sigma_{zz}^3 - 3\sigma_{zz}^2\frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}) + 3\sigma_{zz}\frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &= \sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3 - \mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \frac{1}{3}(\sigma_{xx} + \sigma_{yy} + \sigma_{zz})\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &= \sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3 - \mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &= \sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3 - \mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^3\end{aligned}$$

$$\begin{aligned}\mathcal{I}_3(\boldsymbol{\tau}) &= \frac{1}{3}(\tau_{xx}^3 + \tau_{yy}^2 + \tau_{zz}^3) + \tau_{xx}(\tau_{xy}^2 + \tau_{xz}^2) + \tau_{yy}(\tau_{xy}^2 + \tau_{yz}^2) + \tau_{zz}(\tau_{xz}^2 + \tau_{yz}^2) + 2\tau_{xy}\tau_{yz}\tau_{yz} \\ &= \frac{1}{3}(\sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3) - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2) + \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &\quad + (\sigma_{xx} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}))(\sigma_{xy}^2 + \sigma_{xz}^2) + (\sigma_{yy} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}))(\sigma_{xy}^2 + \sigma_{yz}^2) + (\sigma_{zz} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma}))(\sigma_{xz}^2 + \sigma_{yz}^2) + 2\sigma_{xy}\sigma_{yz}\sigma_{yz} \\ &= \frac{1}{3}(\sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3) - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{xy}^2 + \sigma_{yz}^2 + \sigma_{xz}^2 + \sigma_{yz}^2) + \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &\quad + \sigma_{xx}(\sigma_{xy}^2 + \sigma_{xz}^2) + \sigma_{yy}(\sigma_{xy}^2 + \sigma_{yz}^2) + \sigma_{zz}(\sigma_{xz}^2 + \sigma_{yz}^2) + 2\sigma_{xy}\sigma_{yz}\sigma_{yz} \\ &= \frac{1}{3}(\sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3) - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xy}^2 + 2\sigma_{xz}^2 + 2\sigma_{yz}^2) + \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &\quad + \sigma_{xx}(\sigma_{xy}^2 + \sigma_{xz}^2) + \sigma_{yy}(\sigma_{xy}^2 + \sigma_{yz}^2) + \sigma_{zz}(\sigma_{xz}^2 + \sigma_{yz}^2) + 2\sigma_{xy}\sigma_{yz}\sigma_{yz} \\ &= \frac{1}{3}(\sigma_{xx}^3 + \sigma_{yy}^3 + \sigma_{zz}^3) - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\underbrace{\frac{1}{2}(\sigma_{xx}^2 + \sigma_{yy}^2 + \sigma_{zz}^2 + 2\sigma_{xy}^2 + 2\sigma_{xz}^2 + 2\sigma_{yz}^2)}_{\mathcal{I}_2(\boldsymbol{\sigma})} + \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 \\ &\quad + \sigma_{xx}(\sigma_{xy}^2 + \sigma_{xz}^2) + \sigma_{yy}(\sigma_{xy}^2 + \sigma_{yz}^2) + \sigma_{zz}(\sigma_{xz}^2 + \sigma_{yz}^2) + 2\sigma_{xy}\sigma_{yz}\sigma_{yz} \\ &= \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma})\end{aligned}$$

Then, without doubt

$$\mathcal{I}_3(\boldsymbol{\tau}) = \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma})$$

Let us now rewrite this relationship as a function of the principal invariants using the following relationships:

$$\begin{aligned}\mathcal{I}_1(\boldsymbol{\sigma}) &= \mathcal{K}_1(\boldsymbol{\sigma}) \\ \mathcal{I}_2(\boldsymbol{\sigma}) &= \frac{1}{2}\mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma}) \\ \mathcal{I}_3(\boldsymbol{\sigma}) &= \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma})\end{aligned}$$

$$\begin{aligned}\mathcal{I}_3(\boldsymbol{\tau}) &= \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma}) \\ &= \frac{2}{27}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \frac{2}{3}\mathcal{K}_1(\boldsymbol{\sigma})\left(\frac{1}{2}\mathcal{K}_1(\boldsymbol{\sigma})^2 - \mathcal{K}_2(\boldsymbol{\sigma})\right) + \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma}) \\ &= \frac{2}{27}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^3 + \frac{2}{3}\mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma}) \\ &= \frac{2}{27}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma})\end{aligned}\tag{T.5}$$

$$\mathcal{I}_3(\boldsymbol{\tau}) = \frac{2}{27}\mathcal{K}_1(\boldsymbol{\sigma})^3 - \frac{1}{3}\mathcal{K}_1(\boldsymbol{\sigma})\mathcal{K}_2(\boldsymbol{\sigma}) + \mathcal{K}_3(\boldsymbol{\sigma})\tag{T.6}$$

If one replaces the  $\mathcal{K}$ 's by  $I$ 's then one finds the formula in the literature <sup>3 4</sup>.

$$\begin{aligned}\frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} &= \frac{\partial}{\partial \boldsymbol{\sigma}} \left( \frac{2}{27}\mathcal{I}_1(\boldsymbol{\sigma})^3 - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathcal{I}_2(\boldsymbol{\sigma}) + \mathcal{I}_3(\boldsymbol{\sigma}) \right) \\ &= \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 \underbrace{\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}}}_1 - \frac{2}{3} \underbrace{\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}}}_1 \mathcal{I}_2(\boldsymbol{\sigma}) - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \underbrace{\frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}}}_{\boldsymbol{\sigma}} + \underbrace{\frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}}}_{\boldsymbol{\sigma} \cdot \boldsymbol{\sigma}} \\ &= \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 \mathbf{1} - \frac{2}{3}\mathbf{1}\mathcal{I}_2(\boldsymbol{\sigma}) - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\sigma} + \boldsymbol{\sigma} \cdot \boldsymbol{\sigma} \\ &= \left( \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\sigma}) \right) \mathbf{1} - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\sigma} + \boldsymbol{\sigma} \cdot \boldsymbol{\sigma}\end{aligned}\tag{T.7}$$

Using  $\mathcal{I}_2(\boldsymbol{\tau}) = -\frac{1}{6}\mathcal{I}_1(\boldsymbol{\sigma})^2 + \mathcal{I}_2(\boldsymbol{\sigma})$ :

$$\begin{aligned}\frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} &= \left( \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\tau}) - \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 \right) \mathbf{1} - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\sigma} + \boldsymbol{\sigma} \cdot \boldsymbol{\sigma} \\ &= \left( \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\tau}) \right) \mathbf{1} - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\sigma} + \boldsymbol{\sigma} \cdot \boldsymbol{\sigma} \\ &= \left( \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\tau}) \right) \mathbf{1} - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma}) \left( \boldsymbol{\tau} + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathbf{1} \right) + \left( \boldsymbol{\tau} + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathbf{1} \right) \cdot \left( \boldsymbol{\tau} + \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathbf{1} \right) \\ &= \left( \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2 - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\tau}) \right) \mathbf{1} - \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\tau} - \frac{2}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2\mathbf{1} + \boldsymbol{\tau} \cdot \boldsymbol{\tau} + \frac{2}{3}\mathcal{I}_1(\boldsymbol{\sigma})\boldsymbol{\tau} + \frac{1}{9}\mathcal{I}_1(\boldsymbol{\sigma})^2\mathbf{1} \\ &= \boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3}\mathcal{I}_2(\boldsymbol{\tau})\mathbf{1}\end{aligned}\tag{T.8}$$

<sup>3</sup><https://www.pantelisliolios.com/deviatoric-stress-and-invariants/>

<sup>4</sup>[https://en.wikipedia.org/wiki/Cauchy\\_stress\\_tensor](https://en.wikipedia.org/wiki/Cauchy_stress_tensor)

which is the so-called Hill tensor<sup>5</sup>.

Note that this tensor is deviatoric:

$$\mathrm{tr} \left[ \boldsymbol{\tau} \cdot \boldsymbol{\tau} - \frac{2}{3} \mathcal{I}_2(\boldsymbol{\tau}) \mathbf{1} \right] = \mathrm{tr}[\boldsymbol{\tau} \cdot \boldsymbol{\tau}] - 2\mathcal{I}_2(\boldsymbol{\tau}) = 2\mathcal{I}_2(\boldsymbol{\tau}) - 2\mathcal{I}_2(\boldsymbol{\tau}) = 0$$

---

<sup>5</sup>[https://en.wikipedia.org/wiki/Lode\\_coordinates](https://en.wikipedia.org/wiki/Lode_coordinates)

## Derivatives

The derivatives of the invariants with respect to the stress tensor are tensors given as follows:

$$\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \begin{pmatrix} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{xx}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{xy}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{xz}} \\ \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{yx}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{yy}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{yz}} \\ \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{zx}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{zy}} & \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \sigma_{zz}} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \mathbf{1} \quad (\text{T.9})$$

$$\frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} = \begin{pmatrix} \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{xx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{xy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{xz}} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{yx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{yy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{yz}} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{zx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{zy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\sigma})}{\partial \sigma_{zz}} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 2\sigma_{xx} & 2\sigma_{xy} & 2\sigma_{xz} \\ 2\sigma_{yx} & 2\sigma_{yy} & 2\sigma_{yz} \\ 2\sigma_{zx} & 2\sigma_{zy} & 2\sigma_{zz} \end{pmatrix} = \boldsymbol{\sigma} \quad (\text{T.10})$$

$$\begin{aligned} \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \boldsymbol{\sigma}} &= \begin{pmatrix} \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{xx}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{xy}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{xz}} \\ \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{yx}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{yy}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{yz}} \\ \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{zx}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{zy}} & \frac{\partial \mathcal{I}_3(\boldsymbol{\sigma})}{\partial \sigma_{zz}} \end{pmatrix} \\ &= \begin{pmatrix} \sigma_{xx}^2 + \sigma_{xy}^2 + \sigma_{xz}^2 & \sigma_{xx}\sigma_{xy} + \sigma_{yy}\sigma_{xy} + \sigma_{xz}\sigma_{yz} & \sigma_{xx}\sigma_{xz} + \sigma_{zz}\sigma_{xz} + \sigma_{xy}\sigma_{yz} \\ \dots & \sigma_{yy}^2 + \sigma_{xy}^2 + \sigma_{yz}^2 & \dots \\ \dots & \dots & \dots \sigma_{zz}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \end{pmatrix} \\ &= \boldsymbol{\sigma} \cdot \boldsymbol{\sigma} \end{aligned} \quad (\text{T.11})$$

where we have used the generic form of the second and third invariants, i.e. not assuming the tensors to be symmetric so that (for example)  $\sigma_{xz}$  and  $\sigma_{zx}$  are distinct quantities.

The Lodé angle  $\theta_L(\boldsymbol{\tau})$  is actually a function of  $\mathcal{I}_2(\boldsymbol{\tau})$  and  $\mathcal{I}_3(\boldsymbol{\tau})$  as follows:

$$\sin 3\theta_L(\boldsymbol{\tau}) = -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}}$$

Since this quantity unambiguously depends on the deviatoric stress tensor, I will omit the ' $\boldsymbol{\tau}$ ' dependency in what follows. Then

$$\begin{aligned} \frac{\partial}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \sin 3\theta_L &= 3 \cos 3\theta_L(\boldsymbol{\tau}) \frac{\partial \theta_L}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \\ \frac{\partial}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \sin 3\theta_L &= 3 \cos 3\theta_L(\boldsymbol{\tau}) \frac{\partial \theta_L}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \end{aligned}$$

so that

$$\begin{aligned}
\frac{\partial \theta_L}{\partial \mathcal{I}_2(\boldsymbol{\tau})} &= \frac{1}{3 \cos 3\theta_L} \frac{\partial}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \sin 3\theta_L \\
&= \frac{1}{3 \cos 3\theta_L} \frac{\partial}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \left( -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \\
&= \frac{1}{3 \cos 3\theta_L} \left( -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \left( -\frac{3}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \right) \\
&= \frac{1}{3 \cos 3\theta_L} \sin 3\theta_L \left( -\frac{3}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \right) \\
&= -\frac{1}{2} \tan 3\theta_L \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \\
\\
\frac{\partial \theta_L}{\partial \mathcal{I}_3(\boldsymbol{\tau})} &= \frac{1}{3 \cos 3\theta_L} \frac{\partial}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \sin 3\theta_L \\
&= \frac{1}{3 \cos 3\theta_L} \frac{\partial}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \left( -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \\
&= \frac{1}{3 \cos 3\theta_L} \left( -\frac{3\sqrt{3}}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \\
&= \frac{1}{3 \cos 3\theta_L} \left( -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \\
&= \frac{1}{3 \cos 3\theta_L} \sin 3\theta_L \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \\
&= \frac{1}{3} \tan 3\theta_L \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})}
\end{aligned}$$

We have just established the useful relationships

$$\frac{\partial \theta_L}{\partial \mathcal{I}_2(\boldsymbol{\tau})} = -\frac{1}{2} \tan 3\theta_L \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \quad (\text{T.12})$$

$$\frac{\partial \theta_L}{\partial \mathcal{I}_3(\boldsymbol{\tau})} = \frac{1}{3} \tan 3\theta_L \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \quad (\text{T.13})$$

and in the end we can write

$$\begin{aligned}
\frac{\partial \theta_L}{\partial \boldsymbol{\sigma}} &= \frac{\partial \theta_L(\boldsymbol{\tau})}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{\partial \theta_L(\boldsymbol{\tau})}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\
&= \left( -\frac{1}{2} \tan 3\theta_L \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \right) \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \left( \frac{1}{3} \tan 3\theta_L \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \right) \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \\
&= \tan 3\theta_L \left[ -\frac{1}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{3} \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \\
&= \frac{\sin 3\theta_L}{\cos 3\theta_L} \left[ -\frac{1}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{3} \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \\
&= \frac{1}{\cos 3\theta_L} \left( -\frac{3\sqrt{3}}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \right) \left[ -\frac{1}{2} \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{3} \frac{1}{\mathcal{I}_3(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \\
&= -\frac{\sqrt{3}}{2 \cos 3\theta_L} \left[ -\frac{3}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{5/2}} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \quad (\text{T.14})
\end{aligned}$$

i.e.

$$\frac{\partial \theta_L}{\partial \boldsymbol{\sigma}} = -\frac{\sqrt{3}}{2 \cos 3\theta_L} \left[ -\frac{3}{2} \frac{\mathcal{I}_3(\boldsymbol{\tau})}{\mathcal{I}_2(\boldsymbol{\tau})^{5/2}} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} + \frac{1}{\mathcal{I}_2(\boldsymbol{\tau})^{3/2}} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} \right] \quad (\text{T.15})$$

which is Eq. (7.68) of Owen & Hinton:

$$\frac{\partial \theta}{\partial \boldsymbol{\sigma}} = \frac{-\sqrt{3}}{2 \cos 3\theta} \left[ \frac{1}{(J_2')^{3/2}} \frac{\partial J_3}{\partial \boldsymbol{\sigma}} - \frac{3J_3}{(J_2')^2} \frac{\partial (J_2')^{1/2}}{\partial \boldsymbol{\sigma}} \right]. \quad (7.68)$$

Taken from Owen and Hinton [967]

# Appendix U

## The $\Gamma$ tensor in plasticity

WARNING: this is not finished.

Let us start from

$$\vec{\dot{\epsilon}} = \mathbf{\Gamma}(\vec{\sigma}) \cdot \vec{\sigma}$$

Note that we will later need  $\mathbf{\Gamma}^{-1}$  which begs the question of it being invertible... In Zienkiewicz [1422], the author states: “In many forms of the visco-plastic law the relationship  $\vec{\dot{\epsilon}} = \mathbf{\Gamma}(\vec{\sigma}) \cdot \vec{\sigma}$  is such that no volumetric strain rate exists i.e. the material is incompressible. Now  $\mathbf{\Gamma}$  does not posses an inverse”.

### U.0.1 Computing the $\Gamma$ matrix

Let us first establish that we can write quite generally in the isotropic case

$$\frac{\partial Q}{\partial \vec{\sigma}} = \mathbf{\Gamma}_0 \cdot \vec{\sigma}$$

By applying the chain rule we can write

$$\begin{aligned} \frac{\partial Q}{\partial \vec{\sigma}} &= \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \vec{\sigma}} + \frac{\partial Q}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \vec{\sigma}} + \frac{\partial Q}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \vec{\sigma}} \\ &= \left( \frac{\partial Q}{\partial \mathcal{I}_1(\boldsymbol{\sigma})} \mathbf{M}_1(\boldsymbol{\sigma}) + \frac{\partial Q}{\partial \mathcal{I}_2(\boldsymbol{\tau})} \mathbf{M}_2(\boldsymbol{\sigma}) + \frac{\partial Q}{\partial \mathcal{I}_3(\boldsymbol{\tau})} \mathbf{M}_3(\boldsymbol{\sigma}) \right) \cdot \vec{\sigma} \end{aligned} \quad (\text{U.1})$$

All we have to do now is to compute the three symmetric matrices  $\mathbf{M}_{1,2,3}(\boldsymbol{\sigma})$  which independent of  $F$  or  $Q$ .



### Computing matrix $M_1$

$$\begin{aligned}
\frac{\partial \mathcal{I}_1(\boldsymbol{\sigma})}{\partial \vec{\sigma}} &= \frac{\partial}{\partial \vec{\sigma}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\
&= \begin{pmatrix} \frac{\partial}{\partial \sigma_{xx}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\ \frac{\partial}{\partial \sigma_{yy}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\ \frac{\partial}{\partial \sigma_{zz}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\ \frac{\partial}{\partial \sigma_{xy}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\ \frac{\partial}{\partial \sigma_{xz}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \\ \frac{\partial}{\partial \sigma_{yz}} (\sigma_{xx} + \sigma_{yy} + \sigma_{zz}) \end{pmatrix} \\
&= \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \\
&= \frac{1}{\sigma_{xx} + \sigma_{yy} + \sigma_{zz}} \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \cdot \vec{\sigma} \\
&= \underbrace{\frac{1}{\mathcal{I}_1(\boldsymbol{\sigma})} \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}}_{M_1(\boldsymbol{\sigma})} \cdot \vec{\sigma} \tag{U.2}
\end{aligned}$$

### Computing matrix $M_2$

We start from

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{6} [(\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{yy} - \sigma_{zz})^2 + (\sigma_{xx} - \sigma_{zz})^2] + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \tag{U.3}$$

Then

$$\begin{aligned}
\frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \vec{\sigma}} &= \begin{pmatrix} \frac{\partial}{\partial \sigma_{xx}} \mathcal{I}_2(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{yy}} \mathcal{I}_2(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{zz}} \mathcal{I}_2(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{xy}} \mathcal{I}_2(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{xz}} \mathcal{I}_2(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{yz}} \mathcal{I}_2(\boldsymbol{\tau}) \end{pmatrix} \\
&= \begin{pmatrix} \frac{1}{6}(2(\sigma_{xx} - \sigma_{yy}) + 2(\sigma_{xx} - \sigma_{zz})) \\ \frac{1}{6}(-2(\sigma_{xx} - \sigma_{yy}) + 2(\sigma_{yy} - \sigma_{zz})) \\ \frac{1}{6}(-2(\sigma_{yy} - \sigma_{zz}) - 2(\sigma_{xx} - \sigma_{zz})) \\ 2\sigma_{xy} \\ 2\sigma_{xz} \\ 2\sigma_{yz} \end{pmatrix} \\
&= \begin{pmatrix} \frac{2}{3}\sigma_{xx} - \frac{1}{3}\sigma_{yy} - \frac{1}{3}\sigma_{zz} \\ -\frac{2}{3}\sigma_{xx} + \frac{4}{3}\sigma_{yy} - \frac{2}{3}\sigma_{zz} \\ -\frac{2}{3}\sigma_{xx} - \frac{2}{3}\sigma_{yy} + \frac{4}{3}\sigma_{zz} \\ 2\sigma_{xy} \\ 2\sigma_{xz} \\ 2\sigma_{yz} \end{pmatrix} \\
&= \underbrace{\begin{pmatrix} 2/3 & -1/3 & -1/3 & 0 & 0 & 0 \\ -1/3 & 2/3 & -1/3 & 0 & 0 & 0 \\ -1/3 & -1/3 & 2/3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{pmatrix}}_{M_2(\boldsymbol{\sigma})} \cdot \vec{\sigma} \tag{U.4}
\end{aligned}$$

This is the same matrix as in Eq. (13.11) in Zienkiewicz [1422] (1975).

Another look at it using tensors: we start this time from

$$\mathcal{I}_2(\boldsymbol{\tau}) = \frac{1}{6} [(\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{yy} - \sigma_{zz})^2 + (\sigma_{xx} - \sigma_{zz})^2] + \frac{1}{2} (\sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 + \sigma_{yx}^2 + \sigma_{zx}^2 + \sigma_{zy}^2) \tag{U.5}$$

Then

$$\begin{aligned}
\frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \boldsymbol{\sigma}} &= \begin{pmatrix} \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{xz}} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{yz}} \\ \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{zx}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{zy}} & \frac{\partial \mathcal{I}_2(\boldsymbol{\tau})}{\partial \sigma_{zz}} \end{pmatrix} \\
&= \begin{pmatrix} \frac{2}{3}\sigma_{xx} - \frac{1}{3}\sigma_{yy} - \frac{1}{3}\sigma_{zz} & \sigma_{xy} & \sigma_{xz} \\ \sigma_{yx} & -\frac{2}{3}\sigma_{xx} + \frac{4}{3}\sigma_{yy} - \frac{2}{3}\sigma_{zz} & \sigma_{yz} \\ \sigma_{zx} & \sigma_{zy} & -\frac{2}{3}\sigma_{xx} - \frac{2}{3}\sigma_{yy} + \frac{4}{3}\sigma_{zz} \end{pmatrix} \\
&= \boldsymbol{\sigma} - \frac{1}{3} \mathcal{I}_1(\boldsymbol{\sigma}) \mathbf{1} \\
&= \boldsymbol{\tau} \tag{U.6}
\end{aligned}$$

### Computing matrix $M_3$

The third invariant proves to be the most annoying:

$$\mathcal{I}_3(\boldsymbol{\tau}) = \frac{1}{3} \sum_{i,j,k} \tau_{ij} \tau_{jk} \tau_{ki} \quad (\text{U.7})$$

$$\begin{aligned} &= \frac{1}{3} \tau_{xx} (\tau_{xx}^2 + 3\tau_{xy}^2 + 3\tau_{xz}^2) \\ &+ \frac{1}{3} \tau_{yy} (3\tau_{xy}^2 + \tau_{yy}^2 + 3\tau_{yz}^2) \\ &+ \frac{1}{3} \tau_{zz} (3\tau_{xz}^2 + 3\tau_{yz}^2 + \tau_{zz}^2) \\ &+ 2\tau_{xy} \tau_{xz} \tau_{yz} \end{aligned} \quad (\text{U.8})$$

Then

$$\begin{aligned} \frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \vec{\sigma}} &= \begin{pmatrix} \frac{\partial}{\partial \sigma_{xx}} \mathcal{I}_3(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{yy}} \mathcal{I}_3(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{zz}} \mathcal{I}_3(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{xy}} \mathcal{I}_3(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{xz}} \mathcal{I}_3(\boldsymbol{\tau}) \\ \frac{\partial}{\partial \sigma_{yz}} \mathcal{I}_3(\boldsymbol{\tau}) \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial}{\partial \tau_{xx}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{xx}}{\partial \sigma_{xx}} \\ \frac{\partial}{\partial \tau_{yy}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{yy}}{\partial \sigma_{yy}} \\ \frac{\partial}{\partial \tau_{zz}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{zz}}{\partial \sigma_{zz}} \\ \frac{\partial}{\partial \tau_{xy}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{xy}}{\partial \sigma_{xy}} \\ \frac{\partial}{\partial \tau_{xz}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{xz}}{\partial \sigma_{xz}} \\ \frac{\partial}{\partial \tau_{yz}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{\partial \tau_{yz}}{\partial \sigma_{yz}} \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial}{\partial \tau_{xx}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ \frac{\partial}{\partial \tau_{yy}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ \frac{\partial}{\partial \tau_{zz}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ \frac{\partial}{\partial \tau_{xy}} \mathcal{I}_3(\boldsymbol{\tau}) 1 \\ \frac{\partial}{\partial \tau_{xz}} \mathcal{I}_3(\boldsymbol{\tau}) 1 \\ \frac{\partial}{\partial \tau_{yz}} \mathcal{I}_3(\boldsymbol{\tau}) 1 \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial}{\partial \tau_{xx}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ \frac{\partial}{\partial \tau_{yy}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ \frac{\partial}{\partial \tau_{zz}} \mathcal{I}_3(\boldsymbol{\tau}) \frac{2}{3} \\ 2\tau_{xx} \tau_{xy} + 2\tau_{yy} \tau_{xy} + 2\tau_{xz} \tau_{yz} \\ 2\tau_{xx} \tau_{xz} + 2\tau_{zz} \tau_{xz} + 2\tau_{xy} \tau_{yz} \\ 2\tau_{yy} \tau_{yz} + 2\tau_{zz} \tau_{yz} + 2\tau_{xy} \tau_{xz} \end{pmatrix} \quad (\text{U.9}) \end{aligned}$$

$$\begin{aligned} &= \begin{pmatrix} (\tau_{xx}^2 + \tau_{xy}^2 + \tau_{xz}^2) \frac{2}{3} \\ (\tau_{xy}^2 + \tau_{yy}^2 + \tau_{yz}^2) \frac{2}{3} \\ (\tau_{xz}^2 + \tau_{yz}^2 + \tau_{zz}^2) \frac{2}{3} \\ 2\tau_{xx} \tau_{xy} + 2\tau_{yy} \tau_{xy} + 2\tau_{xz} \tau_{yz} \\ 2\tau_{xx} \tau_{xz} + 2\tau_{zz} \tau_{xz} + 2\tau_{xy} \tau_{yz} \\ 2\tau_{yy} \tau_{yz} + 2\tau_{zz} \tau_{yz} + 2\tau_{xy} \tau_{xz} \end{pmatrix} \quad (\text{U.10}) \end{aligned}$$

We have  $\boldsymbol{\tau} = \boldsymbol{\sigma} - \frac{1}{3}\mathcal{I}_1(\boldsymbol{\sigma})\mathbf{1}$  so

$$\tau_{xx} = \sigma_{xx} - \frac{1}{3}\mathcal{I}_1 \quad \Rightarrow \quad \tau_{xx}^2 = (\sigma_{xx} - \frac{1}{3}\mathcal{I}_1)^2 = \sigma_{xx}^2 - \frac{2}{3}\sigma_{xx}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2 \quad (\text{U.11})$$

$$\tau_{yy} = \sigma_{yy} - \frac{1}{3}\mathcal{I}_1 \quad \Rightarrow \quad \tau_{yy}^2 = (\sigma_{yy} - \frac{1}{3}\mathcal{I}_1)^2 = \sigma_{yy}^2 - \frac{2}{3}\sigma_{yy}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2 \quad (\text{U.12})$$

$$\tau_{zz} = \sigma_{zz} - \frac{1}{3}\mathcal{I}_1 \quad \Rightarrow \quad \tau_{zz}^2 = (\sigma_{zz} - \frac{1}{3}\mathcal{I}_1)^2 = \sigma_{zz}^2 - \frac{2}{3}\sigma_{zz}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2 \quad (\text{U.13})$$

Finally

$$\frac{\partial \mathcal{I}_3(\boldsymbol{\tau})}{\partial \vec{\sigma}} = \begin{pmatrix} (\sigma_{xx}^2 - \frac{2}{3}\sigma_{xx}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2 + \sigma_{xy}^2 + \sigma_{xz}^2)\frac{2}{3} \\ (\sigma_{xy}^2 + \sigma_{yy}^2 - \frac{2}{3}\sigma_{yy}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2 + \sigma_{yz}^2)\frac{2}{3} \\ (\sigma_{xz}^2 + \sigma_{yz}^2 + \sigma_{zz}^2 - \frac{2}{3}\sigma_{zz}\mathcal{I}_1 + \frac{1}{9}\mathcal{I}_1^2)\frac{2}{3} \\ 2(\sigma_{xx} - \frac{1}{3}\mathcal{I}_1)\sigma_{xy} + 2(\sigma_{yy} - \frac{1}{3}\mathcal{I}_1)\sigma_{xy} + 2\sigma_{xz}\sigma_{yz} \\ 2(\sigma_{xx} - \frac{1}{3}\mathcal{I}_1)\sigma_{xz} + 2(\sigma_{zz} - \frac{1}{3}\mathcal{I}_1)\sigma_{xz} + 2\sigma_{xy}\sigma_{yz} \\ 2(\sigma_{yy} - \frac{1}{3}\mathcal{I}_1)\sigma_{yz} + 2(\sigma_{zz} - \frac{1}{3}\mathcal{I}_1)\sigma_{yz} + 2\sigma_{xy}\sigma_{xz} \end{pmatrix} \quad (\text{U.14})$$

$$= \begin{pmatrix} (\sigma_{xx}^2 - \frac{2}{3}\sigma_{xx}\mathcal{I}_1 + \sigma_{xy}^2 + \sigma_{xz}^2)\frac{2}{3} \\ (\sigma_{xy}^2 + \sigma_{yy}^2 - \frac{2}{3}\sigma_{yy}\mathcal{I}_1 + \sigma_{yz}^2)\frac{2}{3} \\ (\sigma_{xz}^2 + \sigma_{yz}^2 + \sigma_{zz}^2 - \frac{2}{3}\sigma_{zz}\mathcal{I}_1)\frac{2}{3} \\ 2\sigma_{xx}\sigma_{xy} + 2\sigma_{yy}\sigma_{xy} + 2\sigma_{xz}\sigma_{yz} \\ 2\sigma_{xx}\sigma_{xz} + 2\sigma_{zz}\sigma_{xz} + 2\sigma_{xy}\sigma_{yz} \\ 2\sigma_{yy}\sigma_{yz} + 2\sigma_{zz}\sigma_{yz} + 2\sigma_{xy}\sigma_{xz} \end{pmatrix} + \begin{pmatrix} \frac{1}{9}\mathcal{I}_1^2\frac{2}{3} \\ \frac{1}{9}\mathcal{I}_1^2\frac{2}{3} \\ \frac{1}{9}\mathcal{I}_1^2\frac{2}{3} \\ -\frac{4}{3}\mathcal{I}_1\sigma_{xy} \\ -\frac{4}{3}\mathcal{I}_1\sigma_{xz} \\ -\frac{4}{3}\mathcal{I}_1\sigma_{yz} \end{pmatrix}$$

$$= \begin{pmatrix} (\sigma_{xx}^2 - \frac{2}{3}\sigma_{xx}\mathcal{I}_1 + \sigma_{xy}^2 + \sigma_{xz}^2)\frac{2}{3} \\ (\sigma_{xy}^2 + \sigma_{yy}^2 - \frac{2}{3}\sigma_{yy}\mathcal{I}_1 + \sigma_{yz}^2)\frac{2}{3} \\ (\sigma_{xz}^2 + \sigma_{yz}^2 + \sigma_{zz}^2 - \frac{2}{3}\sigma_{zz}\mathcal{I}_1)\frac{2}{3} \\ 2(\sigma_{xx} + \sigma_{yy})\sigma_{xy} + 2\sigma_{xz}\sigma_{yz} \\ 2(\sigma_{xx} + \sigma_{zz})\sigma_{xz} + 2\sigma_{xy}\sigma_{yz} \\ 2(\sigma_{yy} + \sigma_{zz})\sigma_{yz} + 2\sigma_{xy}\sigma_{xz} \end{pmatrix} + \frac{2}{3}\mathcal{I}_1 \begin{pmatrix} \frac{1}{9}\mathcal{I}_1 \\ \frac{1}{9}\mathcal{I}_1 \\ \frac{1}{9}\mathcal{I}_1 \\ -2\sigma_{xy} \\ -2\sigma_{xz} \\ -2\sigma_{yz} \end{pmatrix} \quad (\text{U.15})$$

$$= \begin{pmatrix} (\sigma_{xx}^2 - \frac{2}{3}\sigma_{xx}\mathcal{I}_1 + \sigma_{xy}^2 + \sigma_{xz}^2)\frac{2}{3} \\ (\sigma_{xy}^2 + \sigma_{yy}^2 - \frac{2}{3}\sigma_{yy}\mathcal{I}_1 + \sigma_{yz}^2)\frac{2}{3} \\ (\sigma_{xz}^2 + \sigma_{yz}^2 + \sigma_{zz}^2 - \frac{2}{3}\sigma_{zz}\mathcal{I}_1)\frac{2}{3} \\ 2(\mathcal{I}_1 - \sigma_{zz})\sigma_{xy} + 2\sigma_{xz}\sigma_{yz} \\ 2(\mathcal{I}_1 - \sigma_{yy})\sigma_{xz} + 2\sigma_{xy}\sigma_{yz} \\ 2(\mathcal{I}_1 - \sigma_{xx})\sigma_{yz} + 2\sigma_{xy}\sigma_{xz} \end{pmatrix} + \frac{2}{3}\mathcal{I}_1 \begin{pmatrix} \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{yy} & \frac{1}{3}\sigma_{zz} \\ \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{yy} & \frac{1}{3}\sigma_{zz} \\ \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{yy} & \frac{1}{3}\sigma_{zz} \\ & & -2 \\ & & -2 \\ & & -2 \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} \quad (\text{U.16})$$

FINISH!!!! not complicated but no rush. Also try to see whether it matches table I of Zienkiewicz

and Cormeau [1423] - note the different ordering of terms in vector

$$\begin{pmatrix} \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{zz} & \frac{1}{3}\sigma_{yy} & -\frac{2}{3}\sigma_{yz} & \frac{1}{3}\sigma_{xz} & \frac{1}{3}\sigma_{xy} \\ \frac{1}{3}\sigma_{zz} & \frac{1}{3}\sigma_{yy} & \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{yz} & -\frac{2}{3}\sigma_{xz} & \frac{1}{3}\sigma_{xy} \\ \frac{1}{3}\sigma_{yy} & \frac{1}{3}\sigma_{xx} & \frac{1}{3}\sigma_{zz} & \frac{1}{3}\sigma_{yz} & \frac{1}{3}\sigma_{xz} & -\frac{2}{3}\sigma_{xy} \\ -\frac{2}{3}\sigma_{yz} & \frac{1}{3}\sigma_{yz} & \frac{1}{3}\sigma_{yz} & -\sigma_{xx} & \sigma_{xy} & \sigma_{xz} \\ \frac{1}{3}\sigma_{xz} & -\frac{2}{3}\sigma_{xz} & \frac{1}{3}\sigma_{xz} & \sigma_{xy} & -\sigma_{yy} & \sigma_{yz} \\ \frac{1}{3}\sigma_{xy} & \frac{1}{3}\sigma_{xy} & -\frac{2}{3}\sigma_{xy} & \sigma_{xz} & \sigma_{yz} & -\sigma_{zz} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{yz} \\ \sigma_{xz} \\ \sigma_{xy} \end{pmatrix} \quad (\text{U.17})$$

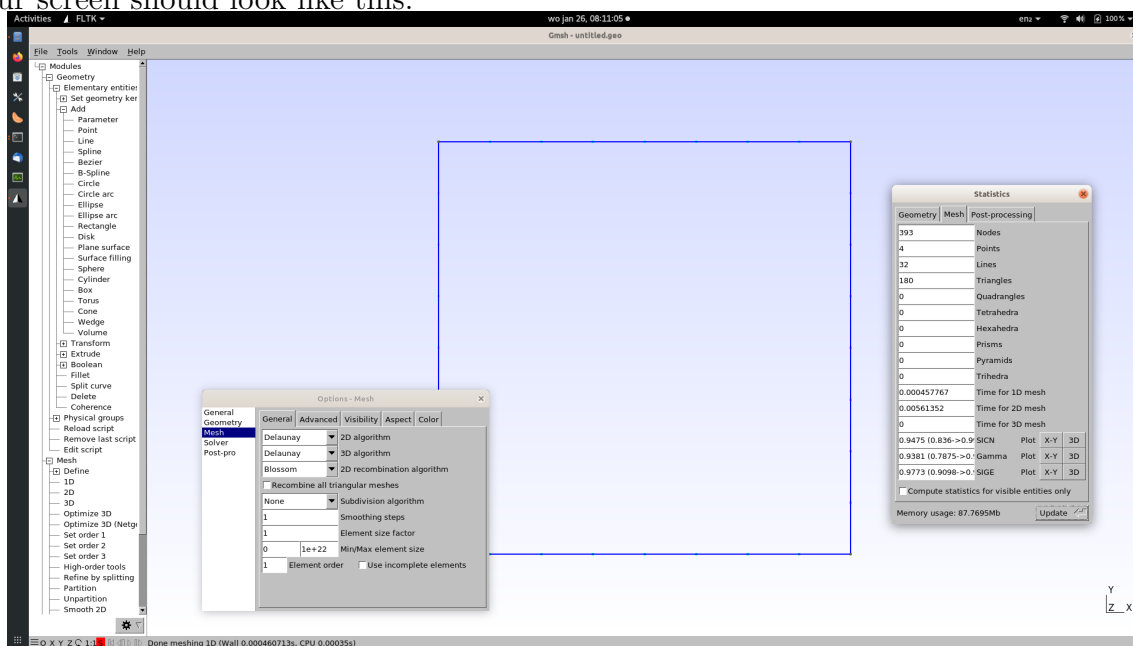
$$= \frac{1}{3} \begin{pmatrix} \sigma_{xx} & \sigma_{zz} & \sigma_{yy} & -2\sigma_{yz} & \sigma_{xz} & \sigma_{xy} \\ \sigma_{zz} & \sigma_{yy} & \sigma_{xx} & \sigma_{yz} & -2\sigma_{xz} & \sigma_{xy} \\ \sigma_{yy} & \sigma_{xx} & \sigma_{zz} & \sigma_{yz} & \sigma_{xz} & -2\sigma_{xy} \\ -2\sigma_{yz} & \sigma_{yz} & \sigma_{yz} & -3\sigma_{xx} & 3\sigma_{xy} & 3\sigma_{xz} \\ \sigma_{xz} & -2\sigma_{xz} & \sigma_{xz} & 3\sigma_{xy} & -3\sigma_{yy} & 3\sigma_{yz} \\ \sigma_{xy} & \sigma_{xy} & -2\sigma_{xy} & 3\sigma_{xz} & 3\sigma_{yz} & -3\sigma_{zz} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{yz} \\ \sigma_{xz} \\ \sigma_{xy} \end{pmatrix} \quad (\text{U.18})$$

$$= \begin{pmatrix} \sigma_{xx}^2 + 2\sigma_{yy}\sigma_{zz} - \sigma_{xy}\sigma_{yz} + \sigma_{xz}^2 \\ \sigma_{yy}^2 + 2\sigma_{xx}\sigma_{zz} + 2\sigma_{xy}\sigma_{yz} - 2\sigma_{xz}^2 \\ \sigma_{zz}^2 + 2\sigma_{xx}\sigma_{yy} - \sigma_{xy}\sigma_{yz} + \sigma_{xz}^2 \\ (-2\sigma_{xx} + \sigma_{yy} + \sigma_{zz})\sigma_{yz} - 3\sigma_{xx}\sigma_{xy} + 3\sigma_{xy}\sigma_{xz} + 3\sigma_{xz}\sigma_{yz} \\ (\sigma_{xx} - 2\sigma_{yy} + \sigma_{zz})\sigma_{xz} + 3\sigma_{xy}^2 - 3\sigma_{yy}\sigma_{xz} + 3\sigma_{yz}^2 \\ (\sigma_{xx} + \sigma_{yy} - 2\sigma_{zz})\sigma_{xy} \\ + 3\sigma_{xy}\sigma_{xz} + 3\sigma_{xz}\sigma_{xz} - 3\sigma_{xz}\sigma_{yz} \end{pmatrix} \quad (\text{U.19})$$

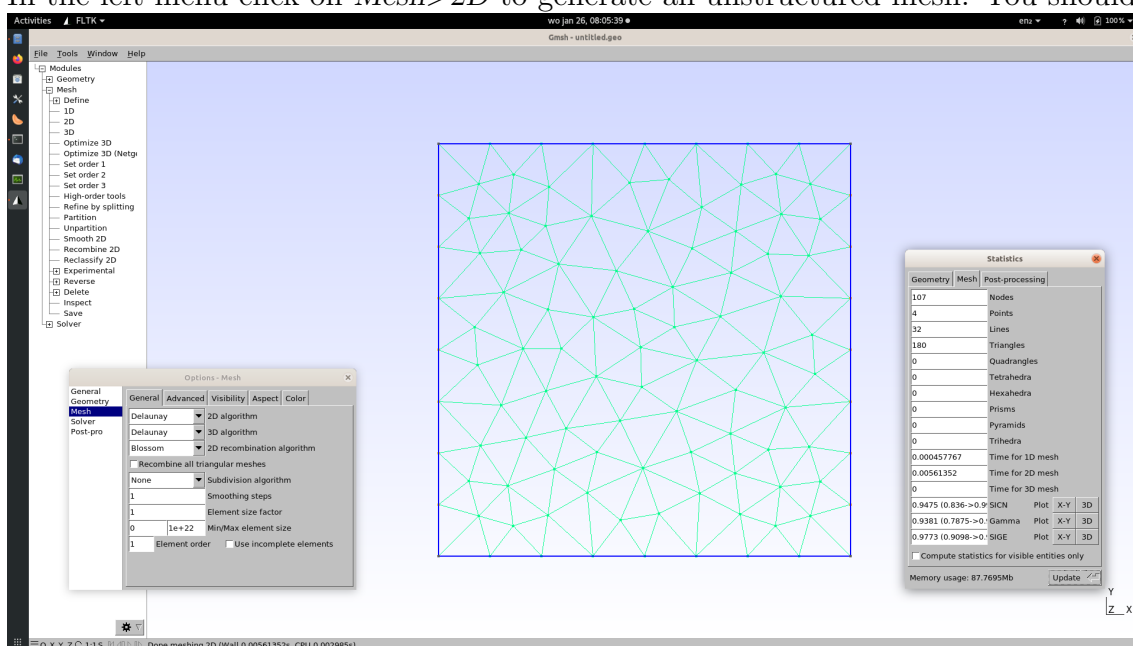
# Appendix V

## Using gmsh

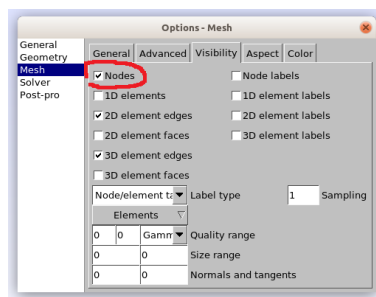
In the left menu *Modules>Geometry>Elementary entities>Add* choose *Rectangle* and input the coordinates of the lower left corner and its size. Then click on *Tools>Options* and *Tools>Statistics*. Your screen should look like this:



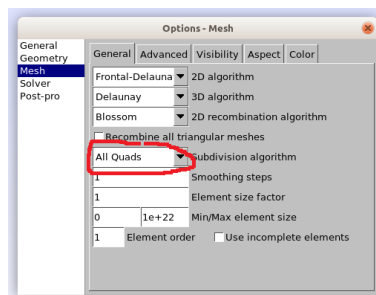
In the left menu click on *Mesh>2D* to generate an unstructured mesh. You should get this:



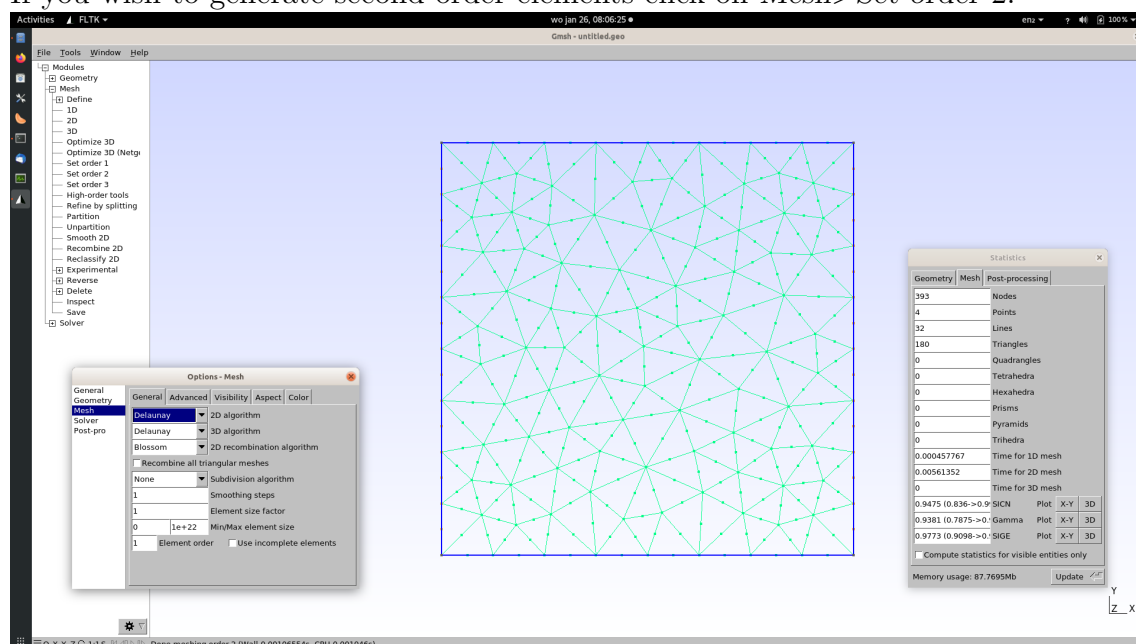
Make sure you can visualise the nodes by setting:



If you wish to use quadrilaterals



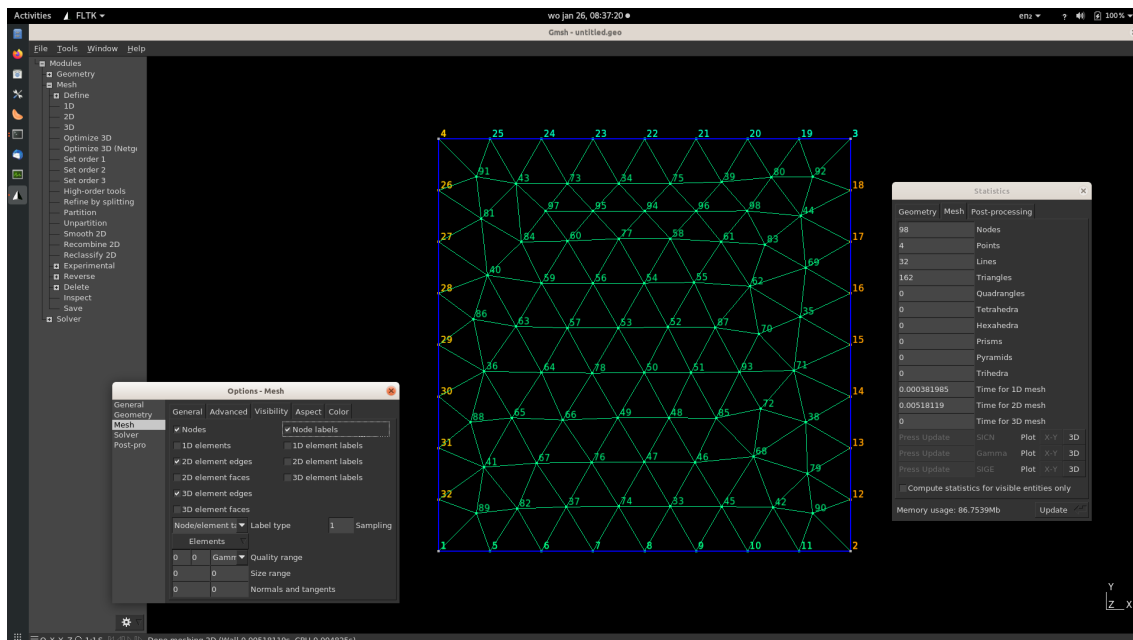
If you wish to generate second order elements click on *Mesh>Set order 2*:



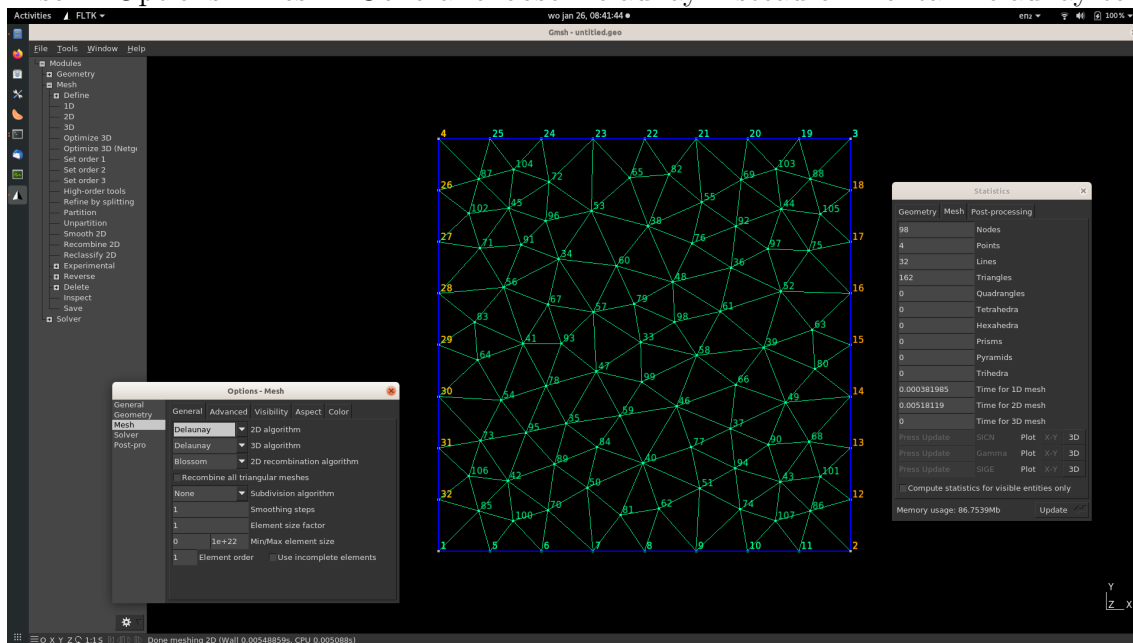
You can refine the mesh by clicking on *Refine by splitting*

Export the mesh: *File>Export*. Choose .mesh format.

In *Options - General - General tab* click on *Use dark interface* and in *Options - Mesh - Visibility tab* click on *Node labels*. Your screen now looks like this:



Also in Options - Mesh - General choose Delaunay instead of Frontal Delaunay to get this



In order to re-generate a mesh, you must wish to click on 1D and then on 2D again.



# Appendix W

## Directional derivative, total and material derivative

The following is a brief review of some basics of differential calculus which underlie many derivations in continuum mechanics.

### W.0.1 Directional derivative

Let's start with a 2-dimensional example of a function  $f$  of the two variables  $x$  and  $y$ , hence  $f(x, y)$ .

Consider an arbitrary point  $(x_0, y_0)$  in the domain of  $f$ . We want to determine the rate of change of  $f$  in any direction in  $(x_0, y_0)$ .

Let  $(x_1, y_1)$  be another point and define the unit vector points from  $(x_0, y_0)$  in the direction of  $(x_1, y_1)$  as:

$$\vec{n} = \begin{pmatrix} n_1 \\ n_2 \end{pmatrix} = \frac{1}{d} \begin{pmatrix} x_1 - x_0 \\ y_1 - y_0 \end{pmatrix} \quad \text{with} \quad d = \sqrt{(x_1 - x_0)^2 + (y_1 - y_0)^2}$$

The line segment connecting the two points can be parameterised as:

$$\begin{aligned} x &= x_0 + sn_1 \\ y &= y_0 + sn_2 \end{aligned} \quad , \quad s \in [0, d] \quad (\text{W.1})$$

Note that  $s$ , the arclength parameter, has the same dimension as the coordinates. On this line the function  $f$  is described by the 1-D function  $f(s) = f(x(s), y(s))$ ,  $s \in [0; d]$ . The rate of change of  $f(x, y)$  in the direction of  $\vec{n}$  at some point  $(x(s), y(s))$  on the line is then  $\partial f(s)/\partial s$ . This derivative can be related to the original coordinates as follows using the change rule of partial differentiation:

$$\frac{df(s)}{ds} = \frac{f(x(s), y(s))}{ds} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial s}$$

Using Eq. (W.1) this gives:

$$\frac{df(s)}{ds} = \frac{\partial f}{\partial x} n_1 + \frac{\partial f}{\partial y} n_2 = \vec{\nabla} f \cdot \vec{n} \quad (\text{W.2})$$

This is the so-called directional derivative which can be computed in any point  $(x, y)$  and direction  $\vec{n}$  as long as the two basic partial derivatives  $\partial f/\partial x$  and  $\partial f/\partial y$ , which give the rate of change in the positive direction of the two axes, respectively, exist in  $(x, y)$ .

Importantly,  $df(s)/ds$  can be directly compared to  $\partial f/\partial x$  and  $\partial f/\partial y$  because all derivatives have the same physical dimension in any application by virtue of the parameterisation (W.1). A change of parameterisation parameter affects the l.h.s. of (W.2) but not the r.h.s because the latter depends

on the  $x$  and  $y$  coordinates and the unit vector  $\mathbf{n}$  (which by definition (W.1) is dimensionless). A change of parameterisation variable from arc-length  $s$  to time  $t = s/v$  will change the l.h.s. into

$$\frac{df}{ds} = \frac{df}{dt} \frac{dt}{ds} = \frac{1}{v} \frac{df}{dt}$$

Substituting this result in (W.2) and rewriting gives another type of directional derivative:

$$\frac{df(t)}{dt} = \vec{\nabla} \cdot \vec{v}$$

where  $u\vec{p}\vec{n}u = \mathbf{v}\vec{n}$  can be interpreted as the local velocity vector, but only if this would be useful in the context of what  $f$  physically represents. Such re-parameterisation is useful in case one explicitly wants to determine the rate-of-change of  $f$  with respect to a parameter different from the arc-length  $s$ . The material derivative of continuum mechanics is an example of such a scaled directional derivative (see below).

## W.0.2 Total differential

We can make the result (W.2) independent of parameterization as follows. By differentiating (W.1) to  $s$ , we find

$$dx(s) = n_1 ds \quad \text{and} \quad dy(s) = n_2 ds \quad (\text{W.3})$$

If we align the differential vector  $d\vec{r} = (dx, dy)^T$  with the line segment we can write  $d\vec{r} = (dx(s), dy(s))^T$ . By using (W.3) we get  $d\vec{r}(s) = \vec{n}ds$  and  $|d\vec{r}| = ds = \sqrt{dx^2 + dy^2}$  because  $\vec{n}$  is of unit length. These relations between  $d\vec{r}$ ,  $dx$ ,  $dy$ , and  $ds$ , all with the same physical dimension, are general used. Next, rewriting (W.2) as  $df(s) = (\vec{\nabla} \cdot \vec{n})ds = \vec{\nabla} \cdot (\vec{n}ds)$  gives:

$$df(s) = \vec{\nabla} f \cdot d\vec{r}(s)$$

Because the line segment  $(x_0, y_0) \rightarrow (x_1, y_1)$  is arbitrarily chosen we can as well write

$$df = \vec{\nabla} f \cdot d\vec{r} = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy \quad (\text{W.4})$$

Equation (W.4) is called the **total differential** of  $f(x, y)$  which holds in each point  $(x, y)$  where the partial derivatives are calculated.

This leads to the following interpretation: Given a function  $f(x, y)$  then in any point  $(x, y)$  in which the partial derivatives  $\partial f/\partial x$  and  $\partial f/\partial y$  exist we can compute the change  $df$  in  $f$  that occurs when going from  $(x, y) \rightarrow (x + dx, y + dy)$  as (W.4), where  $df = f(x + dx, y + dy) - f(x, y)$ . This holds for every choice, including 0 or negative, of the differential steps  $dx$  and  $dy$ .

Generalisation to  $N$ -dimensional space: For any multi-parameter function  $f(x_1, \dots, x_N)$  equation (W.4) generalizes to the total differential

$$df = \frac{\partial f}{\partial x_1} dx_1 + \dots \frac{\partial f}{\partial x_N} dx_N = \vec{\nabla} f^T \cdot d\vec{r}$$

Similarly, equations (W.1) and (W.2) can be generalized to functions of any number of parameters by parameterising the line connecting points  $(x_1^0, \dots, x_N^0)$  and  $(x_1, \dots, x_N)$ :

$$\begin{aligned} x_1 &= x_1^0 + sn_1 \\ &\dots \\ x_N &= x_N^0 + sn_N \end{aligned}$$

with  $s \in [0; d]$  and  $d = \sqrt{(x_1 - x_1^0)^2 + (x_N - x_N^0)^2}$ .

Direction of maximal change: It follows from (W.2) or (W.4) that the change of a function is largest if  $\vec{\nabla} f \cdot d\vec{r}$  is maximum which occurs in any chosen point when  $\vec{\nabla} f$  is parallel to  $d\vec{r}$ . This implies that in every point the gradient vector  $\vec{\nabla} f$  always points in the direction of maximum change of  $f$  and that  $|\vec{\nabla} f| = |df(s)/ds|$  is that maximum change.

*Calculating a normal vector.* Suppose that  $f(x_1, \dots, x_N) = k$  is the level surface of function  $f$  for the constant  $k$  (for example the irregular and time-dependent temperature surface  $T(t, x_1, x_2, x_3) = 20$  degrees in a room full of people). The equation  $f(x_1, \dots, x_N) = k$  implicitly defines the  $(N - 1)$ -dimensional surface in  $N$ -dimensional space of all points for which  $f = k$ . We want to determine in any chosen point of this surface the vector  $\vec{n}$  that is perpendicular to the surface. This is done as follows: Consider (W.4):  $df = \vec{\nabla} f \cdot d\vec{r}$  and take  $d\vec{r}$  to be a step from a point  $(x_1, \dots, x_N)$  on  $f = k$  along the level surface, i.e.  $d\vec{r}$  lies in the level surface. In this case  $df = 0$  because  $f = k$  on the level surface. We find from (W.4) that  $\vec{\nabla} f \cdot d\vec{r} = 0$ , implying that  $\vec{\nabla} f$  is perpendicular to  $d\vec{r}$ . Hence,  $\vec{\nabla} f$  is a vector which is always normal to a level surface. The unit normal in any point  $(x_1, \dots, x_N)$  on the level surface is then calculated as:  $\vec{n} = \vec{\nabla} f / |\vec{\nabla} f|$  where  $\vec{\nabla} f$  is the gradient in that point.

A corollary of this result is that if one calculates  $\vec{\nabla} f$  in some point  $(x_1^0, \dots, x_N^0)$  in  $N$ -space, then one also knows the local direction of the level surface  $f(x_1, \dots, x_N) = f(x_1^0, \dots, x_N^0)$  that passes through  $(x_1^0, \dots, x_N^0)$ .

### W.0.3 The material derivative

In continuum mechanics we distinguish the spatial coordinates  $x_1, x_2, x_3$  and time  $t$ . Hence any function defined on this 4-parameter space is written as  $f(t, x_1, x_2, x_3)$ . The total differential (10) is then

$$df = \frac{\partial f}{\partial t} dt + \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \frac{\partial f}{\partial x_3} dx_3 \quad (\text{W.5})$$

In principle, the time differential  $dt$  and spatial differentials  $dx_i$  can be arbitrarily chosen. For instance, taking  $dt = 0$  that only the spatial changes in the function at fixed time are considered, while taking  $dx_i = 0$  focuses on the temporal variation in a chosen fixed point. Generally, in continuum mechanics a special choice is made for the directional derivative, which involves the local direction of the flow. This direction is given at any point  $(x_1, x_2, x_3)$  and any time  $t$  by the velocity vector

$$\vec{v}(t, x_1, x_2, x_3) = \frac{d\vec{r}}{dt} \quad (\text{W.6})$$

where  $d\vec{r}$  is the spatial step taken by a flow particle from  $(x_1, x_2, x_3) \rightarrow (x_1 + dx_1, x_2 + dx_2, x_3 + dx_3)$  during the time interval  $t \rightarrow t + dt$ . Hence, the flow direction  $d\vec{r}$  in point  $(x_1, x_2, x_3)$  depends on the time  $t$  such that  $d\vec{r}(t) = v(t, x_1, x_2, x_3)dt$ .

Taking  $\vec{v} = (\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3)^T$  equation (W.5) becomes

$$df = \frac{\partial f}{\partial t} dt + \frac{\partial f}{\partial x_1} \mathbf{v}_1 dt + \frac{\partial f}{\partial x_2} \mathbf{v}_2 dt + \frac{\partial f}{\partial x_3} \mathbf{v}_3 dt$$

Dividing by  $\delta t$  yields

$$\frac{Df}{Dt} = \frac{\partial f}{\partial t} + \frac{\partial f}{\partial x_1} \mathbf{v}_1 + \frac{\partial f}{\partial x_2} \mathbf{v}_2 + \frac{\partial f}{\partial x_3} \mathbf{v}_3 = \frac{\partial f}{\partial t} + \vec{v} \cdot \vec{\nabla} f$$

This equation is called the material derivative of  $f$  and describes the rate of change of  $f$  with time in the local direction of the flow.

FINISH

## W.0.4 Material derivative of a volume integral

Let  $F(\vec{r}, t)$  be some scalar function depending on spatial coordinates  $\vec{r}$  and time  $t$  and  $V(t)$  a volume that may also depend on  $t$ . Define the volume integral

$$I(t) = \int_{V(t)} F(\vec{r}, t) dV.$$

For example, if  $F$  is density, then  $I(t)$  is the mass contained in the volume.

Assume a deforming medium with incremental displacement field  $\vec{s}(\vec{r}, t)$ . Consider the deformation that occurs between  $t$  and  $t + \Delta t$  in which  $\Delta t$  is a very small time step such that we can write that a particle at position  $\vec{r}$  at time  $t$  will be displaced to  $\vec{r} + \Delta\vec{r}$  at  $t + \Delta t$ .

Then

$$\Delta\vec{r} = \vec{s}(\vec{r}, t + \Delta t) - \vec{s}(\vec{r}, t) = \frac{\vec{s}(\vec{r}, t + \Delta t) - \vec{s}(\vec{r}, t)}{\Delta t} \Delta t = \vec{v} \Delta t$$

where  $\vec{v} = d\vec{s}/dt$  is the velocity vector at  $(\vec{r}, t)$ .

The volume  $V(t)$  will deform to  $V'(t + \Delta t)$ . The material derivative of  $I(t)$  is defined as:

$$\frac{DI}{Dt} = \frac{D}{Dt} \int_{V(t)} F(\vec{r}, t) dV = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \left[ \int_{V'(t+\Delta t)} F(\vec{r} + \Delta\vec{r}, t + \Delta t) dV' - \int_{V(t)} F(\vec{r}, t) dV \right] \quad (\text{W.7})$$

For incremental  $\Delta t$  we can approximate

$$F(\vec{r} + \Delta\vec{r}, t + \Delta t) = F(\vec{r}, t) + \frac{\partial F}{\partial x_j} v_j \Delta t + \frac{\partial F}{\partial t} \Delta t = F(\vec{r}, t) + \frac{DF}{Dt} \Delta t$$

Further, from continuum mechanics we have for the volume change associated with the incremental displacement field  $\vec{s}(\vec{r}, t)$ :

$$\frac{dV' - dV}{dV'} = \vec{\nabla} \cdot \vec{s} = \vec{\nabla} \cdot (\vec{v} \Delta t)$$

or

$$dV' = \left( 1 + \frac{\partial v_j}{\partial x_j} \Delta t \right) dV$$

Using these results

$$F(\vec{r} + \Delta\vec{r}, t + \Delta t) dV' = \left( F(\vec{r}, t) + \frac{DF}{Dt} \Delta t \right) \left( 1 + \frac{\partial v_j}{\partial x_j} \Delta t \right) dV = F(\vec{r}, t) dV + \frac{DF}{Dt} \Delta t dV + F(\vec{r}, t) \frac{\partial v_j}{\partial x_j} \Delta t dV + \frac{DF}{Dt} \frac{\partial v_j}{\partial x_j} \Delta t^2 dV$$

such that now the integration over  $V'$  can be replaced by an integration over  $V$ :

$$\int_{V'(t+\Delta t)} F(\vec{r} + \Delta\vec{r}, t + \Delta t) dV' \simeq \int_{V(t)} F(\vec{r}, t) dV + \int_{V(t)} \left[ \left( \frac{DF}{Dt} + \frac{\partial v_j}{\partial x_j} \right) \Delta t + \frac{DF}{Dt} \frac{\partial v_j}{\partial x_j} \Delta t^2 \right] dV$$

Substituting this result in the above definition of the material derivative

$$\frac{DI}{Dt} = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \left[ \int_{V(t)} \left( \frac{DF}{Dt} + \frac{\partial v_j}{\partial x_j} \right) \Delta t + \frac{DF}{Dt} \frac{\partial v_j}{\partial x_j} \Delta t^2 \right] dV$$

This leads to material derivative of a volume integral:

$$\frac{D}{Dt} \int_{V(t)} F dV = \int_{V(t)} \frac{DF}{Dt} + F \frac{\partial v_j}{\partial x_j} dV$$

# Bibliography

- [1] B.T. Aagaard, M.G. Knepley, and C.A. Williams. “A domain decomposition approach to implementing fault slip in finite-element models of quasi-static and dynamic crustal deformation”. In: *J. Geophys. Res.* 118 (2013), pp. 3059–3079. DOI: 10.1002/jgrb.50217.
- [2] F.J Adewale, A.P. Lucky, A.P. Oluwabunmi, and E.F. Boluwaji. “Selecting the most appropriate model for rheological characterization of synthetic based drilling mud”. In: *International Journal of Applied Engineering Research* 12.18 (2017), pp. 7614–7649. DOI: xxxx.
- [3] Deniz Tolga Akcabay, David R Dowling, and William W Schultz. “Clothes washing simulations”. In: *Computers & Fluids* 100 (2014), pp. 79–94. DOI: 10.1016/j.compfluid.2014.05.005.
- [4] Francis Albarède and Rob D van der Hilst. “Zoned mantle convection”. In: *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences* 360.1800 (2002), pp. 2569–2592. DOI: 10.1098/rsta.2002.1081.
- [5] M. Albers. “A local mesh refinement multigrid method for 3D convection problems with strongly variable viscosity”. In: *J. Comp. Phys.* 160 (2000), pp. 126–150.
- [6] Leandro R Alejano and Antonio Bobet. “Drucker–prager criterion”. In: *The ISRM Suggested Methods for Rock Characterization, Testing and Monitoring: 2007-2014*. 2012, pp. 247–252. DOI: 10.1007/s00603-012-0278-2.
- [7] L. Alisic, M. Gurnis, G. Stadler, C. Burstedde, and O. Ghattas. “Multi-scale dynamics and rheology of mantle flow with plates”. In: *J. Geophys. Res.* 117 (2012). DOI: 10.1029/2012JB009234.
- [8] Richard B Alley. “Flow-law hypotheses for ice-sheet modeling”. In: *Journal of Glaciology* 38.129 (1992), pp. 245–256. DOI: 10.3189/S0022143000003658.
- [9] V. Allken, R. Huismans, and C. Thieulot. “Factors controlling the mode of rift interaction in brittle-ductile coupled systems: a 3D numerical study”. In: *Geochem. Geophys. Geosyst.* 13.5 (2012), Q05010. DOI: 10.1029/2012GC004077.
- [10] V. Allken, R. Huismans, and C. Thieulot. “Three dimensional numerical modelling of upper crustal extensional systems”. In: *J. Geophys. Res.* 116 (2011), B10409. DOI: 10.1029/2011JB008319.
- [11] V. Allken, R.S. Huismans, H. Fossen, and C. Thieulot. “3D numerical modelling of graben interaction and linkage: a case study of the Canyonlands grabens, Utah”. In: *Basin Research* 25 (2013), pp. 1–14. DOI: 10.1111/bre.12010.
- [12] P.R. Amestoy, A. Buttari, J.-Y. L’Excellent, and T. Mary. “Performance and Scalability of the Block Low-Rank Multifrontal Factorization on Multicore Architectures”. In: *ACM Transactions on Mathematical Software* 45 (1 2019), 2:1–2:26. DOI: 10.1145/3242094.
- [13] P.R. Amestoy and I.S. Duff. “Vectorization of a multiprocessor multifrontal code”. In: *International Journal of Supercomputer Applications* 3 (1989), pp. 41–59. DOI: 10.1177/109434208900300303.

- [14] P.R. Amestoy, I.S. Duff, J.Koster, and J.-Y. L'Excellent. "A fully asynchronous multifrontal solver using distributed dynamic scheduling". In: *SIAM Journal of Matrix Analysis and Applications* 23.1 (2001), pp. 15–41.
- [15] P.R. Amestoy, I.S. Duff, and J.-Y. L'Excellent. "Multifrontal parallel distributed symmetric and unsymmetric solvers". In: *Computer Methods in Applied Mechanics and Engineering* 184 (2000), pp. 501–520.
- [16] P.R. Amestoy, A. Guermouche, J.-Y. L'Excellent, and S. Pralet. "Hybrid scheduling for the parallel solution of linear systems". In: *Parallel Computing* 32.2 (2006), pp. 136–156. DOI: 10.1016/j.parco.2005.07.004.
- [17] Christophe Ancey and Steve Cochard. "The dam-break problem for Herschel–Bulkley viscoplastic fluids down steep flumes". In: *Journal of Non-Newtonian Fluid Mechanics* 158.1-3 (2009), pp. 18–35. DOI: 10.1016/j.jnnfm.2008.08.008.
- [18] Edward Anders and Nicolas Grevesse. "Abundances of the elements: Meteoritic and solar". In: *Geochimica et Cosmochimica acta* 53.1 (1989), pp. 197–214. DOI: 10.1016/0016-7037(89)90286-X.
- [19] B.D. Anderson. "AUTOMATED ALL-QUADRILATERAL MESH ADAPTATION THROUGH REFINEMENT AND COARSENING". PhD thesis. Department of Civil and Environmental Engineering, Brigham Young University, 2009.
- [20] Bret D Anderson, Steven E Benzley, and Steven J Owen. "Automatic all quadrilateral mesh adaption through refinement and coarsening". In: *Proceedings of the 18th international meshing roundtable*. Springer, 2009, pp. 557–574. DOI: 10.1007/978-3-642-04319-2\_32.
- [21] C.A. Anderson and R.J. Bridwell. "A finite element method for studying the transient non-linear thermal creep of geological structures". In: *International Journal for Numerical and Analytical Methods in Geomechanics* 4 (1980), pp. 255–276.
- [22] J.D. Anderson. *Computational Fluid Dynamics*. McGraw-Hill, 1995.
- [23] Thomas Andolfsson. "Analyses of thermal conductivity from mineral composition and analyses by use of Thermal Conductivity Scanner: a study of thermal properties of Scanian rock types". In: *MSc thesis - Dissertations in Geology at Lund University* (2013).
- [24] Miguel Andrés-Martínez, Jason P Morgan, Marta Pérez-Gussinyé, and Lars Rüpke. "A new free-surface stabilization algorithm for geodynamical modelling: Theory and numerical tests". In: *Physics of the Earth and Planetary Interiors* 246 (2015), pp. 41–51. DOI: 10.1016/j.pepi.2015.07.003.
- [25] Miguel Andrés-Martínez, Marta Pérez-Gussinyé, John Armitage, and Jason P Morgan. "Thermomechanical Implications of Sediment Transport for the Architecture and Evolution of Continental Rifts and Margins". In: *Tectonics* 38.2 (2019), pp. 641–665. DOI: 10.1029/2018TC005346.
- [26] D.N. Arnold, F. Brezzi, and M. Fortin. "A stable finite element for the Stokes equation". In: *Calcolo* XXI.IV (1984), pp. 337–344. DOI: 10.1007/BF02576171.
- [27] Douglas Arnold, Daniele Boffi, and Richard Falk. "Approximation by quadrilateral finite elements". In: *Mathematics of computation* 71.239 (2002), pp. 909–922. DOI: 10.1090/S0025-5718-02-01439-4.
- [28] Douglas N Arnold. "On nonconforming linear-constant elements for some variants of the Stokes equations". In: *Istit. Lombardo Accad. Sci. Lett. Rend. A* 127.1 (1993). DOI: xxxx.
- [29] P.A. Arrial, N. Flyer, G.B. Wright, and L.H. Kellogg. "On the sensitivity of 3-D thermal convection codes to numerical discretization: a model intercomparison". In: *Geosci. Model Dev.* 7 (2014), pp. 2065–2076. DOI: 10.5194/gmd-7-2065-2014.

- [30] M. Arroyo, C. Couder-Castaneda, A. Trujillo-Alcantara, I.E. Herrera-Diaz, and N. Vera-Chavez. “A Performance Study of a Dual Xeon-Phi Cluster for the Forward Modelling of Gravitational Fields”. In: *Scientific Programming* (2015). DOI: 10.1155/2015/316012.
- [31] Alireza Asgari and LN Moresi. “Multiscale particle-in-cell method: from fluid to solid mechanics”. In: *Advanced Methods for Practical Applications in Fluid Mechanics. InTech, Croatia* (2012), pp. 185–208.
- [32] M.F. Asgharzadeh, R.R.B. Von Frese, H.R. Kim, T.E. Leftwich, and J.W. Kim. “Spherical prism gravity effects by Gauss-Legendre quadrature integration”. In: *Geophysical Journal International* 169.1 (2007), pp. 1–11. DOI: 10.1111/j.1365-246X.2007.03214.x.
- [33] C Auth and H Harder. “Multigrid solution of convection problems with strongly variable viscosity”. In: *Geophysical Journal International* 137.3 (1999), pp. 793–804. DOI: 10.1046/j.1365-246x.1999.00833.x.
- [34] EH Ayachour. “A fast implementation for GMRES method”. In: *Journal of Computational and Applied Mathematics* 159.2 (2003), pp. 269–283.
- [35] A Yu Babeyko, Stephan V Sobolev, RB Trumbull, Onno Oncken, and LL Lavier. “Numerical models of crustal scale convection and partial melting beneath the Altiplano–Puna plateau”. In: *Earth and Planetary Science Letters* 199.3-4 (2002), pp. 373–388. DOI: 10.1016/S0012-821X(02)00597-6.
- [36] Andrey Y Babeyko, Stephan V Sobolev, Tim Vietor, Onno Oncken, and Robert B Trumbull. “Numerical study of weakening processes in the central Andean back-arc”. In: *The Andes*. 2006, pp. 495–512. DOI: 10.1007/978-3-540-48684-8\_24.
- [37] Ivo Babuška. “Error-bounds for finite element method”. In: *Numerische Mathematik* 16.4 (1971), pp. 322–333. DOI: 10.1007/BF02165003.
- [38] W. Bai. “The quadrilateral Mini finite element for the Stokes problem”. In: *Computer Methods in Applied Mechanics and Engineering* 143 (1997), pp. 41–47. DOI: 10.1016/S0045-7825(96)01146-2.
- [39] R.R. Bakker, M. Frehner, and M. Lupi. “How temperature-dependent elasticity alters host rock/magmatic reservoir models: A case study on the effects of ice-cap unloading on shallow volcanic systems”. In: *epsl* 456 (2016), pp. 16–25. DOI: 10.1016/j.epsl.2016.09.039.
- [40] N.J. Balmforth, Y. Forterre, and O. Pouliquen. “The viscoplastic Stokes layer”. In: *Journal of Non-Newtonian Fluid Mechanics* 158 (2009), pp. 46–53. DOI: 10.1016/j.jnnfm.2008.07.008.
- [41] N.J. Balmforth and A.C. Rust. “Weakly nonlinear viscoplastic convection”. In: *Journal of Non-Newtonian Fluid Mechanics* 158 (2009), pp. 36–45. DOI: 10.1016/j.jnnfm.2008.07.012.
- [42] NJ Balmforth, RV Craster, P Perona, AC Rust, and R Sassi. “Viscoplastic dam breaks and the Bostwick consistometer”. In: *Journal of non-newtonian fluid mechanics* 142.1-3 (2007), pp. 63–78. DOI: 10.1016/j.jnnfm.2006.06.005.
- [43] Wolfgang Bangerth, Imbunm Kim, Dongwoo Sheen, and Jaeryun Yim. “On Hanging Node Constraints for Nonconforming Finite Elements using the Douglas–Santos–Sheen–Ye Element as an Example”. In: *SIAM Journal on Numerical Analysis* 55.4 (2017), pp. 1719–1739. DOI: 10.1137/16M1071432.
- [44] Wolfgang Bangerth et al. *ASPECT: Advanced Solver for Problems in Earth’s ConvecTion, User Manual*. doi:10.6084/m9.figshare.4865333. July 2022. DOI: 10.6084/m9.figshare.4865333. URL: <https://doi.org/10.6084/m9.figshare.4865333>.

- [45] H.A. Barnes. “The yield stress - everything flows?” In: *J. Non-Newtonian Fluid Mech.* 81 (1999), pp. 133–178.
- [46] HA Barnes and K Walters. “The yield stress myth?” In: *Rheologica acta* 24.4 (1985), pp. 323–326.
- [47] Terence D Barr and Gregory A Houseman. “Deformation fields around a fault embedded in a non-linear ductile medium”. In: *Geophysical Journal International* 125.2 (1996), pp. 473–490. DOI: 10.1111/j.1365-246X.1996.tb00012.x.
- [48] R. Barrett et al. *Templates for the solution of linear systems: building blocks for iterative methods*. SIAM, 1994.
- [49] TL Barry et al. “Whole-mantle convection with tectonic plates preserves long-term global patterns of upper mantle geochemistry”. In: *Scientific reports* 7.1 (2017), pp. 1–9. DOI: 10.1038/s41598-017-01816-y.
- [50] T. Barth, P. Bochev, M. Gunzburger, and J. Shadid. “A taxonomy of consistently stabilised finite element methods for the Stokes problem”. In: *SIAM J. Sci. Comput.* 25.5 (2004), pp. 1585–1607. DOI: 10.1137/S1064827502407718.
- [51] PJ Barton. “The relationship between seismic velocity and density in the continental crust - a useful constraint?” In: *Geophysical Journal International* 87.1 (1986), pp. 195–208.
- [52] G.K. Batchelor. *An introduction to fluid dynamics*. Cambridge University Press, 1967.
- [53] K.-J. Bathe. *Finite Element Procedures in Engineering Analysis*. Prentice-Hall, 1982.
- [54] GE Batt and J Braun. “On the thermomechanical evolution of compressional orogens”. In: *Geophysical Journal International* 128.2 (1997), pp. 364–382. DOI: 10.1111/j.1365-246X.1997.tb01561.x.
- [55] L. Battaglia, M.A. Storti, and J. D’Elia. “An interface capturing finite element approach for free surface flows using unstructured grids”. In: *Mecanica Computacional XXVII* (2008), pp. 33–48. DOI: xxxxx.
- [56] Simon Bauer et al. “TerraNeo—Mantle Convection Beyond a Trillion Degrees of Freedom”. In: *Software for Exascale Computing - SPPEXA 2016-2019*. Ed. by Hans-Joachim Bungartz, Severin Reiz, Benjamin Uekermann, Philipp Neumann, and Wolfgang E. Nagel. Springer International Publishing, 2020, pp. 569–610. DOI: 10.1007/978-3-030-47956-5\_19.
- [57] J.R. Baumgardner. “Three-Dimensional treatment of convective flow in the Earth’s mantle”. In: *Journal of Statistical Physics* 39.5/6 (1985), pp. 501–511. DOI: 10.1007/BF01008348.
- [58] J.R. Baumgardner and P.O. Frederickson. “Isocahedral discretisation of the two-sphere”. In: *SIAM J. Numer Anal.* 22.6 (1985), pp. 1107–1115. DOI: 10.1137/0722066.
- [59] C. Beaumont, P. Fullsack, and J. Hamilton. “Erosional control of active compressional orogens”. In: *Thrust Tectonics* 99 (1992), pp. 1–18. DOI: 10.1007/978-94-011-3066-0\_1.
- [60] C. Beaumont, P. Fullsack, and J. Hamilton. “Styles of crustal deformation in compressional orogens caused by subduction of the underlying lithosphere”. In: *Tectonophysics* 232 (1994), pp. 119–132. DOI: 10.1016/0040-1951(94)90079-5.
- [61] C. Beaumont, J.A. Munoz, J. Hamilton, and P. Fullsack. “Factors controlling the Alpine evolution of the central Pyrenees inferred from a comparison of observations and geodynamical models”. In: *J. Geophys. Res.* 105 (2000), pp. 8121–8145. DOI: 10.1029/1999JB900390.
- [62] Christopher Beaumont and Anthony Lambert. “Crustal structure from surface load tilts, using a finite element model”. In: *Geophysical Journal International* 29.2 (1972), pp. 203–226.



- [63] Christopher Beaumont and Garry Quinlan. “A geodynamic framework for interpreting crustal-scale seismic-reflectivity patterns in compressional orogens”. In: *Geophysical Journal International* 116.3 (1994), pp. 754–783. DOI: 10.1111/j.1365-246X.1994.tb03295.x.
- [64] J.M. Becker and M. Bevis. “Love’s problem”. In: *Geophy. J. Int.* 156 (2004), pp. 171–178. DOI: 10.1111/j.1365-246X.2003.02150.x.
- [65] T.W. Becker. “On the effect of temperature and strain-rate dependent viscosity on global mantle flow, net rotation, and plate-driving forces”. In: *Geophy. J. Int.* 167 (2006), pp. 943–957.
- [66] T.W. Becker and B.J.P. Kaus. *Numerical Geodynamics. v1.1*. Tech. rep. University Southern California, 2010.
- [67] Thorsten W Becker and Richard J O’Connell. “Predicting plate velocities with mantle circulation models”. In: *Geochemistry, Geophysics, Geosystems* 2.12 (2001).
- [68] Mark D Behn, David L Goldsby, and Greg Hirth. “The role of grain size evolution in the rheology of ice: implications for reconciling laboratory creep data and the Glen flow law”. In: *The Cryosphere* 15.9 (2021), pp. 4589–4605. DOI: 10.5194/tc-15-4589-2021.
- [69] M. Behr. “On the Application of Slip Boundary Condition on Curved Boundaries”. In: *Int. J. Num. Meth. Fluids* 45 (2004), pp. 43–51. DOI: 10.1002/flid.663.
- [70] Léa Bello, Nicolas Coltice, Tobias Rolf, and Paul J Tackley. “On the predictability limit of convection models of the Earth’s mantle”. In: *Geochemistry, Geophysics, Geosystems* 15.6 (2014), pp. 2319–2328. DOI: 10.1002/2014GC005254.
- [71] David J Benson. “An efficient, accurate, simple ALE method for nonlinear finite element programs”. In: *Computer methods in applied mechanics and engineering* 72.3 (1989), pp. 305–350.
- [72] M. Benzi, G.H. Golub, and J. Liesen. “Numerical solution of saddle point problems”. In: *Acta Numerica* 14 (2005), pp. 1–137. DOI: 10.1017/S0962492904000212.
- [73] M. Benzi and A.J. Wathen. “Some Preconditioning Techniques for Saddle Point Problems”. In: *Model Order Reduction: Theory, Research Aspects and Applications*. Ed. by W.H.A. Schilders, H.A. van der Vorst, and Joost Rommes. Springer, 2008, pp. 195–211.
- [74] D. Bercovici and G. Schubert. *Treatise on geophysics: Mantle dynamics. Vol. 7*. Elsevier, 2007.
- [75] D. Bercovici and G. Schubert. *Treatise on geophysics: Mantle dynamics. Vol. 7 - Second Edition*. Elsevier, 2015.
- [76] David Bercovici. “A simple model of plate generation from mantle flow”. In: *Geophysical Journal International* 114.3 (1993), pp. 635–650. DOI: 10.1111/j.1365-246X.1993.tb06993.x.
- [77] David Bercovici. “A source-sink model of the generation of plate tectonics from non-Newtonian mantle flow”. In: *Journal of Geophysical Research: Solid Earth* 100.B2 (1995), pp. 2013–2030.
- [78] David Bercovici. “Plate generation in a simple model of lithosphere-mantle flow with dynamic self-lubrication”. In: *Earth and Planetary Science Letters* 144.1-2 (1996), pp. 41–51.
- [79] M. Bercovier and M. Engelman. “A finite-element for the numerical solution of viscous incompressible flows”. In: *J. Comp. Phys.* 30 (1979), pp. 181–201. DOI: 10.1016/0021-9991(79)90098-6.
- [80] M. Bercovier and M. Engelman. “A finite-element method for incompressible Non-Newtonian flows”. In: *J. Comp. Phys.* 36 (1980), pp. 313–326. DOI: 10.1016/0021-9991(80)90163-1.

- [81] Arie P van den Berg, Peter E van Keken, and David A Yuen. “The effects of a composite non-Newtonian and Newtonian rheology on mantle convection”. In: *Geophysical Journal International* 115.1 (1993), pp. 62–78. DOI: 10.1111/j.1365-246X.1993.tb05588.x.
- [82] Arie P van den Berg and David A Yuen. “Is the lower-mantle rheology Newtonian today?”. In: *Geophysical research letters* 23.16 (1996), pp. 2033–2036. DOI: 10.1029/96GL02065.
- [83] Christine Bernardi and Genevieve Raugel. “Analysis of Some Finite Elements for the Stokes Problem”. In: *Mathematics of Computation* 44.169 (1985), pp. 71–79.
- [84] FH Bertrand, MR Gadbois, and PA Tanguy. “Tetrahedral elements for fluid flow”. In: *International Journal for Numerical Methods in Engineering* 33.6 (1992), pp. 1251–1267. DOI: 10.1002/nme.1620330610.
- [85] B. Deglo de Besses, A. Magnin, and P. Jay. “Sphere drag in a viscoplastic fluid”. In: *AIChE Journal* 50.10 (2004), pp. 2627–2629. DOI: 10.1002/aic.10252.
- [86] M.J. Beuchert and Y.Y. Podladchikov. “Viscoelastic mantle convection and lithospheric stresses”. In: *GJI* 183 (2010), pp. 35–63. DOI: 10.1111/j.1365-246X.2010.04708.x.
- [87] Marcus J Beuchert, Yuri Y Podladchikov, Nina SC Simon, and Lars H Rüpkke. “Modeling of craton stability using a viscoelastic rheology”. In: *Journal of Geophysical Research: Solid Earth* 115.B11 (2010). DOI: 10.1029/2009JB006482.
- [88] E.C. Bingham. *Fluidity and Plasticity*. McGraw-Hill, New York, 1922.
- [89] Francis Birch. “Elasticity and constitution of the Earth’s interior”. In: *Journal of Geophysical Research* 57.2 (1952), pp. 227–286. DOI: 10.1029/JZ057i002p00227.
- [90] P. Bird and D.A. Yuen. “The use of the minimum-dissipation principle in tectonophysics”. In: *Earth Planet. Sci. Lett.* 45 (1979), pp. 214–217. DOI: 10.1016/0012-821X(79)90122-5.
- [91] Peter Bird. “Finite element modeling of lithosphere deformation: the Zagros collision orogeny”. In: *Tectonophysics* 50.2-3 (1978), pp. 307–336.
- [92] Peter Bird. “Thin-plate and thin-shell finite-element programs for forward dynamic modeling of plate deformation and faulting”. In: *Computers & Geosciences* 25.4 (1999), pp. 383–394.
- [93] D. Bittner and H. Schmeling. “Numerical modelling of melting processes and induced diapirism in the lower crust”. In: *Geophys. J. Int.* 123 (1995), pp. 59–70. DOI: 10.1111/j.1365-246X.1995.tb06661.x.
- [94] J Blackery and E Mitsoulis. “Creeping motion of a sphere in tubes filled with a Bingham plastic material”. In: *Journal of non-newtonian fluid mechanics* 70.1-2 (1997), pp. 59–77.
- [95] B. Blankenbach et al. “A benchmark comparison for mantle convection codes”. In: *Geophys. J. Int.* 98 (1989), pp. 23–38. DOI: 10.1111/j.1365-246X.1989.tb05511.x.
- [96] A Blazquez, B Meyssignac, JM Lemoine, Etienne Berthier, A Ribes, and A Cazenave. “Exploring the uncertainty in GRACE estimates of the mass redistributions at the Earth surface: implications for the global water and sea level budgets”. In: *Geophysical Journal International* 215.1 (2018), pp. 415–430. DOI: 10.1093/gji/ggy293.
- [97] Irina Blinova, Ilya Makeev, and Igor Popov. “Benchmark solutions for stokes flows in cylindrical and spherical geometry”. In: *Bulletin of the Transilvania University of Brasov. Mathematics, Informatics, Physics. Series III* 9.1 (2016), p. 11. DOI: xxxx.
- [98] AM Bobrov and AA Baranov. “Thermochemical Mantle Convection with Drifting Deformable Continents: Main Features of Supercontinent Cycle”. In: *Pure and Applied Geophysics* (2019), pp. 1–21. DOI: 10.1007/s00024-019-02164-w.

- [99] M Bocher, Nicolas Coltice, Alexandre Fournier, and Paul J Tackley. “A sequential data assimilation approach for the joint reconstruction of mantle convection and surface tectonics”. In: *Geophysical Journal International* 204.1 (2016), pp. 200–214.
- [100] Marie Bocher, Alexandre Fournier, and Nicolas Coltice. “Ensemble Kalman filter for the reconstruction of the Earth’s mantle circulation”. In: *Nonlinear Processes in Geophysics* 25.1 (2018), pp. 99–123.
- [101] P. B. Bochev, C. R. Dohrmann, and M. D. Gunzburger. “Stabilization of Low-order Mixed Finite Elements for the Stokes Equations”. In: *SIAM Journal on Numerical Analysis* 44.1 (2006), pp. 82–101. DOI: 10.1137/s0036142905444482.
- [102] Pavel B Bochev and Clark R Dohrmann. “A computational study of stabilized, low-order C 0 finite element approximations of Darcy equations”. In: *Computational Mechanics* 38.4-5 (2006), pp. 323–333. DOI: 10.1007/s00466-006-0036-y.
- [103] Pavel B Bochev and Max D Gunzburger. *Least-squares finite element methods*. Vol. 166. Springer Science & Business Media, 2009.
- [104] Pavel B Bochev, Max D Gunzburger, and John N Shadid. “Stability of the SUPG finite element method for transient advection–diffusion problems”. In: *Computer methods in applied mechanics and engineering* 193.23-26 (2004), pp. 2301–2323. DOI: 10.1016/j.cma.2004.01.026.
- [105] L. Bodri and B. Bodri. “Flow, stress and temperature in island arc areas”. In: *Geophysical & Astrophysical Fluid Dynamics* 13.1 (1979), pp. 95–105. DOI: 10.1080/03091927908243763.
- [106] Reinhard Boehler. “Melting temperature of the Earth’s mantle and core: Earth’s thermal structure”. In: *Annual Review of Earth and Planetary Sciences* 24.1 (1996), pp. 15–40.
- [107] D. Boffi and L. Gastaldi. “On the quadrilateral  $Q_2 - P_1$  element for the Stokes problem”. In: *Int. J. Num. Meth. Fluids* 39 (2002), pp. 1001–1011. DOI: 10.1002/flid.358.
- [108] Daniele Boffi, Franco Brezzi, and Michel Fortin. “Finite elements for the Stokes problem”. In: *Mixed finite elements, compatibility conditions, and applications: Lectures given at the CIME Summer School held in Cetraro, Italy, June 26-July 1, 2006*. Springer, 2008, pp. 45–100. DOI: 10.1007/978-3-540-78319-0\_2.
- [109] Daniele Boffi, Franco Brezzi, and Michel Fortin. *Mixed Finite Element Methods and Applications*. Springer, 2013. ISBN: 978-3-642-36518-8.
- [110] Daniele Boffi, Nicola Cavallini, Francesca Gardini, and Lucia Gastaldi. “Local mass conservation of Stokes finite elements”. In: *Journal of scientific computing* 52.2 (2012), pp. 383–400. DOI: 10.1007/s10915-011-9549-4.
- [111] JM Boland and RA Nicolaides. “On the stability of bilinear-constant velocity-pressure finite elements”. In: *Numerische Mathematik* 44.2 (1984), pp. 219–222. DOI: 10.1007/BF01410106.
- [112] JM Boland and RA Nicolaides. “Stable and semistable low order finite elements for viscous flows”. In: *SIAM journal on numerical analysis* 22.3 (1985), pp. 474–492. DOI: 10.1137/0722028.
- [113] Alessandro Bonaccorso and Paul M Davis. “Models of ground deformation from vertical volcanic conduits with application to eruptions of Mount St. Helens and Mount Etna”. In: *Journal of Geophysical Research: Solid Earth* 104.B5 (1999), pp. 10531–10542.
- [114] M.-A. Bonnardot, R. Hassani, and E. Tric. “Numerical modelling of lithosphere-asthenosphere interaction in a subduction zone”. In: *Earth Planet. Sci. Lett.* 272 (2008), pp. 698–708. DOI: 10.1016/j.epsl.2008.06.009.

- [115] Alain Bonneville and Patrick Capolsini. “THERMIC: a 2-D finite-element tool to solve conductive and advective heat transfer problems in Earth sciences”. In: *Computers & Geosciences* 25.10 (1999), pp. 1137–1148.
- [116] O. Botella and R. Peyret. “Benchmark spectral results on the lid-driven cavity flow”. In: *Computers and Fluids* 27.4 (1998), pp. 421–433.
- [117] MHP Bott, GD Waghorn, and A Whittaker. “Plate boundary forces at subduction zones and trench-arc compression”. In: *Tectonophysics* 170.1-2 (1989), pp. 1–15.
- [118] Lorenzo Botti and Daniele A Di Pietro. “A pressure-correction scheme for convection-dominated incompressible flows with discontinuous velocity and continuous pressure”. In: *Journal of computational physics* 230.3 (2011), pp. 572–585. DOI: 10.1016/j.jcp.2010.10.004.
- [119] Mathieu Bouffard, Stéphane Labrosse, Gaël Choblet, Alexandre Fournier, Julien Aubert, and Paul J Tackley. “A particle-in-cell method for studying double-diffusive convection in the liquid layers of planetary interiors”. In: *Journal of Computational Physics* 346 (2017), pp. 552–571. DOI: 10.1016/j.jcp.2017.06.028.
- [120] Edouard Boujo, Claire Bourquard, Y Xiong, and Nicolas Noiray. “Processing time-series of randomly forced self-oscillators: The example of beer bottle whistling”. In: *Journal of Sound and Vibration* 464 (2020), p. 114981. DOI: 10.1016/j.jsv.2019.114981.
- [121] Johannes Bouman et al. “GOCE gravitational gradients along the orbit”. In: *Journal of Geodesy* 85.11 (2011), p. 791. DOI: 10.1007/s00190-011-0464-0.
- [122] Johannes Bouman et al. “GOCE gravity gradient data for lithospheric modeling”. In: *International Journal of Applied Earth Observation and Geoinformation* 35 (2015), pp. 16–30.
- [123] L. Bourgouin, H.-B. Mühlhaus, A.J. Hale, and A. Arsac. “Studying the influence of a solid shell on lava dome growth and evolution using the level set method”. In: *Geophy. J. Int.* 170 (2007), pp. 1431–1438. DOI: 10.1111/j.1365-246X.2007.03471.x.
- [124] L. Bourgouin, H.-B. Mühlhaus, A.J. Hale, and A. Arsac. “Towards realistic simulations of lava dome growth using the level set method”. In: *Acta Geotechnica* 1 (2006), pp. 225–236. DOI: 10.1007/s11440-006-0016-6.
- [125] Allan F Bower. *Applied mechanics of solids*. CRC press, 2009. ISBN: 978-1-4398-0247-2.
- [126] D.J. Bower, M. Gurnis, and D. Sun. “Dynamic origins of seismic wavespeed variation in D”. In: *Phys. Earth. Planet. Inter.* 214 (2013), pp. 74–86. DOI: 10.1016/j.pepi.2012.10.004.
- [127] Malte Braack and Gert Lube. “Finite elements with local projection stabilization for incompressible flow problems”. In: *Journal of Computational Mathematics* (2009), pp. 116–147. DOI: xxx.
- [128] D. Braess. *Finite Elements: Theory, Fast Solvers, and Applications in Solid Mechanics*. Cambridge, 2007. ISBN: 978-0-521705189.
- [129] D. Braess and R. Sarazin. “An Efficient Smoother for the Stokes Problem”. In: *Applied Numerical Math.* 23 (1997), pp. 3–20.
- [130] Henning Braess and Peter Wriggers. “Arbitrary Lagrangian Eulerian finite element analysis of free surface flow”. In: *Computer Methods in Applied Mechanics and Engineering* 190.1-2 (2000), pp. 95–109. DOI: 10.1016/S0045-7825(99)00416-8.
- [131] Carla Braitenberg. “Exploration of tectonic structures with GOCE in Africa and across-continent”. In: *International Journal of Applied Earth Observation and Geoinformation* 35 (2015), pp. 88–95. DOI: 10.1016/j.jag.2014.01.013.

- [132] James H Bramble and Joseph E Pasciak. “A preconditioning technique for indefinite systems resulting from mixed approximations of elliptic problems”. In: *Mathematics of Computation* 50.181 (1988), pp. 1–17.
- [133] J.P. Brandenburg, E.H. Hauri, P.E. van Keken, and C.J. Ballentine. “A multiple-system study of the geochemical evolution of the mantle with force-balanced plates and thermochemical effects”. In: *Earth Planet. Sci. Lett.* 276 (2008), pp. 1–13. DOI: 10.1016/j.epsl.2008.08.027.
- [134] J.P. Brandenburg and P.E. van Keken. “Methods for thermochemical convection in Earth’s mantle with force-balanced plates”. In: *Geochem. Geophys. Geosyst.* 8.11 (2007).
- [135] J. Braun. “Pecube: a new finite-element code to solve the 3D heat transport equation including the effects of a time-varying, finite amplitude surface topography”. In: *Computers and Geosciences* 29 (2003), pp. 787–794. DOI: 10.1016/S0098-3004(03)00052-9.
- [136] J. Braun, C. Thieulot, P. Fullsack, M. DeKool, and R.S. Huismans. “DOUAR: a new three-dimensional creeping flow model for the solution of geological problems”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 76–91. DOI: 10.1016/j.pepi.2008.05.003.
- [137] J. Braun and P. Yamato. “Structural evolution of a three-dimensional, finite-width crustal wedge”. In: *Tectonophysics* 484 (2010), pp. 181–192. DOI: 10.1016/j.tecto.2009.08.032.
- [138] Jean Braun. “The many surface expressions of mantle dynamics”. In: *Nature Geoscience* 3.12 (2010), p. 825. DOI: 10.1038/ngeo1020.
- [139] Jean Braun. “Three-dimensional numerical modeling of compressional orogenies: Thrust geometry and oblique convergence”. In: *Geology* 21.2 (1993), pp. 153–156. DOI: 10.1130/0091-7613(1993)021<0153:TDNMOC>2.3.CO;2.
- [140] Jean Braun. “Three-dimensional numerical simulations of crustal-scale wrenching using a non-linear failure criterion”. In: *Journal of Structural Geology* 16.8 (1994), pp. 1173–1186. DOI: 10.1016/0191-8141(94)90060-4.
- [141] Jean Braun and Christopher Beaumont. “A physical explanation of the relation between flank uplifts and the breakup unconformity at rifted continental margins”. In: *Geology* 17.8 (1989), pp. 760–764.
- [142] Jean Braun and Christopher Beaumont. “Styles of continental rifting: results from dynamic models of lithospheric extension”. In: *Canadian Society of Petroleum Geologists, Memoir* 12 (1987), pp. 241–258.
- [143] Jean Braun and Christopher Beaumont. “Three-dimensional numerical experiments of strain partitioning at oblique plate boundaries: Implications for contrasting tectonic styles in the southern Coast Ranges, California, and central South Island, New Zealand”. In: *J. Geophys. Res.* 100.B9 (1995), pp. 18, 059–18, 074. DOI: 10.1029/95JB01683.
- [144] Susanne C Brenner. “Korn’s inequalities for piecewise H1 vector fields”. In: *Mathematics of Computation* (2004), pp. 1067–1087. DOI: xxxx. URL: <http://www.jstor.org/stable/4099887>.
- [145] D Breuer, H Zhou, David A Yuen, and T Spohn. “Phase transitions in the Martian mantle: Implications for the planet’s volcanic history”. In: *Journal of Geophysical Research: Planets* 101.E3 (1996), pp. 7531–7542. DOI: 10.1029/96JE00117.
- [146] Doris Breuer, Dave A Yuen, and Tilman Spohn. “Phase transitions in the Martian mantle: Implications for partially layered convection”. In: *Earth and Planetary Science Letters* 148.3-4 (1997), pp. 457–469. DOI: 10.1016/S0012-821X(97)00049-6.

- [147] M Breuer, A Manglik, J Wicht, T Trümper, H Harder, and U Hansen. “Thermochemically driven convection in a rotating spherical shell”. In: *Geophysical Journal International* 183.1 (2010), pp. 150–162. DOI: 10.1111/j.1365-246X.2010.04722.x.
- [148] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, 1991.
- [149] Franco Brezzi. “On the existence, uniqueness and approximation of saddle-point problems arising from Lagrangian multipliers”. In: *Publications mathématiques et informatique de Rennes* S4 (1974), pp. 1–26. DOI: xxxx.
- [150] Franco Brezzi, Marie-Odile Bristeau, Leopoldo P Franca, Michel Mallet, and Gilbert Rogé. “A relationship between stabilized finite element methods and the Galerkin method with bubble functions”. In: *Computer Methods in Applied Mechanics and Engineering* 96.1 (1992), pp. 117–129. DOI: 10.1016/0045-7825(92)90102-P.
- [151] RJ Bridwell and CA Anderson. *Thermomechanical models of the Rio Grande rift*. Tech. rep. Los Alamos Scientific Lab., NM (USA), 1980.
- [152] RJ Bridwell and C Potzick. “Thermal regimes, mantle diapirs and crustal stresses of continental rifts”. In: *Tectonophysics* 73.1-3 (1981), pp. 15–32.
- [153] William L Briggs, Steve F McCormick, et al. *A multigrid tutorial*. Vol. 72. Siam, 2000.
- [154] A.N. Brooks and T.J.R. Hughes. “Streamline Upwind/Petrov-Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier-Stokes equations”. In: *Computer Methods in Applied Mechanics and Engineering* 32 (1982), pp. 199–259. DOI: 10.1016/0045-7825(82)90071-8.
- [155] Geoffrey C Brown and Alan E Mussett. *The Inaccessible Earth: An Integrated View of its Structure and Composition*, 276 pp. Chapman and Hall, London, 1993.
- [156] S. Brune and J. Autin. “The rift to break-up evolution of the Gulf of Aden: Insights from 3D numerical lithospheric-scale modelling”. In: *Tectonophysics* 607 (2013), pp. 65–79. DOI: 10.1016/j.tecto.2013.06.029.
- [157] C.-H. Bruneau and M. Saad. “The 2D lid-driven cavity problem revisited”. In: *Computers & Fluids* 35 (2006), pp. 326–348. DOI: 10.1016/j.compfluid.2004.12.004.
- [158] W Roger Buck. “Small-scale convection induced by passive rifting: the cause for uplift of rift shoulders”. In: *Earth and Planetary Science Letters* 77.3-4 (1986), pp. 362–372. DOI: 10.1016/0012-821X(86)90146-9.
- [159] B.A. Buffet, C.W. Gable, and R.J. R.J. O’Connell. “Linear stability of a layered fluid with mobile surface plates”. In: *J. Geophys. Res.* 99(B10) (1994), pp. 19, 885–19, 900.
- [160] H.H. Bui, R. Fukugawa, K. Sako, and S. Ohno. “Lagrangian meshfree particles method (SPH) for large deformation and failure flows of geomaterial using elastic-plastic soil constitutive model”. In: *Int. J. Numer. Anal. Geomech.* 32.12 (2008), pp. 1537–1570.
- [161] S. Buiter et al. “The numerical sandbox: comparison of model results for a shortening and an extension experiment”. In: *Analogue and Numerical Modelling of Crustal-Scale Processes. Geological Society, London. Special Publications* 253 (2006), pp. 29–64.
- [162] S.J.H. Buiter, R. Govers, and M.J.R. Wortel. “A modelling study of vertical surface displacements at convergent plate margins”. In: *Geophy. J. Int.* 147 (2001), pp. 415–427. DOI: 10.1046/j.1365-246X.2001.00545.x.
- [163] S.J.H. Buiter, R.S. Huismans, and C. Beaumont. “Dissipation analysis as a guide to mode selection during crustal extension and implications for the styles of sedimentary basins”. In: *J. Geophys. Res.* 113.B06406 (2008), B06406. DOI: 10.1029/2007JB005272.

- [164] S.J.H. Buiter et al. “Benchmarking numerical models of brittle thrust wedges”. In: *Journal of Structural Geology* 92 (2016), pp. 140–177. DOI: 10.1016/j.jsg.2016.03.003.
- [165] A.L. Bull, M. Domeier, and T.H. Torsvik. “The effect of plate motion history on the longevity of deep mantle heterogeneities”. In: *Earth Planet. Sci. Lett.* 401 (2014), pp. 172–182. DOI: 10.1016/j.epsl.2014.06.008.
- [166] A.L. Bull, A.K. McNamara, T.W. Becker, and J. Ritsema. “Global scale models of the mantle flow field predicted by synthetic tomography models”. In: *Phys. Earth. Planet. Inter.* 182 (2010), pp. 129–138. DOI: 10.1016/j.pepi.2010.03.004.
- [167] Abigail L Bull, Allen K McNamara, and Jeroen Ritsema. “Synthetic tomography of plume clusters and thermochemical piles”. In: *Earth and Planetary Science Letters* 278.3-4 (2009), pp. 152–162. DOI: 10.1016/j.epsl.2008.11.018.
- [168] Keith Edward Bullen. *The Earth’s density*. Springer Science & Business Media, 1975.
- [169] P.S. Bullen. *Handbook of Means and Their Inequalities*. Springer; 2nd edition, 2003.
- [170] H.-P. Bunge. “Low plume excess temperature and high core heat flux inferred from non-adiabatic geotherms in internally heated mantle circulation models”. In: *Physics of the Earth and Planetary Interiors* 153.1-3 (2005), pp. 3–10. DOI: 10.1016/j.pepi.2005.03.017.
- [171] H.-P. Bunge, Y. Ricard, and J. Matas. “Non-adiabaticity in mantle convection”. In: *Geophysical Research Letters* 28.5 (2001), pp. 879–882. DOI: 10.1029/2000GL011864.
- [172] H.-P. Bunge, M. Richards, C. Lithgow-Bertelloni, J.R. Baumgardner, S.P. Grand, and B. Romanowicz. “Time scales and heterogeneous structure in geodynamic Earth models”. In: *Science* 280 (1998), pp. 91–95.
- [173] H.-P. Bunge, M.A. Richards, and J.R. Baumgardner. “A sensitivity study of three-dimensional spherical mantle convection at  $10^8$  Rayleigh number: Effects of depth-dependent viscosity, heating mode, and endothermic phase change”. In: *J. Geophys. Res.* 102.B6 (1997), pp. 11, 991–12, 007. DOI: 10.1029/96JB03806.
- [174] H.-P. Bunge, M.A. Richards, and J.R. Baumgardner. “Effect of depth-dependent viscosity on the planform of mantle convection”. In: *Nature* 379 (1996), pp. 436–438. DOI: 10.1038/379436a0.
- [175] Martí Burcet, Beñat Oliveira, Juan Carlos Afonso, and Sergio Zlotnik. “A face-centred finite volume approach for coupled transport phenomena and fluid flow”. In: *Applied Mathematical Modelling* 125 (2024), pp. 293–312. DOI: 10.1016/j.apm.2023.08.031.
- [176] Guido Buresti. “A note on Stokes’ hypothesis”. In: *Acta Mechanica* 226 (2015), pp. 3555–3559. DOI: 10.1007/s00707-015-1380-9.
- [177] J-P Burg and Yu Podladchikov. “Lithospheric scale folding: numerical modelling and application to the Himalayan syntaxes”. In: *International Journal of Earth Sciences* 88.2 (1999), pp. 190–200. DOI: 10.1007/s005310050259.
- [178] J.-P. Burg and S.M. Schmalholz. “Viscous heating allows thrusting to overcome crustal-scale buckling: Numerical investigation with application to the Himalayan syntaxes”. In: *Earth Planet. Sci. Lett.* 274 (2008), pp. 189–203. DOI: 10.1016/j.epsl.2008.07.022.
- [179] Roland Bürgmann and Georg Dresen. “Rheology of the lower crust and upper mantle: Evidence from rock mechanics, geodesy, and field observations”. In: *Annu. Rev. Earth Planet. Sci.* 36 (2008), pp. 531–567. DOI: 10.1146/annurev.earth.36.031207.124326.
- [180] Gilmer R Burgos, Andreas N Alexandrou, and Vladimir Entov. “On the determination of yield surfaces in Herschel–Bulkley fluids”. In: *Journal of Rheology* 43.3 (1999), pp. 463–483. DOI: 10.1122/1.550992.

- [181] E. Burman and P. Hansbo. “Edge stabilization for the generalized Stokes problem: A continuous interior penalty method”. In: *Computer Methods in Applied Mechanics and Engineering* 195 (2006), pp. 2393–2410. DOI: 10.1016/j.cma.2005.05.009.
- [182] Erik Burman and Peter Hansbo. “A unified stabilized method for Stokes’ and Darcy’s equations”. In: *Journal of Computational and Applied Mathematics* 198.1 (2007), pp. 35–51. DOI: 10.1016/j.cam.2005.11.022.
- [183] E. Burov and S. Cloetingh. “Erosion and rift dynamics: new thermomechanical aspects of post-rift evolution of extensional basins”. In: *Earth Planet. Sci. Lett.* 150 (1997), pp. 7–26. DOI: 10.1016/S0012-821X(97)00069-1.
- [184] E. Burov and A. Poliakov. “Erosion and rheology controls on synrift and postrift evolution: Verifying old and new ideas using a fully coupled numerical model”. In: *J. Geophys. Res.* 106.B8 (2001), pp. 16, 461–16, 481. DOI: 10.1029/2001JB000433.
- [185] EB Burov and P Molnar. “Small and large-amplitude gravitational instability of an elastically compressible viscoelastic Maxwell solid overlying an inviscid incompressible fluid: dependence of growth rates on wave number and elastic constants at low Deborah numbers”. In: *Earth and Planetary Science Letters* 275.3-4 (2008), pp. 370–381. DOI: 10.1016/j.epsl.2008.08.032.
- [186] EB Burov, AB Watts, et al. “The long-term strength of continental lithosphere: “jelly sandwich” or “crème brûlée”?” In: *GSA today* 16.1 (2006), p. 4. DOI: 10.1130/1052-5173(2006)016<4:TLTSC>2.0.CO;2.
- [187] Evgene B Burov. “Rheology and strength of the lithosphere”. In: *Marine and Petroleum Geology* 28.8 (2011), pp. 1402–1443. DOI: 10.1016/j.marpetgeo.2011.05.008.
- [188] C. Burstedde, O. Ghattas, G. Stadler, T. Tu, and L.C. Wilcox. “Parallel scalable adjoint-based adaptive solution of variable-viscosity Stokes flow problems”. In: *Computer Methods in Applied Mechanics and Engineering* 198 (2009), pp. 1691–1700. DOI: 10.1016/j.cma.2008.12.015.
- [189] C. Burstedde et al. “Large-scale adaptive mantle convection simulation”. In: *Geophys. J. Int.* 192 (2013), pp. 889–906. DOI: 10.1093/gji/ggs070.
- [190] C. Burstedde et al. “Scalable Adaptive Mantle Convection Simulation on Petascale Supercomputers”. In: *ACM/IEEE SC Conference Series, 2008* (2008).
- [191] Carsten Burstedde, Lucas C. Wilcox, and Omar Ghattas. “p4est: Scalable Algorithms for Parallel Adaptive Mesh Refinement on Forests of Octrees”. In: *SIAM Journal on Scientific Computing* 33.3 (2011), pp. 1103–1133. DOI: 10.1137/100791634.
- [192] Carsten Burstedde et al. “Extreme-scale AMR”. In: *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE Computer Society. 2010, pp. 1–12. DOI: 10.1109/SC.2010.25.
- [193] F.H. Busse et al. “3D convection at infinite Prandtl number in Cartesian geometry - a benchmark comparison”. In: *Geophys. Astrophys. Fluid Dynamics* 75.1 (1994), pp. 39–59. DOI: 10.1080/03091929408203646.
- [194] Friedrich H Busse. “Fundamentals of thermal convection”. In: (1989).
- [195] J. Butcher. *Numerical Methods for Ordinary Differential Equations*. New York: John Wiley & Sons, 2003.
- [196] James Byerlee. “Friction of rocks”. In: *Rock friction and earthquake prediction*. 1978, pp. 615–626.



- [197] Mauro Cacace and Antoine B Jacquey. “Flexible parallel implicit modelling of coupled thermal–hydraulic–mechanical processes in fractured rocks”. In: *Solid Earth* 8.5 (2017), pp. 921–941. DOI: 10.5194/se-2017-33.
- [198] O. Cadek and D.A. Yuen. “GLOBAL GEODYNAMICS: Geodynamics workshop in the Czech Republic”. In: *Terra Nova* 5.6 (1993), pp. 573–590. DOI: 10.1111/j.1365-3121.1993.tb00308.x.
- [199] O Čadek and AP van den Berg. “Radial profiles of temperature and viscosity in the Earth’s mantle inferred from the geoid and lateral seismic structure”. In: *Earth and Planetary Science Letters* 164.3-4 (1998), pp. 607–615. DOI: 10.1016/S0012-821X(98)00244-1.
- [200] C Cadio, I Panet, A Davaille, M Diamant, L Métivier, and O de Viron. “Pacific geoid anomalies revisited in light of thermochemical oscillating domes in the lower mantle”. In: *Earth and Planetary Science Letters* 306.1-2 (2011), pp. 123–135. DOI: 10.1016/j.epsl.2011.03.040.
- [201] J. Cahouet and J.-P. Chabard. “Some fast 3D finite element solvers for the generalized Stokes problem”. In: *Int. J. Num. Meth. Fluids* 8 (1988), pp. 869–895.
- [202] Z Cai, Jim Douglas Jr, Juan E Santos, Dongwoo Sheen, and Xiu Ye. “Nonconforming quadrilateral finite elements: a correction”. In: *Calcolo* 37.4 (2000), pp. 253–254. DOI: 10.1007/s100920070004.
- [203] Zhiqiang Cai, Jim Douglas, and Xiu Ye. “A stable nonconforming quadrilateral finite element method for the stationary Stokes and Navier–Stokes equations”. In: *Calcolo* 36.4 (1999), pp. 215–232. DOI: 10.1007/s100920050031.
- [204] Marco Camesasca, Miron Kaufman, and Ica Manas-Zloczower. “Quantifying fluid mixing with the Shannon entropy”. In: *Macromolecular theory and simulations* 15.8 (2006), pp. 595–607. DOI: 10.1002/mats.200600037.
- [205] Lorenzo G Candiotti, Stefan M Schmalholz, and Thibault Duretz. “Impact of upper mantle convection on lithosphere hyperextension and subsequent horizontally forced subduction initiation”. In: *Solid Earth* 11.6 (2020), pp. 2327–2357. DOI: 10.5194/se-11-2327-2020.
- [206] T-S Cao, Pierre Montmitonnet, and P-O Bouchard. “A detailed description of the Gurson–Tvergaard–Needleman model within a mixed velocity–pressure finite element formulation”. In: *International journal for numerical methods in engineering* 96.9 (2013), pp. 561–583. DOI: 10.1002/nme.4571.
- [207] Z.-H. Cao. “Fast Uzawa algorithm for generalized saddle point problems”. In: *Applied Numerical Mathematics* 46 (2003), pp. 157–171.
- [208] G.F. Carey and J.T. Oden. *Finite Elements: Fluid Mechanics. Vol. VI*. Englewood Cliffs, Prentice-Hall, 1986. ISBN: 0-13-317132-9.
- [209] P Carré, AJF Metherell, and TJ Quinn. “Apropos of “The Gravitational Field of a 111 Tetrahedron””. In: *Metrologia* 23.2 (1986), p. 119. DOI: 10.1088/0026-1394/23/2/007.
- [210] P.J. Carreau. “Rheological Equations from Molecular Network Theories”. In: *Transactions of the Society of rheology* 16.1 (1972), pp. 99–127. DOI: 10.1122/1.549276.
- [211] J. Carrero, B. Cockburn, and D. Schötzau. “Hybridized globally divergence-free LDG methods. Part I: The Stokes problem”. In: *Mathematics of Computation* 75.254 (2005), pp. 533–563.
- [212] Carsten Carstensen, Karonline Köhler, Daniel Peterseim, and Mira Schedensack. “Comparison results for the Stokes equations”. In: *Applied Numerical Mathematics* 95 (2015), pp. 118–129. DOI: 10.1016/j.apnum.2013.12.005.

- [213] J.R. Cash and A.H. Karp. “A variable order Runge-Kutta method for initial value problems with rapidly varying right-hand sides”. In: *ACM Transactions on Mathematical Software* 16.3 (1990), pp. 201–222.
- [214] Patrick Cassen and Ray T Reynolds. “Convection in the Moon: Effect of variable viscosity”. In: *Journal of Geophysical Research* 79.20 (1974), pp. 2937–2944. DOI: 10.1029/JB079i020p02937.
- [215] P. Castillo, B. Cockburn, I. Perugia, and D. Schötzau. “Local discontinuous Galerkin methods for elliptic problems”. In: *Commun. Numer. Meth. Engng* 18 (2002), pp. 69–75.
- [216] N.G. Cerpa, R. Hassani, M. Gerbault, and J.-H. Prévost. “A fictitious domain method for lithosphere-asthenosphere interaction: Application to periodic slab folding in the upper mantle”. In: *Geochemistry, Geophysics, Geosystems* 15.5 (2014), pp. 1852–1877. DOI: 10.1002/2014GC005241.
- [217] Jagabanduhu Chakrabarty. *Theory of plasticity*. 2012. ISBN: 0-7506-6638-2.
- [218] TGG Chanut, Safwan Aljbaae, and V Carruba. “Mascon gravitation model using a shaped polyhedral source”. In: *Monthly Notices of the Royal Astronomical Society* 450.4 (2015), pp. 3742–3749. DOI: 10.1093/mnras/stv845.
- [219] D.S. Chapman. “Thermal gradients in the continental crust”. In: *Geological Society, London, Special Publications* 24 (1986), pp. 63–70. DOI: 10.1144/GSL.SP.1986.024.01.07.
- [220] D. Chappelle and K.J. Bathe. “The Inf-Sup test”. In: *Computers & Structures* 47.4/5 (1993), pp. 537–545.
- [221] J.S. Chen, C. Pan, and T.Y.P. Chang. “On the control of pressure oscillation in bilinear-displacement constant-pressure element”. In: *Comput. Methods Appl. Mech. Engrg.* 128 (1995), pp. 137–152.
- [222] Qingshan Chen, Max Gunzburger, and Mauro Perego. “Well-posedness results for a nonlinear Stokes problem arising in glaciology”. In: *SIAM Journal on Mathematical Analysis* 45.5 (2013), pp. 2710–2733. DOI: 10.1137/110848694.
- [223] Z. Chen. “Analysis of mixed methods using conforming and nonconforming finite element methods”. In: *Modélisation mathématique et analyse numérique* 27.1 (1993), pp. 9–34. DOI: xxxx.
- [224] Z. Chen. *Finite Element Methods and Their Applications*. Springer, 2005.
- [225] Z. Chen. “Projection Finite Element methods for semiconductor device equations”. In: *Computers Math. Applic.* 25.8 (1993), pp. 81–88. DOI: 10.1016/0898-1221(93)90173-S.
- [226] Z. Chen and P. Oswald. “Multigrid and multilevel methods for nonconforming  $Q_1$  elements”. In: *Mathematic of Computation* 67.222 (1998), pp. 667–693.
- [227] P. Chenin and C. Beaumont. “Influence of offset weak zones on the development of rift basins: Activation and abandonment during continental extension and breakup”. In: *J. Geophys. Res.* 118 (2013), pp. 1–23. DOI: 10.1002/jgrb.50138.
- [228] P. Chenin, S.M. Schmalholz, G. Manatschal, and G.D. Karner. “Necking of the Lithosphere: A Reappraisal of Basic Concepts With Thermo-Mechanical Numerical Modeling”. In: *Journal of Geophysical Research: Solid Earth* 123.6 (2018), pp. 5279–5299. DOI: 10.1029/2017JB014155.
- [229] Pauline Chenin, Gianreto Manatschal, Alessandro Decarlis, Stefan M Schmalholz, Thibault Duretz, and Marco Beltrando. “Emersion of distal domains in advanced stages of continental rifting explained by asynchronous crust and mantle necking”. In: *Geochemistry, Geophysics, Geosystems* (2019). DOI: 10.1029/2019GC008357.

- [230] Pauline Chenin, Stefan M Schmalholz, Gianreto Manatschal, and Thibault Duretz. “Impact of crust–mantle mechanical coupling on the topographic and thermal evolutions during the necking phase of ‘magma-poor’ and ‘sediment-starved’ rift systems: A numerical modeling study”. In: *Tectonophysics* (2020), p. 228472. DOI: 10.1016/j.tecto.2020.228472.
- [231] M.V. Chertova, T. Geenen, A. van den Berg, and W. Spakman. “Using open sidewalls for modelling self-consistent lithosphere subduction dynamics”. In: *Solid Earth* 3 (2012), pp. 313–326. DOI: 10.5194/se-3-313-2012.
- [232] Jean Chery, Francis Lucazeau, Marc Daignieres, and Jean-Pierre Vilotte. “Large uplift of rift flanks: A genetic link with lithospheric rigidity?”. In: *Earth and Planetary Science Letters* 112.1-4 (1992), pp. 195–211. DOI: 10.1016/0012-821X(92)90016-0.
- [233] Alima Chibani and Nasserline Kechkar. “Minimal Locally Stabilized Q1-Q0 Schemes for the Generalized Stokes Problem”. In: *Journal of the Korean Mathematical Society* 57.5 (2020), pp. 1239–1266. DOI: 10.4134/JKMS.j190628.
- [234] S. Chiu-Webster, E.J. Hinch, and J.R. Lister. “Very viscous horizontal convection”. In: *J. Fluid Mech.* 611 (2008), pp. 395–426.
- [235] G. Choblet. “Modelling thermal convection with large viscosity gradients in one block of the ‘cubed sphere’”. In: *J. Comp. Phys.* 205 (2005), pp. 269–291.
- [236] G. Choblet, O. Čadek, F. Couturier, and C. Dumoulin. “OEDIPUS: a new tool to study the dynamics of planetary interiors”. In: *Geophys. J. Int.* 170 (2007), pp. 9–30.
- [237] E. Choi and K.D. Petersen. “Making Coulomb angle-oriented shear bands in numerical tectonic models”. In: *Tectonophysics* 657 (2015), pp. 94–101. DOI: 10.1016/j.tecto.2015.06.026.
- [238] E. Choi, E. Tan, L.L. Lavier, and V.M. Calo. “DynEarthSol2D: An efficient unstructured finite element method to study long-term tectonic deformation”. In: *J. Geophys. Res.* 118 (2013), pp. 1–16. DOI: 10.1002/jgrb.50148.
- [239] A.J. Chorin. “A numerical method for solving incompressible viscous flow problems”. In: *J. Comp. Phys.* 2 (1967), pp. 12–26. DOI: 10.1006/jcph.1997.5716.
- [240] Nikolas I Christensen and Walter D Mooney. “Seismic velocity structure and composition of the continental crust: A global view”. In: *Journal of Geophysical Research: Solid Earth* 100.B6 (1995), pp. 9761–9788. DOI: /10.1029/95JB00259.
- [241] R.R. Christensen. “An Eulerian technique for thermomechanical modeling of lithospheric extension”. In: *J. Geophys. Res.* 97.B2 (1992), pp. 2015–2036. DOI: 10.1029/91JB02642.
- [242] U Christensen. “Convection in a variable-viscosity fluid: Newtonian versus power-law rheology”. In: *Earth and Planetary Science Letters* 64.1 (1983), pp. 153–162. DOI: 10.1016/0012-821X(83)90060-2.
- [243] U Christensen. “Convection with pressure-and temperature-dependent non-Newtonian rheology”. In: *Geophysical Journal International* 77.2 (1984), pp. 343–384. DOI: 10.1111/j.1365-246X.1984.tb01939.x.
- [244] U Christensen. “Mixing by time-dependent convection”. In: *Earth and planetary science letters* 95.3-4 (1989), pp. 382–394. DOI: 10.1016/0012-821X(89)90112-X.
- [245] U.R. Christensen and D.A. Yuen. “The interaction of a subducting lithospheric slab with a chemical or phase boundary”. In: *J. Geophys. Res.* 89(B6) (1984), pp. 4389–4402. DOI: 10.1029/JB089iB06p04389.

- [246] U.R. Christensen and D.A. Yuen. “Time-dependent convection with non-Newtonian viscosity”. In: *Journal of Geophysical Research* 94.B1 (1989), pp. 814–820. DOI: 10.1029/JB094iB01p00814.
- [247] Ulrich Christensen. “A numerical model of coupled subcontinental and oceanic convection”. In: *Tectonophysics* 95.1-2 (1983), pp. 1–23. DOI: 10.1016/0040-1951(83)90256-1.
- [248] Ulrich Christensen. “Phase boundaries in finite amplitude mantle convection”. In: *Geophysical Journal International* 68.2 (1982), pp. 487–497. DOI: 10.1111/j.1365-246X.1982.tb04911.x.
- [249] Ulrich R Christensen. “The influence of trench migration on slab penetration into the lower mantle”. In: *Earth and Planetary Science Letters* 140.1-4 (1996), pp. 27–39. DOI: 10.1016/0012-821X(96)00023-4.
- [250] Ulrich R Christensen. “Time-dependent convection in elongated Rayleigh-Benard cells”. In: *Geophysical Research Letters* 14.3 (1987), pp. 220–223. DOI: 10.1029/GL014i003p00220.
- [251] Ulrich R Christensen and David A Yuen. “Layered convection induced by phase transitions”. In: *Journal of Geophysical Research: Solid Earth* 90.B12 (1985), pp. 10291–10300. DOI: 10.1029/JB090iB12p10291.
- [252] Edmund Christiansen and Knud D. Andersen. “Computation of collapse states with von Mises type yield condition”. In: *International Journal for Numerical Methods in Engineering* 46 (1999), pp. 1185–1202.
- [253] Edmund Christiansen and Ole S. Pedersen. “Automatic mesh refinement in limit analysis”. In: *International Journal for Numerical Methods in Engineering* 50 (2001), pp. 1331–1346.
- [254] M.A. Christon. “Dealing with pressure: FEM solution strategies for the pressure in the time-dependent Navier-Stokes equations”. In: *Int. J. Num. Meth. Fluids* 38 (2002), pp. 1177–1198.
- [255] M.A. Christon and G.O. Cook. *LS-DYNA’s Incompressible Flow Solver User’s Manual*. Livermore Software Technology Corporation.
- [256] M.A. Christon, P.M. Gresho, and S.B. Sutton. “Computational predictability of time-dependent natural convection flows in enclosures (including a benchmark solution)”. In: *Int. J. Num. Meth. Fluids* 40 (2002), pp. 953–980.
- [257] Mark A Christon. “The new incompressible flow capabilities in LS-DYNA”. In: *6th International LS-DYNA Users Conference*. 2000.
- [258] Andrea Cioncolini and Daniele Boffi. “The MINI mixed finite element for the Stokes problem: An experimental investigation”. In: *Computers & Mathematics with Applications* 77.9 (2019), pp. 2432–2446. DOI: 10.1016/j.camwa.2018.12.028.
- [259] H. Čížková, A.P. van den Berg, W. Spakman, and Ctirad Matyska. “The viscosity of the earth’s lower mantle inferred from sinking speed of subducted lithosphere”. In: *Phys. Earth. Planet. Inter.* 200–201 (2012), pp. 56–62. DOI: 10.1016/j.pepi.2012.02.010.
- [260] T.C. Clevenger, T. Heister, G. Kanschat, and M. Kronbichler. “A flexible, parallel, adaptive geometric multigrid method for FEM”. In: *ACM Trans. Math. Softw.* 47.1 (2020). DOI: 10.1145/3425193.
- [261] Thomas C Clevenger and Timo Heister. “Comparison Between Algebraic and Matrix-free Geometric Multigrid for a Stokes Problem on Adaptive Meshes with Variable Viscosity”. In: *Numer. Linear Algebra Appl.* (2021). DOI: 10.1002/nla.2375.

- [262] S Cochard and C Ancey. “Experimental investigation of the spreading of viscoplastic fluids on inclined planes”. In: *Journal of Non-Newtonian Fluid Mechanics* 158.1-3 (2009), pp. 73–84.
- [263] B. Cockburn and J. Gopalakrishnan. “The derivation of hybridizable discontinuous Galerkin methods for Stokes flow”. In: *SIAM J. Numer. Anal.* 47.2 (2009), pp. 1092–1125.
- [264] B. Cockburn, G. Kanschat, and D. Schötzau. “The local discontinuous Galerkin method for linearized incompressible fluid flow: a review”. In: *Computers and Fluids* 34 (2005), pp. 491–506.
- [265] B. Cockburn, G. Kanschat, D. Schoetzau, and C. Schwab. “Local discontinuous Galerkin methods for the Stokes system”. In: *SIAM J. Numer. Anal.* 40.1 (2002), pp. 319–343.
- [266] B. Cockburn, G.E. Karniadakis, and C.W.Shu. “The Development of Discontinuous Galerkin Methods”. In: *Discontinuous Galerkin Methods. Lecture Notes in Computational Science and Engineering* 11 (2000).
- [267] B. Cockburn, G.E. Karniadakis, and C.-W. Shu. *Discontinuous Galerkin Methods. Theory, Computation and Applications*. Springer, 2000.
- [268] B. Cockburn, N.C. Nguyen, and J. Peraire. “A Comparison of HDG Methods for Stokes Flow”. In: *J. Sci. Comput.* 45 (2010), pp. 215–237.
- [269] Bernardo Cockburn, Guido Kanschat, and Dominik Schötzau. “The local discontinuous Galerkin method for the Oseen equations”. In: *Mathematics of Computation* 73.246 (2004), pp. 569–593.
- [270] Ramon Codina. “On stabilized finite element methods for linear systems of convection–diffusion–reaction equations”. In: *Computer Methods in Applied Mechanics and Engineering* 188.1-3 (2000), pp. 61–82. DOI: 10.1016/S0045-7825(00)00177-8.
- [271] L. Colli, H.-P. Bunge, and B.S.A. Schuberth. “On retrodictions of global mantle flow with assimilated surface velocities”. In: *Geophysical Research Letters* 42.20 (2015), pp. 8341–8348. DOI: 10.1002/2015GL066001.
- [272] Nicolas Coltice, Laurent Husson, Claudio Faccenna, and Maëlis Arnould. “What drives tectonic plates?” In: *Science Advances* 5.10 (2019). DOI: 10.1126/sciadv.aax4295.
- [273] Nicolas Coltice and J Schmalzl. “Mixing times in the mantle of the early Earth derived from 2-D and 3-D numerical simulations of convection”. In: *Geophysical Research Letters* 33.23 (2006). DOI: 10.1029/2006GL027707.
- [274] Nicolas Coltice and Grace E Shephard. “Tectonic predictions with mantle convection models”. In: *Geophysical Journal International* 213.1 (2018), pp. 16–29. DOI: 10.1093/gji/ggx531.
- [275] G Comini, Marco Manzan, and C Nonino. “Finite element solution of the streamfunction–vorticity equations for incompressible two-dimensional flows”. In: *International journal for numerical methods in fluids* 19.6 (1994), pp. 513–525. DOI: 10.1002/flid.1650190605.
- [276] James AD Connolly and Yury Y Podladchikov. “An analytical solution for solitary porosity waves: dynamic permeability and fluidization of nonlinear viscous and viscoplastic rock”. In: *Geofluids* 15.1-2 (2015), pp. 269–292. DOI: 10.1111/gf1.12110.
- [277] Clinton P Conrad and Peter Molnar. “The growth of Rayleigh-Taylor-type instabilities in the lithosphere for various rheological and density structures”. In: *Geophysical Journal International* 129.1 (1997), pp. 95–112.

- [278] CM Cooper, A Lenardic, A Levander, L Moresi, and K Benn. “Creation and preservation of cratonic lithosphere: seismic constraints and geodynamic models”. In: *GEOPHYSICAL MONOGRAPH-AMERICAN GEOPHYSICAL UNION* 164 (2006), p. 75. DOI: **xxxx**.
- [279] C. Couder-Castaneda, J.C. Ortiz-Aleman, M.G. Orozco-del-Castillo, and M. Nava-Flores. “Forward modeling of gravitational fields on hybrid multi-threaded cluster”. In: *Geofisica Internacional* 54.1 (2015), pp. 31–48.
- [280] D Coumou, T Driesner, and Christoph A Heinrich. “The structure and dynamics of mid-ocean ridge hydrothermal systems”. In: *Science* 321.5897 (2008), pp. 1825–1828.
- [281] P. Coussot. “Yield stress fluid flows: A review of experimental data”. In: *Journal of Non-Newtonian Fluid Mechanics* 211 (2014), pp. 31–49.
- [282] Claire Harvey Craig and Dan McKenzie. “The existence of a thin low-viscosity layer beneath the lithosphere”. In: *Earth and Planetary Science Letters* 78.4 (1986), pp. 420–426.
- [283] F. Crameri and P.J. Tackley. “Spontaneous development of arcuate single-sided subduction in global 3-D mantle convection models with a free surface”. In: *J. Geophys. Res.* 119 (2014). DOI: 10.1002/2014JB010939.
- [284] F. Crameri and P.J. Tackley. “Subduction initiation from a stagnant lid and global overturn: new insights from numerical models with a free surface”. In: *Progress in Earth and Planetary Science* 3 (2016).
- [285] F. Crameri et al. “A comparison of numerical surface topography calculations in geodynamic modelling: an evaluation of the ‘sticky air’ method”. In: *Geophy. J. Int.* 189 (2012), pp. 38–54.
- [286] Fabio Crameri, Grace E. Shephard, and Philip J. Heron. “The misuse of colour in science communication”. In: *Nature Communications* 11 (2020), p. 5444. DOI: 10.1038/s41467-020-19160-7.
- [287] Marc Crombaghs, Erik de Min, and Govert Strang van Hees. “The first absolute gravity measurements in The Netherlands”. In: *TU Delft, Delft, Tech. Rep* (2002).
- [288] M.M. Cross. “Rheology of non-Newtonian fluids: a new flow equation for pseudoplastic systems”. In: *Journal of Colloid Science* 20 (1965), pp. 417–437.
- [289] M. Crouzeix and R.S. Falk. “Nonconforming finite elements for the Stokes problem”. In: *Mathematic of Computation* 52.186 (1989), pp. 437–456. DOI: 10.1090/S0025-5718-1989-0958870-8.
- [290] M. Crouzeix and P.-A. Raviart. “Conforming and nonconforming finite element methods for solving the stationary Stokes equations I”. In: *R.A.I.R.O.* 7.3 (1973), pp. 33–75. DOI: ?. URL: <http://eudml.org/doc/193250>.
- [291] AR Cruden. “Deformation around a rising diapir modeled by creeping flow past a sphere”. In: *Tectonics* 7.5 (1988), pp. 1091–1101.
- [292] L. Cserepes and D.A. Yuen. “Dynamical consequences of mid-mantle viscosity stratification on mantle flows with an endothermic phase transition”. In: *Geophysical Research Letters* 24.2 (1997), pp. 181–184. DOI: 10.1029/96GL03917.
- [293] Marco Cuffaro, Edie Miglio, Mattia Penati, and Marco Viganò. “Mantle thermal structure at northern Mid-Atlantic Ridge from improved numerical methods and boundary conditions”. In: *Geophysical Journal International* 220.2 (2020), pp. 1128–1148. DOI: 10.1093/gji/ggz488.

- [294] Gilda Currenti and Charles A Williams. “Numerical modeling of deformation and stress fields around a magma chamber: Constraints on failure conditions and rheology”. In: *Physics of the Earth and Planetary Interiors* 226 (2014), pp. 14–27. DOI: 10.1016/j.pepi.2013.11.003.
- [295] C.A. Currie and C. Beaumont. “Are diamond-bearing Cretaceous kimberlites related to low-angle subduction beneath western North America”. In: *Earth Planet. Sci. Lett.* 303 (2011), pp. 59–70. DOI: 10.1016/j.epsl.2010.12.036.
- [296] Claire A Currie and Roy D Hyndman. “The thermal structure of subduction zone back arcs”. In: *Journal of Geophysical Research: Solid Earth* 111.B8 (2006). DOI: 10.1029/2005JB004024.
- [297] Elizabeth Cuthill and James McKee. “Reducing the bandwidth of sparse symmetric matrices”. In: *Proceedings of the 1969 24th national conference*. 1969, pp. 157–172. DOI: 10.1145/800195.805928.
- [298] C. Cuvelier, A. Segal, and A.A. van Steenhoven. *Finite Element Methods and Navier-Stokes Equations*. D. Reidel Publishing Company, 1986.
- [299] M. Dabrowski, M. Krotkiewski, and D.W. Schmid. “MILAMIN: Matlab based finite element solver for large problems”. In: *Geochem. Geophys. Geosyst.* 9.4 (2008), Q04030. DOI: 10.1029/2007GC001719.
- [300] S.F. Daly and A. Raefsky. “On the penetration of a hot diapir through a strongly temperature-dependent viscosity medium”. In: *Geophys. J. R. astr. Soc* 83 (1985), pp. 657–681. DOI: 10.1111/j.1365-246X.1985.tb04331.x.
- [301] J. Dannberg, Z. Eilon, U. Faul, R. Gassmüller, P. Moulik, and R. Myhill. “The importance of grain size to mantle dynamics and seismological observations”. In: *Geochem. Geophys. Geosyst.* 18 (2017), pp. 3034–3061. DOI: 10.1002/2017GC006944.
- [302] J. Dannberg and T. Heister. “Compressible magma/mantle dynamics: 3-D, adaptive simulations in ASPECT”. In: *Geophys. J. Int.* 207 (2016), pp. 1343–1366. DOI: 10.1093/gji/ggw329.
- [303] Juliane Dannberg, Rene Gassmüller, Ryan Grove, and Timo Heister. “A new formulation for coupled magma/mantle dynamics”. In: *Geophysical Journal International* 219.1 (2019), pp. 94–107. DOI: 10.1093/gji/ggz190.
- [304] KD Danov, TD Gurkov, H Raszillier, and F Durst. “Stokes flow caused by the motion of a rigid sphere close to a viscous interface”. In: *Chemical Engineering Science* 53.19 (1998), pp. 3413–3434.
- [305] Véronique Dansereau, Jérôme Weiss, Pierre Saramito, and Philippe Lattes. “A Maxwell-Elasto-Brittle rheology for sea ice modelling”. In: *The Cryosphere* 10 (2016), pp. 1339–1359. DOI: 10.5194/tc-10-1339-2016.
- [306] D Rhodri Davies, Stephan C Kramer, Sia Ghelichkhan, and Angus Gibson. “Towards automatic finite-element methods for geodynamics via Firedrake”. In: *Geoscientific Model Development* 15.13 (2022), pp. 5127–5166. DOI: 10.5194/gmd-15-5127-2022.
- [307] D. R. Davies, J. H. Davies, O. Hassan, K. Morgan, and P. Nithiarasu. “Investigations into the applicability of adaptive finite element methods to two-dimensional infinite Prandtl number thermal and thermochemical convection”. In: *Geochemistry, Geophysics, Geosystems* 8.5 (2007). DOI: 10.1029/2006GC001470.
- [308] D.R. Davies, J.H. Davies, P.C. Bollada, O. Hassan, K. Morgan, and P. Nithiarasu. “A hierarchical mesh refinement technique for global 3-D spherical mantle convection modelling”. In: *Geosci. Model Dev.* 6 (2013), pp. 1095–1107. DOI: 10.5194/gmd-6-1095-2013.

- [309] D.R. Davies, C.R. Wilson, and S.C. Kramer. “Fluidity: A fully unstructured anisotropic adaptive mesh computational modeling framework for geodynamics”. In: *Geochem. Geophys. Geosyst.* 12.6 (2011). DOI: 10.1029/2011GC003551.
- [310] David Rhodri Davies, John Huw Davies, O Hassan, K Morgan, and P Nithiarasu. “Adaptive finite element methods in geodynamics: Convection dominated mid-ocean ridge and subduction zone simulations”. In: *International Journal of Numerical Methods for Heat & Fluid Flow* 18.7/8 (2008). DOI: 10.1108/09615530810899079.
- [311] G.F. Davies. “Mantle convection model with a dynamic plate: topography, heat flow and gravity anomalies”. In: *Geophysical Journal International* 98.3 (1989), pp. 461–464. DOI: 10.1111/j.1365-246X.1989.tb02283.x.
- [312] G.F. Davies. “Role of the lithosphere in mantle convection”. In: *Journal of Geophysical Research: Solid Earth* 93.B9 (1988), pp. 10451–10466. DOI: 10.1029/JB093iB09p10451.
- [313] Geoffrey F Davies. “Mantle convection under simulated plates: effects of heating modes and ridge and trench migration, and implications for the core-mantle boundary, bathymetry, the geoid and Benioff zones”. In: *Geophys. J. R. astr. Soc* 84.1 (1986), pp. 153–183.
- [314] Geoffrey F Davies. “Viscous mantle flow under moving lithospheric plates and under subduction zones”. In: *Geophysical Journal International* 49.3 (1977), pp. 557–563. DOI: 10.1111/j.1365-246X.1977.tb01303.x.
- [315] Geoffrey F Davies. “Whole-mantle convection and plate tectonics”. In: *Geophysical Journal International* 49.2 (1977), pp. 459–486. DOI: 10.1111/j.1365-246X.1977.tb03717.x.
- [316] Geoffrey F Davies and Michael Gurnis. “Interaction of mantle dregs with convection: Lateral heterogeneity at the core-mantle boundary”. In: *Geophysical Research Letters* 13.13 (1986), pp. 1517–1520.
- [317] J. H. Davies and D. J. Stevenson. “Physical model of source region of subduction zone volcanics”. In: *Journal of Geophysical Research: Solid Earth* 97.B2 (1992), pp. 2037–2070. DOI: 10.1029/91JB02571.
- [318] R.O. Davis and A.P.S. Selvadurai. *Elasticity and Geomechanics*. Cambridge University Press, 1996.
- [319] R.O. Davis and A.P.S. Selvadurai. *Plasticity and Geomechanics*. Cambridge University Press, 2002. ISBN: 0-521-01809-9.
- [320] P. Davy and P. Cobbold. “Indentation tectonics in nature and experiment. 1. Experiments scaled for gravity”. In: *Bulletin of the Geological Institutions of Uppsala* 14 (1988), pp. 129–141.
- [321] Arne De Coninck et al. “Needles: Toward Large-Scale Genomic Prediction with Marker-by-Environment Interaction”. In: *Genetics* 203.1 (2016), pp. 543–555. DOI: 10.1534/genetics.115.179887.
- [322] Albert de Montserrat, Jason P Morgan, and Jörg Hasenclever. “LaCoDe: a Lagrangian two-dimensional thermo-mechanical code for large-strain compressible visco-elastic geodynamical modeling”. In: *Tectonophysics* 767 (2019), p. 228173. DOI: 10.1016/j.tecto.2019.228173.
- [323] Jeroen H De Smet, Arie P van Den Berg, Nico J Vlaar, and David A Yuen. “A characteristics-based method for solving the transport equation and its application to the process of mantle differentiation and continental root growth”. In: *Geophysical Journal International* 140.3 (2000), pp. 651–659. DOI: 10.1046/j.1365-246X.2000.00993.x.
- [324] A. Deb, J.H. Prevost, and B. Loret. “Adaptive meshing for dynamic strain localisation”. In: *Comput. Methods Appl. Mech. Engrg.* 137 (1996), pp. 285–306. DOI: 10.1016/S0045-7825(96)01068-7.



- [325] Herwig Dejonghe. “A completely analytical family of anisotropic Plummer models”. In: *Monthly Notices of the Royal Astronomical Society* 224.1 (1987), pp. 13–39.
- [326] W. DeLandro-Clarke and G.T. Jarvis. “Numerical models of mantle convection with secular cooling”. In: *Geophysical Journal International* 129.1 (1997), pp. 183–193.
- [327] Fábio CG DeMarco, Cláudia R DeAndrade, and Edson L Zaparoli. “The no-slip boundary condition in the stream function-vorticity formulation using the penalty method”. In: *International communications in heat and mass transfer* 30.4 (2003), pp. 495–504. DOI: 10.1016/S0735-1933(03)00078-2.
- [328] SCR Dennis and JD Hudson. “Methods of solution of the velocity-vorticity formulation of the Navier-Stokes equations”. In: *Journal of Computational Physics* 122.2 (1995), pp. 300–306.
- [329] S Dequand et al. “Simplified models of flue instruments: Influence of mouth geometry on the sound source”. In: *The Journal of the Acoustical Society of America* 113.3 (2003), pp. 1724–1735. DOI: 10.1121/1.1543929.
- [330] Y. Deubelbeiss and B.J.P. Kaus. “Comparison of Eulerian and Lagrangian numerical techniques for the Stokes equations in the presence of strongly varying viscosity”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 92–111. DOI: 10.1016/j.pepi.2008.06.023.
- [331] Byron DeVries, Joe Iannelli, Christian Trefftz, Kurt A O’Hearn, and Greg Wolffe. “Parallel implementations of FGMRES for solving large, sparse non-symmetric linear systems”. In: *Procedia Computer Science* 18 (2013), pp. 491–500.
- [332] G. Dhatt and G. Hubert. “A study of penalty elements for incompressible laminar flows”. In: *Int. J. Num. Meth. Fluids* 6 (1986), pp. 1–19.
- [333] E Di Giuseppe et al. “Characterization of Carbopol® hydrogel rheology for experimental tectonics and geodynamics”. In: *Tectonophysics* 642 (2015), pp. 29–45.
- [334] I. Dione, C. Tibirna, and J. Urquiza. “Stokes equations with penalised slip boundary conditions”. In: *International Journal of Computational Fluid Dynamics* 27.6-7 (2013), pp. 283–296. DOI: 10.1080/10618562.2013.821114.
- [335] Nathaniel A Dixon and William B Durham. “Measurement of activation volume for creep of dry olivine at upper-mantle conditions”. In: *Journal of Geophysical Research: Solid Earth* 123.10 (2018), pp. 8459–8473. DOI: 10.1029/2018JB015853.
- [336] C.R. Dohrmann and P.B. Bochev. “A stabilized finite element method for the Stokes problem based on polynomial pressure projections”. In: *Int. J. Num. Meth. Fluids* 46 (2004), pp. 183–201. DOI: 10.1002/flid.752.
- [337] V. Dolejsi and M. Feistauer. *Discontinuous Galerkin Methods. Analysis and Applications to Compressible Flow*. Springer, 2015.
- [338] Vít Dolejší. “Analysis and application of the IIPG method to quasilinear nonstationary convection-diffusion problems”. In: *Journal of Computational and Applied Mathematics* 222.2 (2008), pp. 251–273.
- [339] J Donea, Antonio Huerta, J-Ph Ponthot, and A Rodriguez-Ferran. “Arbitrary Lagrangian-Eulerian Methods, volume 1 of Encyclopedia of Computational Mechanics, chapter 14”. In: *John Wiley & Sons Ltd* 3 (2004), pp. 1–25.
- [340] J. Donea, S. Giuliani, K. Morgan, and L. Quartapelle. “The significance of chequerboarding in a Galerkin finite element solution of the Navier-Stokes equations”. In: *International Journal for Numerical Methods in Engineering* 17.5 (1981), pp. 790–795. DOI: 10.1002/nme.1620170511.

- [341] Jean Donea and Antonio Huerta. *Finite Element Methods for Flow Problems*. John Wiley & Sons, 2003. ISBN: 978-0-471-49666-3.
- [342] J.R. Dormand and P.J. Prince. “A family of embedded Runge-Kutta formulae”. In: *Journal of Computational and Applied Mathematics* 6.1 (1980), pp. 19–26.
- [343] J.R. Dormand and P.J. Prince. “A reconsideration of some embedded Runge-Kutta formulae”. In: *Journal of Computational and Applied Mathematics* 15 (1986), pp. 203–211.
- [344] Jim Douglas and Junping Wang. “An absolutely stabilised finite element method for the Stokes problem”. In: *Mathematics of Computation* 52.186 (1989), pp. 495–508. DOI: 10.1090/S0025-5718-1989-0958871-X.
- [345] Jim Douglas Jr, Juan E Santos, Dongwoo Sheen, and Xiu Ye. “Nonconforming Galerkin methods based on quadrilateral elements for second order elliptic problems”. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 33.4 (1999), pp. 747–770. DOI: 10.1051/m2an:1999161.
- [346] M van Driel, Lion Krischer, Simon C Stähler, Kambod Hosseini, and Tarje Nissen-Meyer. “Instaseis: Instant global seismograms based on a broadband waveform database”. In: *Solid Earth* 6.2 (2015), pp. 701–717. DOI: 10.5194/se-6-701-2015.
- [347] D.C. Drucker and W. Prager. “Soil mechanics and plastic analysis or limit design”. In: *Quarterly of Applied Mathematics* 10.2 (1952), pp. 157–165. DOI: xxxx.
- [348] CP Dubey and VM Tiwari. “Computation of the gravity field and its gradient: Some applications”. In: *Computers & geosciences* 88 (2016), pp. 83–96. DOI: 10.1016/j.cageo.2015.12.007.
- [349] J.K. Dukowicz. “Efficient Volume Computation for Three- Dimensional Hexahedral Cells”. In: *J. Comp. Phys.* 74 (1988), pp. 493–496. DOI: 10.1016/0021-9991(88)90091-5.
- [350] D.A. Dunavant. “High-degree efficient symmetrical Gaussian quadrature rules for the triangle”. In: *Int. J. Num. Meth. Eng.* 21 (1985), pp. 1129–1148. DOI: 10.1002/nme.1620210612.
- [351] J.A. Dunbar and D.S. Sawyer. “Three-dimensional dynamical model of continental rift propagation and margin plateau formation”. In: *J. Geophys. Res.* 101.B12 (1996), pp. 27, 845–27, 863. DOI: 10.1029/96JB01231.
- [352] T. Duretz, D.A. May, T.V. Gerya, and P.J. Tackley. “Discretization errors and free surface stabilisation in the finite difference and marker-in-cell method for applied geodynamics: A numerical study”. In: *Geochem. Geophys. Geosyst.* 12.Q07004 (2011). DOI: 10.1029/2011GC003567.
- [353] Thibault Duretz, Dave A May, and Philippe Yamato. “A free surface capturing discretization for the staggered grid finite difference scheme”. In: *Geophysical Journal International* 204.3 (2016), pp. 1518–1530. DOI: 10.1093/gji/ggv526.
- [354] A. Düster and E. Rank. “The p-version of the finite element method compared to an adaptive h-version for the deformation theory of plasticity”. In: *Computer Methods in Applied Mechanics and Engineering* 190 (2011), pp. 1925–1935. DOI: 10.1016/S0045-7825(00)00215-2.
- [355] Urmi Dutta, Shamik Sarkar, and Nibir Mandal. “Ballooning versus curling of mantle plumes: views from numerical models”. In: *Current Science* 104.7 (2013), pp. 893–903. DOI: xxxx.
- [356] S. Dyksterhuis, P. Rey, R.D. Mueller, and L. Moresi. “Effects of initial weakness on rift architecture”. In: *Geological Society, London, Special Publications* 282 (2007), pp. 443–455. DOI: 10.1144/SP282.18.

- [357] A.M. Dziewonski and D.L. Anderson. “Preliminary reference Earth model”. In: *Phys. Earth. Planet. Inter.* 25 (1981), pp. 297–356. DOI: 10.1016/0031-9201(81)90046-7.
- [358] J Ebbing, J Bouman, M Fuchs, S Gradmann, and R Haagmans. “Sensitivity of GOCE gravity gradients to crustal thickness and density variations: Case study for the Northeast Atlantic Region”. In: *Gravity, Geoid and Height Systems*. Springer, 2014, pp. 291–298. DOI: 10.1007/978-3-319-10837-7\_37.
- [359] Jörg Ebbing et al. “Advancements in satellite gravity gradient data for crustal studies”. In: *The Leading Edge* 32.8 (2013), pp. 900–906. DOI: xxxx.
- [360] David L. Egholm. “A new strategy for discrete element numerical models: 1. Theory”. In: *J. Geophys. Res.* 112 (2007), B05203. DOI: 10.1029/2006JB004557.
- [361] David L. Egholm, Mike Sandiford, Ole R. Clausen, and Søren B. Nielsen. “A new strategy for discrete element numerical models: 2. Sandbox applications”. In: *J. Geophys. Res.* 112 (2007), B05204. DOI: 10.1029/2006JB004558.
- [362] Yuzuru Eguchi. “A new positive-definite regularization of incompressible Navier–Stokes equations discretized with Q1/P0 finite element”. In: *International journal for numerical methods in fluids* 41.8 (2003), pp. 881–904. DOI: 10.1002/d.482.
- [363] Louis W Ehrlich and Murli M Gupta. “Some difference schemes for the biharmonic equation”. In: *SIAM Journal on Numerical Analysis* 12.5 (1975), pp. 773–790. DOI: 10.1137/0712058.
- [364] V. Eijkhout. *Introduction to High Performance Scientific Computing*. Creative Commons, 2013.
- [365] Stanley C Eisenstat, Howard C Elman, and Martin H Schultz. “Variational iterative methods for nonsymmetric systems of linear equations”. In: *SIAM Journal on Numerical Analysis* 20.2 (1983), pp. 345–357.
- [366] Dirk Elbeshhausen and Jay Melosh. “A nonlinear and time-dependent visco-elasto-plastic rheology model for studying shock-physics phenomena”. In: *arXiv preprint arXiv:1805.06453* (2018).
- [367] Y. Elesin, T. Gerya, I.M. Artemieva, and H. Thybo. “Samovar: a thermomechanical code for modeling of geodynamic processes in the lithosphere – application to basin evolution”. In: *Arabian Journal of Geosciences* 3 (2010), pp. 477–497. DOI: 10.1007/s12517-010-0215-1.
- [368] S. Ellis, P. Fullsack, and C. Beaumont. “Oblique convergence of the crust driven by basal forcing: implications for length-scales of deformation and strain partitioning in orogens”. In: *Geophys. J. Int.* 120 (1995), pp. 24–44. DOI: 10.1111/j.1365-246X.1995.tb05909.x.
- [369] Susan Ellis, Guido Schreurs, and Marion Panien. “Comparisons between analogue and numerical models of thrust wedge development”. In: *Journal of Structural Geology* 26.9 (2004), pp. 1659–1675. DOI: 10.1016/j.jsg.2004.02.012.
- [370] Kirk Ellsworth, Gerald Schubert, and Charles G Sammis. “Viscosity profile of the lower mantle”. In: *Geophysical Journal International* 83.1 (1985), pp. 199–213. DOI: 10.1111/j.1365-246X.1985.tb05163.x.
- [371] H. Elman, D. Silvester, and A. Wathen. *Finite Elements and Fast Iterative Solvers*. Oxford Science Publications, 2014. ISBN: 978-0-19-967879-2.
- [372] M.S. Engelman, R.L. Sani, and P.M. Gresho. “The implementation of normal and/or tangential boundary conditions in finite element codes for incompressible fluid flow”. In: *Int. J. Num. Meth. Fluids* 2 (1982), pp. 225–238. DOI: 10.1002/flid.1650020302.

- [373] MS Engelman, Gilbert Strang, and K-J Bathe. “The application of quasi-Newton methods in fluid mechanics”. In: *International Journal for Numerical Methods in Engineering* 17.5 (1981), pp. 707–718.
- [374] P. England. “Some numerical investigations of large scale continental deformation”. In: *Mountain Building Processes*. Academic Press, 1982, pp. 129–189. DOI: xxxx.
- [375] P. England and G. Houseman. “A dynamical model of lithosphere extension and sedimentary basin formation”. In: *J. Geophys. Res.* 91.B3 (1986), pp. 3664–3676. DOI: 10.1029/JB091iB01p00719.
- [376] Philip England and Dan McKenzie. “A thin viscous sheet model for continental deformation”. In: *Geophysical Journal International* 70.2 (1982), pp. 295–321. DOI: 10.1111/j.1365-246X.1982.tb04969.x.
- [377] Philip England and Peter Molnar. “Active deformation of Asia: From kinematics to dynamics”. In: *Science* 278.5338 (1997), pp. 647–650. DOI: 10.1126/science.278.5338.647.
- [378] Philip England and Catherine Wilkins. “A simple analytical approximation to the temperature structure in subduction zones”. In: *Geophysical Journal International* 159.3 (2004), pp. 1138–1154. DOI: 10.1111/j.1365-246X.2004.02419.x.
- [379] I. Ergatoudis, B.M. Irons, and O.C. Zienkiewicz. “Curved, isoparametric, “quadrilateral” elements for finite element analysis”. In: *International journal of solids and structures* 4.1 (1968), pp. 31–42. DOI: 10.1016/0020-7683(68)90031-0.
- [380] VV Ermakov and Nikolai Nikolaevich Kalitkin. “The optimal step and regularization for Newton’s method”. In: *USSR Computational Mathematics and Mathematical Physics* 21.2 (1981), pp. 235–242. DOI: 10.1016/0041-5553(81)90022-7.
- [381] E. Erturk. “Discussions on Driven Cavity Flow”. In: *Int. J. Num. Meth. Fluids* 60 (2009), pp. 275–294.
- [382] Pep Espanol and Cedric Thieulot. “Microscopic derivation of hydrodynamic equations for phase-separating fluid mixtures”. In: *The Journal of chemical physics* 118.20 (2003), pp. 9109–9127. DOI: 10.1063/1.1568333.
- [383] Brian Evans and Christopher Goetze. “The temperature variation of hardness of olivine and its implication for polycrystalline yield stress”. In: *Journal of Geophysical Research: Solid Earth* 84.B10 (1979), pp. 5505–5524. DOI: 10.1029/JB084iB10p05505.
- [384] M. Faccenda, T.V. Gerya, N.S. Mancktelow, and L. Moresi. “Fluid flow during slab unbending and dehydration: Implications for intermediate-depth seismicity, slab weakening and deep water recycling”. In: *Geochem. Geophys. Geosyst.* 13.1 (2012). DOI: 10.1029/2011GC003860.
- [385] Cinzia G Farnetani, Bernard Legras, and Paul J Tackley. “Mixing and deformations in mantle plumes”. In: *Earth and Planetary Science Letters* 196.1-2 (2002), pp. 1–15. DOI: 10.1016/S0012-821X(01)00597-0.
- [386] Cinzia G Farnetani and Henri Samuel. “Lagrangian structures and stirring in the Earth’s mantle”. In: *Earth and Planetary Science Letters* 206.3-4 (2003), pp. 335–348. DOI: 10.1016/S0012-821X(02)01085-3.
- [387] R.J. Farrington, L.-N. Moresi, and F.A. Capitanio. “The role of viscoelasticity in subducting plates”. In: *Geochem. Geophys. Geosyst.* 15 (2014), pp. 4291–4304. DOI: 10.1002/2014GC005507.
- [388] U.H. Faul, J.D. Fitz Gerald, R. J.M. Farlai, R. Ahlefeldt, and I. Jackson. “Dislocation creep of fine-grained olivine”. In: *J. Geophys. Res.* 116.B01203, (2011). DOI: 10.1029/2009JB007174.

- [389] E. Fehlberg. “Some old and new Runge-Kutta formulas with stepsize control and their error coefficients”. In: *Computing* 34 (1985), pp. 265–270. DOI: 10.1007/BF02253322.
- [390] Marc Fehling and Wolfgang Bangerth. “Algorithms for parallel generic hp-adaptive finite element software”. In: *ACM Transactions on Mathematical Software* 49.3 (2023), pp. 1–26. DOI: 10.1145/3603372.
- [391] JA Fernández-Merodo, JC García-Davalillo, G Herrera, P Mira, and M Pastor. “2D viscoplastic finite element modelling of slow landslides: the Portalet case study (Spain)”. In: *Landslides* 11.1 (2014), pp. 29–42. DOI: 10.1007/s10346-012-0370-4.
- [392] Sylvaine Ferrachat and Yanick Ricard. “Mixing properties in the Earth’s mantle: Effects of the viscosity stratification and of oceanic crust segregation”. In: *Geochemistry, Geophysics, Geosystems* 2.4 (2001). DOI: 10.1029/2000GC000092.
- [393] Sylvaine Ferrachat and Yanick Ricard. “Regular vs. chaotic mantle mixing”. In: *Earth and Planetary Science Letters* 155.1-2 (1998), pp. 75–86. DOI: 10.1016/S0012-821X(97)00200-8.
- [394] David A Field. “Qualitative measures for initial meshes”. In: *International Journal for Numerical Methods in Engineering* 47.4 (2000), pp. 887–906. DOI: 10.1002/(SICI)1097-0207(20000210)47:4<887::AID-NME804>3.0.CO;2-H.
- [395] Mark P Fischer, Michael R Gross, Terry Engelder, and Roy J Greenfield. “Finite-element analysis of the stress distribution around a pressurized crack in a layered elastic medium: implications for the spacing of fluid-driven joints in bedded sedimentary rock”. In: *Tectonophysics* 247.1-4 (1995), pp. 49–64. DOI: 10.1016/0040-1951(94)00200-S.
- [396] Nicolas Flament. “Present-day dynamic topography and lower-mantle structure from palaeogeographically constrained mantle flow models”. In: *Geophysical Journal International* 216.3 (2019), pp. 2158–2182. DOI: 10.1093/gji/ggy526.
- [397] Luce Fleitout and David A Yuen. “Secondary convection and the growth of the oceanic lithosphere”. In: *Physics of the earth and planetary interiors* 36.3-4 (1984), pp. 181–212. DOI: 10.1016/0031-9201(84)90046-3.
- [398] Raymond C Fletcher. “Three-dimensional folding of an embedded viscous layer in pure shear”. In: *Journal of Structural Geology* 13.1 (1991), pp. 87–96.
- [399] MGG Foreman and AF Bennett. “On no-slip boundary conditions for the incompressible Navier-Stokes equations”. In: *Dynamics of atmospheres and oceans* 12.1 (1988), pp. 47–70. DOI: 10.1016/0377-0265(88)90014-0.
- [400] A. Fortin, M. Jardak, J.J. Gervais, and R. Pierre. “Old and New Results on the Two-Dimensional Poiseuille Flow”. In: *J. Comp. Phys.* 115.2 (1994), pp. 455–469. DOI: 10.1006/jcph.1994.1210.
- [401] M. Fortin. “Old and new finite elements for incompressible flows”. In: *Int. J. Num. Meth. Fluids* 1 (1981), pp. 347–364. DOI: 10.1002/flid.1650010406.
- [402] M. Fortin and S. Boivin. “Iterative stabilisation of the bilinear velocity-constant pressure element”. In: *Int. J. Num. Meth. Fluids* 10 (1990), pp. 125–140. DOI: 10.1002/flid.1650100202.
- [403] M. Fortin and A. Fortin. “Experiments with several elements for viscous incompressible flows”. In: *Int. J. Num. Meth. Fluids* 5 (1985), pp. 911–928. DOI: 10.1002/flid.1650051005.
- [404] M. Fortin and M. Soulie. “A non-conforming piecewise quadratic finite element on triangles”. In: *Int. J. Num. Meth. Eng.* 19 (1983), pp. 505–520.
- [405] Michel Fortin and F Thomasset. “Mixed finite-element methods for incompressible flow problems”. In: *Journal of Computational Physics* 31.1 (1979), pp. 113–145.

- [406] L. Fourel, S. Goes, and G. Morra. “The role of elasticity in slab bending”. In: *Geochem. Geophys. Geosyst.* 15 (2014), pp. 4507–4525. DOI: 10.1002/2014GC005535.
- [407] Marc Fournier, Laurent Jolivet, Philippe Davy, and Jean-Charles Thomas. “Backarc extension and collision: an experimental approach to the tectonics of Asia”. In: *Geophysical Journal International* 157.2 (2004), pp. 871–889.
- [408] Christine Mary Rutherford Fowler. *The solid earth: an introduction to global geophysics; second edition*. Cambridge University Press, 2005.
- [409] L.P. Franca and T.J.R. Hughes. “Convergence analyses of Galerkin least-squares methods for symmetric advective-diffusive forms of the Stokes and incompressible Navier-Stokes equations”. In: *Computer Methods in Applied Mechanics and Engineering* 105 (1993), pp. 285–298. DOI: 10.1016/0045-7825(93)90126-I.
- [410] L.P. Franca and S.P. Oliveira. “Pressure bubbles stabilization features in the Stokes problem”. In: *Computer Methods in Applied Mechanics and Engineering* 192 (2003), pp. 1929–1937.
- [411] L.P. Franca, S.P. Oliveira, and M. Sarkis. “Continuous Q1-Q1 Stokes elements stabilised with non-conforming null edge average velocity functions”. In: *Mathematical Models and Methods in Applied Sciences* 17 (03 2007), pp. 439–459. DOI: 10.1142/S021820250700198X.
- [412] Leopoldo P Franca, Sergio L Frey, and Thomas JR Hughes. “Stabilized finite element methods: I. Application to the advective-diffusive model”. In: *Computer Methods in Applied Mechanics and Engineering* 95.2 (1992), pp. 253–276. DOI: 10.1016/0045-7825(92)90143-8.
- [413] Leopoldo P Franca, G Hauke, and A Masud. “Stabilized finite element methods”. In: *FINITE ELEMENT METHODS: 1970’s AND BEYOND*. International Center for Numerical Methods in Engineering (CIMNE), Barcelona, 2004.
- [414] M. Fraters, C. Thieulot, A. van den Berg, and W. Spakman. “The Geodynamic World Builder: a solution for complex initial conditions in numerical modelling”. In: *Solid Earth* 10 (2019), pp. 1785–1807. DOI: 10.5194/se-10-1785-2019.
- [415] M.R.T. Fraters, W. Bangerth, C. Thieulot, A.C. Glerum, and W. Spakman. “Efficient and Practical Newton Solvers for Nonlinear Stokes Systems in Geodynamic Problems”. In: *Geophy. J. Int.* 218 (2019), pp. 873–894. DOI: 10.1093/gji/ggz183.
- [416] S. Frederiksen and J. Braun. “Numerical modelling of strain localisation during extension of the continental lithosphere”. In: *Earth Planet. Sci. Lett.* 188 (2001), pp. 241–251. DOI: 10.1016/S0012-821X(01)00323-5.
- [417] M. Frehner. “3D fold growth rates”. In: *Terra Nova* 26 (2014), pp. 417–424. DOI: 10.1111/ter.12116.
- [418] Alfred M. Freudenthal and Hilda Geiringer. “The mathematical theory of the inelastic continuum”. In: *Handbuch der Physik, Encyclopedia of Physics*. Vol. VI. Springer-Verlag, 1958.
- [419] P.J. Frey and P.-L. George. *Mesh generation*. Hermes Science, 2000.
- [420] C. Froidevaux. “Energy dissipation and geometric structure at spreading plate boundaries”. In: *Earth Planet. Sci. Lett.* 20 (1973), pp. 419–424. DOI: 10.1016/0012-821X(73)90020-4.
- [421] L. Fuchs and H. Schmeling. “A new numerical method to calculate inhomogeneous and time-dependent large deformation of two-dimensional geodynamic flows with application to diapirism”. In: *Geophy. J. Int.* 194.2 (2013), pp. 623–639. DOI: 10.1093/gji/ggt142.
- [422] Hiromi Fujimoto and Yoshibumi Tomoda. “Lithospheric thickness anomaly near the trench and possible driving force of subduction”. In: *Tectonophysics* 112.1-4 (1985), pp. 103–110. DOI: 10.1016/0040-1951(85)90174-X.

- [423] T. Fukuchi. “Numerical calculation of fully-developed laminar flows in arbitrary cross-sections using finite difference method”. In: *AIP Advances* 1 (2011), p. 042109. DOI: 10.1063/1.3652881.
- [424] Javier Fulla, Juan Rodríguez-González, María Charco, Zdenek Martinec, A Negredo, and Antonio Villaseñor. “Perturbing effects of sub-lithospheric mass anomalies in GOCE gravity gradient and other gravity data modelling: Application to the Atlantic-Mediterranean transition zone”. In: *International Journal of Applied Earth Observation and Geoinformation* 35 (2015), pp. 54–69. DOI: 10.1016/j.jag.2014.02.003.
- [425] CW Fuller, SD Willett, D Fisher, and CY Lu. “A thermomechanical wedge model of Taiwan constrained by fission-track thermochronometry”. In: *Tectonophysics* 425.1-4 (2006), pp. 1–24. DOI: 10.1016/j.tecto.2006.05.018.
- [426] P. Fullsack. “An arbitrary Lagrangian-Eulerian formulation for creeping flows and its application in tectonic models”. In: *Geophys. J. Int.* 120 (1995), pp. 1–23. DOI: 10.1111/j.1365-246X.1995.tb05908.x.
- [427] Jean Furstoss. “Approche numérique de l’évolution microstructurale des péridotites”. PhD thesis. 2020.
- [428] Mikito Furuichi. “Numerical modeling of three dimensional self-gravitating Stokes flow problem with free surface”. In: *Procedia Computer Science* 4 (2011), pp. 1506–1515. DOI: 10.1016/j.procs.2011.04.163.
- [429] M. Furuichi, M. Kameyama, and A. Kageyama. “Three-dimensional Eulerian method for large deformation of viscoelastic fluid: Toward plate-mantle simulation”. In: *J. Comp. Phys.* 227 (2008), pp. 4977–4997. DOI: 10.1016/j.jcp.2008.01.052.
- [430] M. Furuichi and D. Nishiura. “Robust coupled fluid-particle simulation scheme in Stokes-flow regime: Toward the geodynamic simulation including granular media”. In: *Geochem. Geophys. Geosyst.* 15 (2014), pp. 2865–2882.
- [431] Mikito Furuichi, Dave A May, and Paul J Tackley. “Development of a Stokes flow solver robust to large viscosity jumps using a Schur complement approach with mixed precision arithmetic”. In: *Journal of Computational Physics* 230.24 (2011), pp. 8835–8851. DOI: 10.1016/j.jcp.2011.09.007.
- [432] Yoshitsugu Furukawa. “Depth of the decoupling plate interface and thermal structure under arcs”. In: *Journal of Geophysical Research: Solid Earth* 98.B11 (1993), pp. 20005–20013.
- [433] Carl W Gable, Richard J O’connell, and Bryan J Travis. “Convection in three dimensions with surface plates: Generation of toroidal flow”. In: *Journal of Geophysical Research: Solid Earth* 96.B5 (1991), pp. 8391–8405. DOI: 10.1029/90JB02743.
- [434] Sashikumaar Ganesan, Gunar Matthies, and Lutz Tobiska. “Local projection stabilization of equal order interpolation applied to the Stokes problem”. In: *Mathematics of Computation* 77.264 (2008), pp. 2039–2060. DOI: 10.1090/S0025-5718-08-02130-3.
- [435] F. Garel, S. Goes, D.R. Davies, J.H. Davies, S.C. Kramer, and C.R. Wilson. “Interaction of subducted slabs with the mantle transition-zone: A regime diagram from 2-D thermo-mechanical models with a mobile trench and an overriding plate”. In: *Geochem. Geophys. Geosyst.* 15.1739–1765 (2014). DOI: 10.1002/2014GC005257.
- [436] Rao Garimella. “Conformal refinement of unstructured quadrilateral meshes”. In: *Proceedings of the 18th International Meshing Roundtable*. Springer, 2009, pp. 31–44.
- [437] D.K. Gartling. *Nachos - A finite element computer program for incompressible flow problems*. Tech. rep. Sand77-1333. Sandia Laboratories, 1978.

- [438] R. Gassmüller, H. Lokavarapu, E. M. Heien, E. G. Puckett, and W. Bangerth. “Flexible and scalable particle-in-cell methods with adaptive mesh refinement for geodynamic computations”. In: *Geochem. Geophys. Geosyst.* 19.9 (2018), pp. 3596–3604. DOI: 10.1029/2018GC007508.
- [439] Rene Gassmüller, Juliane Dannberg, Wolfgang Bangerth, Timo Heister, and Robert Myhill. “On formulations of compressible mantle convection”. In: *Geophysical Journal International* 221.2 (2020), pp. 1264–1280. DOI: 10.1093/gji/ggaa078.
- [440] Rene Gassmüller, Harsha Lokavarapu, Wolfgang Bangerth, and Elbridge Gerry Puckett. “Evaluating the accuracy of hybrid finite element/particle-in-cell methods for modelling incompressible Stokes flow”. In: *Geophysical Journal International* 219.3 (2019), pp. 1915–1938. DOI: 10.1093/gji/ggz405.
- [441] L. Gastaldo. “Methodes de correction de pression pour les ecoulements compressibles: Application aux equations de Navier-Stokes barotropes et au modele de derive”. PhD thesis. Universite de Provence, 2007.
- [442] Andrey A Gavrilov, Konstantin A Finnikov, and Evgeny V Podryabinkin. “Modeling of steady Herschel–Bulkley fluid flow over a sphere”. In: *Journal of Engineering Thermophysics* 26.2 (2017), pp. 197–215. DOI: 10.1134/S1810232817020060.
- [443] T. Geenen et al. “Scalable robust solvers for unstructured FE geodynamic modeling applications: Solving the Stokes equation for models with large localized viscosity contrasts”. In: *Geochem. Geophys. Geosyst.* 10.9 (2009).
- [444] P.-L. George and H. Borouchaki. *Delaunay Triangulation and Meshing*. Hermes, 1998.
- [445] I. Georgiev, J. Kraus, and S. Margenov. “Multilevel algorithms for Rannacher-Turek finite element approximations of 3D elliptic problems”. In: *Computing* 82 (2008), pp. 217–239.
- [446] M. Gerault, T.W. Becker, B.J.P. Kaus, C. Faccenna, L. Moresi, and L. Husson. “The role of slabs and oceanic plate geometry in the net rotation of the lithosphere, trench motions, and slab return flow”. In: *Geochem. Geophys. Geosyst.* 13.4 (2012), Q04001. DOI: 10.1029/2011GC003934.
- [447] M. Gerbault. “Pressure conditions for shear and tensile failure around a circular magma chamber; insight from elasto-plastic modelling”. In: *Geological Society, London, Special Publications* 367 (2012), pp. 111–130.
- [448] M. Gerbault, F. Cappa, and R. Hassani. “Elasto-plastic and hydromechanical models of failure around an infinitely long magma chamber”. In: *Geochem. Geophys. Geosyst.* 13.3 (2012). DOI: 10.1029/2011GC003917.
- [449] M. Gerbault, A.N.B. Poliakov, and M. Daignieres. “Prediction of faulting from the theories of elasticity and plasticity: what are the limits?” In: *Journal of Structural Geology* 20 (1998), pp. 301–320.
- [450] T. Gerya and D.A. Yuen. “Robust characteristics method for modelling multiphase visco-elasto-plastic thermo-mechanical problems”. In: *Phys. Earth. Planet. Inter.* 163 (2007), pp. 83–105. DOI: 10.1016/j.pepi.2007.04.015.
- [451] T.V. Gerya and J.-P. Burg. “Intrusion of ultramafic magmatic bodies into the continental crust: Numerical simulation”. In: *Phys. Earth. Planet. Inter.* 160 (2007), pp. 124–142. DOI: 10.1016/j.pepi.2006.10.004.
- [452] T.V. Gerya, D.A. May, and T. Duretz. “An adaptive staggered grid finite difference method for modeling geodynamic Stokes flows with strongly variable viscosity”. In: *Geochem. Geophys. Geosyst.* 14.4 (2013). DOI: 10.1002/ggge.20078.



- [453] T.V. Gerya and D.A. Yuen. “Characteristics-based marker-in-cell method with conservative finite-differences schemes for modeling geological flows with strongly variable transport properties”. In: *Phys. Earth. Planet. Inter.* 140 (2003), pp. 293–318. DOI: 10.1016/j.pepi.2003.09.006.
- [454] Taras Gerya. “Geodynamics of the early Earth: Quest for the missing paradigm”. In: *Geology* 47.10 (2019), pp. 1006–1007.
- [455] Taras Gerya. *Numerical Geodynamic Modelling*. Cambridge University Press, 2010.
- [456] Taras Gerya. *Numerical Geodynamic Modelling - 2nd edition*. Cambridge University Press, 2019. ISBN: 978-1-107-14314-2.
- [457] Christophe Geuzaine and Jean-François Remacle. “Gmsh: A 3-D finite element mesh generator with built-in pre-and post-processing facilities”. In: *International journal for numerical methods in engineering* 79.11 (2009), pp. 1309–1331. DOI: 10.1002/nme.2579.
- [458] Siavash Ghelichkhan and Hans-Peter Bunge. “The compressible adjoint equations in geodynamics: derivation and numerical assessment”. In: *GEM-International Journal on Geomathematics* 7.1 (2016), pp. 1–30. DOI: 10.1007/s13137-016-0080-5.
- [459] U. Ghia, K.N. Ghia, and C.T. Shin. “High-Re Solutions for incompressible flow using the Navier-Stokes equations and a multigrid method”. In: *J. Comp. Phys.* 48 (1982), pp. 387–411.
- [460] Norman E Gibbs. “Algorithm 509: A hybrid profile reduction algorithm [F1]”. In: *ACM Transactions on Mathematical Software (TOMS)* 2.4 (1976), pp. 378–387.
- [461] Norman E Gibbs, William G Poole Jr, and Paul K Stockmeyer. “A comparison of several bandwidth and profile reduction algorithms”. In: *ACM Transactions on Mathematical Software (TOMS)* 2.4 (1976), pp. 322–330. DOI: xxxx.
- [462] Frederic Gibou, Ronald Fedkiw, and Stanley Osher. “A review of level-set methods and some recent applications”. In: *Journal of Computational Physics* 353 (2018), pp. 82–109. DOI: 10.1016/j.jcp.2017.10.006.
- [463] JA Gil and Maria José Jurado. “Geological interpretation and numerical modelling of salt movement in the Barbastro–Balaguer anticline, southern Pyrenees”. In: *Tectonophysics* 293.3–4 (1998), pp. 141–155.
- [464] François Glaisner and Tayfun E Tezduyar. *Finite element techniques for the Navier-Stokes equations in the primitive variable formulation and the vorticity stream-function formulation*. Tech. rep. 1987.
- [465] G.C. Gleason and J. Tullis. “A flow law for dislocation creep of quartz aggregates determined with the molten salt cell”. In: *Tectonophysics* 247 (1995), pp. 1–23.
- [466] John W Glen. “The creep of polycrystalline ice”. In: *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* 228.1175 (1955), pp. 519–538.
- [467] A. Glerum, C. Thieulot, M. Fraters, C. Blom, and W. Spakman. “Nonlinear viscoplasticity in ASPECT: benchmarking and applications to subduction”. In: *Solid Earth* 9.2 (2018), pp. 267–294. DOI: 10.5194/se-9-267-2018.
- [468] P Glišović, AM Forte, and Robert Moucha. “Time-dependent convection models of mantle thermal structure constrained by seismic tomography and geodynamics: implications for mantle plume dynamics and CMB heat flux”. In: *Geophysical Journal International* 190.2 (2012), pp. 785–815. DOI: 10.1111/j.1365-246X.2012.05549.x.

- [469] Björn Gmeiner, Markus Huber, Lorenz John, Ulrich Rüde, and Barbara Wohlmuth. “A quantitative performance study for Stokes solvers at the extreme scale”. In: *Journal of Computational Science* 17 (2016), pp. 509–521. DOI: 10.1016/j.jocs.2016.06.006.
- [470] C. Goetze and B. Evans. “Stress and temperature in the bending lithosphere as constrained by experimental rock mechanics”. In: *GJI* 59.3 (1979), pp. 463–478.
- [471] DL Goldsby and DL Kohlstedt. “Superplastic deformation of ice: Experimental observations”. In: *Journal of Geophysical Research: Solid Earth* 106.B6 (2001), pp. 11017–11030.
- [472] G. Goltz and M. Böse. “Configurational entropy of critical earthquake populations”. In: *Geophys. Res. Lett.* 29.20 (2002). DOI: 10.1029/2002GL015540.
- [473] G.H. Golub and C.F. van Loan. *Matrix Computations, 4th edition*. John Hopkins University Press, 2013.
- [474] H Gossman. “Slope modelling with changing boundary conditions-effects of climate and lithology”. In: *Z. Geomorph. NF* 25 (1976), pp. 72–88.
- [475] K-D Gottschaldt, U Walzer, RF Hendel, David Robert Stegman, JR Baumgardner, and H-B Mühlhaus. “Stirring in 3-d spherical models of convection in the Earth’s mantle”. In: *Philosophical Magazine* 86.21-22 (2006), pp. 3175–3204. DOI: 10.1080/14786430500197991.
- [476] Klaus-D Gottschaldt, Uwe Walzer, Dave R Stegman, John R Baumgardner, and Hans B Mühlhaus. “Mantle Dynamics—A Case Study”. In: *Advances in Geocomputing*. 2009, pp. 139–181. DOI: 10.1007/978-3-540-85879-9\_5.
- [477] S. Gourvenec. “Bearing capacity under combined loading”. In: *9th Australia New Zealand Conference on Geomechanics, Auckland, New Zealand, 8-11 february 2004*. 2004.
- [478] S. Gourvenec, M. Randolph, and O. Kingsnorth. “Undrained bearing capacity of square and rectangular footings”. In: *International Journal of Geomechanics* 6 (2006), pp. 147–157.
- [479] R Govers and MJR Wortel. “Initiation of asymmetric extension in continental lithosphere”. In: *Tectonophysics* 223.1-2 (1993), pp. 75–96.
- [480] R. Govers and M.J.R. Wortel. “Lithosphere tearing at STEP faults: Response to edges of subduction zones ”. In: *Earth Planet. Sci. Lett.* 236 (2005), pp. 505–523.
- [481] J. Grandy. *Efficient Computation of Volume of Hexahedral Cells*. Tech. rep. UCRL-ID-128886. Lawrence Livermore National Laboratory, 1997.
- [482] R. Gray and R.N. Pysklywec. “Geodynamic models of mature continental collision: Evolution of an orogen from lithospheric subduction to continental retreat/delamination”. In: *J. Geophys. Res.* 117.B03408 (2012). DOI: 10.1029/2011JB008692.
- [483] R. Gray and R.N. Pysklywec. “Influence of viscosity pressure dependence on deep lithospheric tectonics during continental collision”. In: *J. Geophys. Res.* 118 (2013). DOI: 10.1002/jgrb.50220.
- [484] Harry W Green. “Shearing instabilities accompanying high-pressure phase transformations and the mechanics of deep earthquakes”. In: *Proceedings of the National Academy of Sciences* 104.22 (2007), pp. 9133–9138. DOI: 10.1073/pnas.0608045104.
- [485] P.M. Gresho, S.T. Chan, M.A. Christon, and A.C. Hindmarsh. “A little more on stabilised  $Q_1Q_1$  for transient viscous incompressible flow”. In: *Int. J. Num. Meth. Fluids* 21 (1995), pp. 837–856. DOI: 10.1002/flid.1650211005.
- [486] P.M. Gresho and R.L. Lee. “Don’t suppress the wiggles - They’re telling you something!” In: *Computers and Fluids* 9 (1981), pp. 223–253.

- [487] P.M. Gresho, R.L. Lee, R.L. Sani, M.K. Maslanik, and B.E. Eaton. “The consistent Galerkin FEM for computing derived boundary quantities in thermal and/or fluid problems”. In: *Int. J. Num. Meth. Fluids* 7 (1987), pp. 371–394.
- [488] P.M. Gresho and R.L. Sani. *Incompressible flow and the Finite Element Method, vol II - Isothermal Laminar Flow*. John Wiley and Sons, Ltd, 2000. ISBN: 978-0471492504.
- [489] P.M. Gresho and S.B. Sutton. “Application of the FIDAP code to the 8:1 thermal cavity problem”. In: *Int. J. Num. Meth. Fluids* 40 (2002), pp. 1083–1092. DOI: 10.1002/d.394.
- [490] Philip M Gresho. “Some current CFD issues relevant to the incompressible Navier-Stokes equations”. In: *Computer Methods in Applied Mechanics and Engineering* 87.2-3 (1991), pp. 201–252. DOI: 10.1016/0045-7825(91)90006-R.
- [491] PM Gresho and RL Lee. “Partial vindication of the bilinear velocity, piecewise constant pressure element”. In: *Journal of Computational Physics* 60.1 (1985), pp. 161–164. DOI: 10.1016/0021-9991(85)90023-3.
- [492] Ralf Greve. “Application of a polythermal three-dimensional ice sheet model to the Greenland ice sheet: response to steady-state and transient climate scenarios”. In: *Journal of Climate* 10.5 (1997), pp. 901–918.
- [493] Ralf Greve and Heinz Blatter. *Dynamics of ice sheets and glaciers*. Springer Science & Business Media, 2009.
- [494] M. Griebel, T. Dornseifer, and T. Neunhoffer. *Numerical simulation in Fluid Dynamics*. SIAM, 1997.
- [495] D. Griffiths and D. Silvester. *Unstable modes of the Q1-P0 element*. Tech. rep. 257. University of MAnchester/UMIST, 1994.
- [496] D.F. Griffiths. “Finite Elements for Incompressible Flow”. In: *Math. Meth. in the Appl. Sci.* 1 (1979), pp. 16–31.
- [497] Piotr P Grinevich and Maxim A Olshanskii. “An iterative method for the Stokes-type problem with variable viscosity”. In: *SIAM Journal on Scientific Computing* 31.5 (2009), pp. 3959–3978. DOI: 10.1137/08744803.
- [498] Thomas Grombein, Kurt Seitz, and Bernhard Heck. “Optimized formulas for the gravitational field of a tesseract”. In: *Journal of Geodesy* 87.7 (2013), pp. 645–660. DOI: 10.1007/s00190-013-0636-1.
- [499] L. Gross, L. Bourguin, A. Hale, and H.-B. Mühlhaus. “Interface modeling in incompressible media using level sets in Escript”. In: *Phys. Earth. Planet. Inter.* 163 (2007), pp. 23–34.
- [500] Ólafur Gudmundsson and Malcolm Sambridge. “A regionalized upper mantle (RUM) seismic model”. In: *Journal of Geophysical Research: Solid Earth* 103.B4 (1998), pp. 7121–7136.
- [501] J.-L. Guermond, R. Pasquetti, and Bojan Popov. “Entropy viscosity method for nonlinear conservation laws”. In: *J. Comp. Phys.* (2011). DOI: 10.1016/j.jcp.2010.11.043.
- [502] Jean Luc Guermond, Richard Pasquetti, and Bojan Popov. “Entropy viscosity for conservation equations”. In: *V European Conference on Computational Fluid Dynamics (Eccomas CFD 2010)*. 2010.
- [503] Jean-Luc Guermond and Richard Pasquetti. “Entropy viscosity method for high-order approximations of conservation laws”. In: *Spectral and high order methods for partial differential equations*. Springer, 2011, pp. 411–418.

- [504] JM Guerrero, Julian P Lowman, and Paul J Tackley. “Spurious transitions in convective regime due to viscosity clipping: ramifications for modeling planetary secular cooling”. In: *Geochemistry, Geophysics, Geosystems* 20.7 (2019), pp. 3450–3468. DOI: 10.1029/2019GC008385.
- [505] F. Gueydan, C. Morency, and J.-P. Brun. “Continental rifting as a function of lithosphere mantle strength”. In: *Tectonophysics* 460 (2008), pp. 83–93. DOI: 10.1016/j.tecto.2008.08.012.
- [506] G Guj and F Stella. “A vorticity-velocity method for the numerical solution of 3D incompressible flows”. In: *Journal of Computational Physics* 106.2 (1993), pp. 286–298.
- [507] M. Gunzburger. *Finite Element Methods for Viscous Incompressible Flows: A Guide to Theory, Practice and Algorithms*. Academic, Boston, 1989.
- [508] Max D Gunzburger and Janet S Peterson. “On finite element approximations of the streamfunction-vorticity and velocity-vorticity equations”. In: *International journal for numerical methods in fluids* 8.10 (1988), pp. 1229–1240. DOI: 10.1002/flid.1650081010.
- [509] Peng Guo, Leihua Yao, and Desheng Ren. “Simulation of three-dimensional tectonic stress fields and quantitative prediction of tectonic fracture within the Damintun Depression, Liaohe Basin, northeast China”. In: *Journal of Structural Geology* 86 (2016), pp. 211–223. DOI: 10.1016/j.jsg.2016.03.007.
- [510] Anshul Gupta, George Karypis, and Vipin Kumar. “Highly Scalable Parallel Algorithms for Sparse Matrix Factorization”. In: *IEEE Transactions on Parallel and Distributed Systems* 8.5 (1997), pp. 502–520.
- [511] Anshul Gupta, Seid Koric, and Thomas George. “Sparse Matrix Factorization on Massively Parallel Computers”. In: *SC09 (International Conference for High Performance Computing, Networking, Storage and Analysis)*. 2009.
- [512] P.S. Gupta and A.S. Gupta. “Squeezing flow between parallel plates”. In: *Wear* 45.2 (1977), pp. 177–185. DOI: 10.1016/0043-1648(77)90072-2.
- [513] Michael Gurnis. “Stirring and mixing in the mantle by plate-scale flow: Large persistent blobs and long tendrils coexist”. In: *Geophysical Research Letters* 13.13 (1986), pp. 1474–1477. DOI: 10.1029/GL013i013p01474.
- [514] Michael Gurnis and Geoffrey F Davies. “Mixing in numerical models of mantle convection incorporating plate kinematics”. In: *Journal of Geophysical Research: Solid Earth* 91.B6 (1986), pp. 6375–6395. DOI: 10.1029/JB091iB06p06375.
- [515] Michael Gurnis and Geoffrey F Davies. “The effect of depth-dependent viscosity on convective mixing in the mantle and the possible survival of primitive mantle”. In: *Geophysical Research Letters* 13.6 (1986), pp. 541–544. DOI: 10.1029/GL013i006p00541.
- [516] Michael Gurnis, Christophe Eloy, and Shijie Zhong. “Free-surface formulation of mantle convection - II. Implication for subduction-zone observables”. In: *Geophysical Journal International* 127.3 (1996), pp. 719–727.
- [517] M-A Gutscher et al. “Thermal modeling of the SW Ryukyu forearc (Taiwan): Implications for the seismogenic zone and the age of the subducting Philippine Sea Plate (Huatung Basin)”. In: *Tectonophysics* 692 (2016), pp. 131–142.
- [518] Pierre Guyot and John E Dorn. “A critical review of the Peierls mechanism”. In: *Canadian Journal of Physics* 45.2 (1967), pp. 983–1016. DOI: 10.1139/p67-073.
- [519] E. Hairer, S.P. Nørsett, and G. Wanner. *Solving Ordinary Differential Equations I*. Springer, 1993.

- [520] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*. Springer-Verlag, Berlin, 1991.
- [521] A.J. Hale, L. Bourgouin, and H.B. Muehlhaus. “Using the level set method to model endogenous lava dome growth”. In: *J. Geophys. Res.* 112 (2007), B03213. DOI: 10.1029/2006JB004445.
- [522] A.J. Hale, K.-D. Gottschaldt, G. Rosenbaum, L. Bourgouin, M. Bauchy, and Hans Mühlhaus. “Dynamics of slab tear faults: Insights from numerical modelling”. In: *Tectonophysics* 483 (2010), pp. 58–70. DOI: 10.1016/j.tecto.2009.05.019.
- [523] Andrea Hampel, Jens Lüke, Thomas Krause, and Ralf Hetzel. “Finite-element modelling of glacial isostatic adjustment (GIA): Use of elastic foundations at material boundaries versus the geometrically non-linear formulation”. In: *Computers & geosciences* 122 (2019), pp. 1–14. DOI: 10.1016/j.cageo.2018.08.002.
- [524] H. Han. “A finite element approximation of Navier-Stokes equations using nonconforming elements”. In: *Journal of Computational Mathematics* 2.1 (1984), pp. 77–88. DOI: xxxx.
- [525] J. Handin. “On the Coulomb-Mohr failure criterion”. In: *J. Geophys. Res.* 74.22 (1969), p. 5343.
- [526] Emmanuel Hanert, Vincent Legat, and Éric Deleersnijder. “A comparison of three finite elements to solve the linear shallow water equations”. In: *Ocean Modelling* 5.1 (2003), pp. 17–35. DOI: 10.1016/S1463-5003(02)00012-4.
- [527] P. Hansbo. “A nonconforming rotated  $Q_1$  approximation on tetrahedra”. In: *Computer Methods in Applied Mechanics and Engineering* 200 (2011), pp. 1311–1316. DOI: 10.1016/j.cma.2010.11.002.
- [528] Peter Hansbo, Mats Larson, and Mats G Larson. “A simple nonconforming bilinear element for the elasticity problem”. In: *Trends in Computational Structural Mechanics*. 2001.
- [529] Peter Hansbo and Mats G Larson. “Discontinuous Galerkin and the Crouzeix–Raviart element: application to elasticity”. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 37.1 (2003), pp. 63–72. DOI: 10.1051/m2an:2003020.
- [530] D.L. Hansen. “A meshless formulation for geodynamic modeling”. In: *J. Geophys. Res.* 108 (2003). DOI: 10.1029/2003JB002460.
- [531] D.L. Hansen and S.B. Nielsen. “Why rifts invert in compression”. In: *Tectonophysics* 373 (2003), pp. 5–24.
- [532] D.L. Hansen, S.B. Nielsen, and H. Lykke-Andersen. “The post-Triassic evolution of the Sorgenfrei-Tornquist Zone - results from thermo-mechanical modelling”. In: *Tectonophysics* 328 (2000), pp. 245–267.
- [533] U Hansen and A Ebel. “Experiments with a numerical model related to mantle convection: boundary layer behaviour of small-and large scale flows”. In: *Physics of the earth and planetary interiors* 36.3-4 (1984), pp. 374–390. DOI: 10.1016/0031-9201(84)90058-X.
- [534] U. Hansen and D.A. Yuen. “Effects of depth-dependent thermal expansivity on the interaction of thermal-chemical plumes with a compositional boundary”. In: *Physics of the Earth and Planetary Interiors* 86.1-3 (1994), pp. 205–221. DOI: 10.1016/0031-9201(94)05069-4.
- [535] U. Hansen and D.A. Yuen. “Evolutionary structures in double-diffusive convection in magma chambers”. In: *Geophysical Research Letters* 14.11 (1987), pp. 1099–1102. DOI: 10.1029/GL014i011p01099.

- [536] U. Hansen and D.A. Yuen. “High Rayleigh number regime of temperature-dependent viscosity convection and the Earth’s early thermal history”. In: *Geophysical Research Letters* 20.20 (1993), pp. 2191–2194. DOI: 10.1029/93GL02416.
- [537] U. Hansen and D.A. Yuen. “Nonlinear physics of double-diffusive convection in geological systems”. In: *Earth Science Reviews* 29.1-4 (1990), pp. 385–399. DOI: 10.1016/0012-8252(90)90050-6.
- [538] U. Hansen, D.A. Yuen, and S.E. Kroening. “Mass and Heat Transport in Strongly Time-Dependent Thermal Convection at Infinite Prandtl Number”. In: *Geophysical & Astrophysical Fluid Dynamics* 63.1-4 (1992), pp. 67–89. DOI: 10.1080/03091929208228278.
- [539] U. Hansen, D.A. Yuen, and A.V. Malevsky. “Comparison of steady-state and strongly chaotic thermal convection at high Rayleigh number”. In: *Physical Review A* 46.8 (1992), pp. 4742–4754. DOI: 10.1103/PhysRevA.46.4742.
- [540] Ulrich Hansen and Adolf Ebel. “Time-dependent thermal convection-a possible explanation for a multiscale flow in the Earth’s mantle”. In: *Geophysical Journal International* 94.2 (1988), pp. 181–191. DOI: 10.1111/j.1365-246X.1988.tb05895.x.
- [541] Ulrich Hansen and David A Yuen. “Dynamical influences from thermal-chemical instabilities at the core-mantle boundary”. In: *Geophysical Research Letters* 16.7 (1989), pp. 629–632. DOI: 10.1029/GL016i007p00629.
- [542] Ulrich Hansen, David A Yuen, SE Kroening, and TB Larsen. “Dynamical consequences of depth-dependent thermal expansivity and viscosity on mantle circulations and thermal structure”. In: *Physics of the earth and planetary interiors* 77.3-4 (1993), pp. 205–223. DOI: 10.1016/0031-9201(93)90099-U.
- [543] Ulrich Hansen, David A Yuen, and Sherri E Kroening. “Effects of depth-dependent thermal expansivity on mantle circulations and lateral thermal anomalies”. In: *Geophysical Research Letters* 18.7 (1991), pp. 1261–1264. DOI: 10.1029/91GL01288.
- [544] L. Hanyk, J. Moser, D.A. Yuen, and C. Matyska. “Time-domain approach for the transient responses in stratified viscoelastic Earth models”. In: *Geophysical Research Letters* 22.10 (1995), pp. 1285–1288. DOI: 10.1029/95GL01087.
- [545] Helmut Harder. “Numerical simulation of thermal convection with Maxwellian viscoelasticity”. In: *Journal of non-newtonian fluid mechanics* 39.1 (1991), pp. 67–88. DOI: 10.1016/0377-0257(91)80004-4.
- [546] Helmut Harder. “Phase transitions and the three-dimensional planform of thermal convection in the Martian mantle”. In: *Journal of Geophysical Research: Planets* 103.E7 (1998), pp. 16775–16797. DOI: 10.1029/98JE01543.
- [547] Doug P Hardin, Edward B Saff, et al. “Discretizing manifolds via minimum energy points”. In: *Notices of the AMS* 51.10 (2004), pp. 1186–1194. DOI: xxxx.
- [548] F.H. Harlow and J.E. Welch. “Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface”. In: *The physics of fluids* 8.12 (1965), p. 2182.
- [549] Nathan J Harris, Steven E Benzley, and Steven J Owen. “Conformal Refinement of All-Hexahedral Element Meshes Based on Multiple Twist Plane Insertion”. In: *IMR*. 2004, pp. 157–168.
- [550] E.H. Hartz and Y.Y. Podlachikov. “Toasting the jelly sandwich: The effect of shear heating on lithospheric geotherms and strength”. In: *Geology* 36.4 (2008), pp. 331–334.
- [551] N.A. Haskell. “The Motion of a Viscous Fluid Under a Surface Load”. In: *Physics* 6 (1935), pp. 265–269. DOI: 10.1063/1.1745329.

- [552] R. Hassan, N. Flament, M. Gurnis, D.J. Bower, and D. Müller. “Provenance of plumes in global convection models”. In: *Geochem. Geophys. Geosyst.* 16 (2015), pp. 1465–1489. DOI: 10.1002/2015GC005751.
- [553] R. Hassani and J. Chéry. “Anelasticity explains topography associated with Basin and Range normal faulting”. In: *Geology* 24.12 (1996), pp. 1095–1098. DOI: 10.1130/0091-7613(1996)024<1095:AETAWB>2.3.CO;2.
- [554] Hiroshi Hayashi and Akira Kageyama. “Yin–Yang–Zhong grid: An overset grid system for a sphere”. In: *Journal of Computational Physics* 305 (2016), pp. 895–905. DOI: 10.1016/j.jcp.2015.11.016.
- [555] Y. He, E.G. Puckett, and M.I. Billen. “A discontinuous Galerkin method with a bound preserving limiter for the advection of non-diffusive fields in solid Earth geodynamics”. In: *Phys. Earth. Planet. Inter.* 263 (2017), pp. 23–37.
- [556] B. Heck and K. Seitz. “A comparison of the tesseroid, prism and point-mass approaches for mass reductions in gravity field modelling”. In: *J. Geodesy* 81 (2007), pp. 121–136. DOI: 10.1007/s00190-006-0094-0.
- [557] HJ van Heck and PJ Tackley. “Plate tectonics on super-Earths: equally or more likely than on Earth”. In: *Earth and Planetary Science Letters* 310.3-4 (2011), pp. 252–261. DOI: 10.1016/j.epsl.2011.07.029.
- [558] Otto M Heeres, Akke SJ Suiker, and René de Borst. “A comparison between the Perzyna viscoplastic model and the consistency viscoplastic model”. In: *European Journal of Mechanics-A/Solids* 21.1 (2002), pp. 1–12. DOI: 10.1016/S0997-7538(01)01188-3.
- [559] Juan C Heinrich, Peter S Huyakorn, Olgierd C Zienkiewicz, and AR0353 Mitchell. “An ‘upwind’ finite element scheme for two-dimensional convective transport equation”. In: *International Journal for Numerical Methods in Engineering* 11.1 (1977), pp. 131–143. DOI: 10.1002/nme.1620110113.
- [560] T. Heister, J. Dannberg, R. Gassmüller, and W. Bangerth. “High Accuracy Mantle Convection Simulation through Modern Numerical Methods. II: Realistic Models and Problems”. In: *Geophys. J. Int.* 210.2 (2017), pp. 833–851. DOI: 10.1093/gji/ggx195.
- [561] Christian Helanow and Josefin Ahlkrone. “Stabilized equal low-order finite elements in ice sheet modeling—accuracy and robustness”. In: *Computational Geosciences* 22.4 (2018), pp. 951–974. DOI: 10.1007/s10596-017-9713-5.
- [562] Heinrich Hencky. “Zur Theorie plastischer Deformationen und der hierdurch im Material hervorgerufenen Nachspannungen”. In: *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik* 4.4 (1924), pp. 323–334. DOI: 10.1002/zamm.19240040405.
- [563] P. Hénon, P. Ramet, and J. Roman. “PaStiX: A High-Performance Parallel Direct Solver for Sparse Symmetric Definite Systems”. In: *Parallel Computing* 28.2 (2002), pp. 301–321.
- [564] K.J. Bathe, ed. *Three-dimensional spherical shell convection at infinite Prandtl number using the ‘cubed sphere’ method*. Proceedings Second MIT Conference on Computational Fluid and Solid Mechanics June 17-20, 2003. 2003, pp. 931–933. DOI: xxxx.
- [565] John W Hernlund, Christine Thomas, and Paul J Tackley. “A doubling of the post-perovskite phase boundary and structure of the Earth’s lowermost mantle”. In: *Nature* 434.7035 (2005), p. 882. DOI: 10.1038/nature03472.

- [566] Marco Herwegh, T Poulet, Ali Karrech, and Klaus Regenauer-Lieb. “From transient to steady state deformation and grain size: A thermodynamic approach using elasto-visco-plastic numerical modeling”. In: *Journal of Geophysical Research: Solid Earth* 119.2 (2014), pp. 900–918. DOI: 10.1002/2013JB010701.
- [567] M.R. Hestenes and E. Stiefel. “Methods of Conjugate Gradients for Solving Linear Systems”. In: *Journal of Research of the National Bureau of Standards* 49.6 (1952), pp. 409–436.
- [568] J.S. Hesthaven and T. Warbuton. “Nodal Discontinuous Galerkin Methods. Algorithms, Analysis, and Applications”. In: *Texts in Applied Mathematics*. Ed. by J.E. Marsden, L. Sirovich, and S.S. Antman. Springer, 2008.
- [569] J.M. Hewitt, D.P. McKenzie, and N.O. Weiss. “Dissipative heating in convective flows”. In: *J. Fluid Mech.* 68.4 (1975), pp. 721–738. DOI: 10.1017/S002211207500119X.
- [570] N. Hilairet et al. “High-Pressure creep of Serpentine, interseismic deformation, and initiation of subduction”. In: *Science* 318 (2007), pp. 1910–1913.
- [571] B. Hillebrand, C. Thieulot, T. Geenen, A. P. van den Berg, and W. Spakman. “Using the level set method in geodynamical modeling of multi-material flows and Earth’s free surface”. In: *Solid Earth* 5.2 (2014), pp. 1087–1098. DOI: 10.5194/se-5-1087-2014.
- [572] RD van der Hilst, MV e Hoop, P Wang, S-H Shim, P Ma, and L Tenorio. “Seismostratigraphy and thermal structure of Earth’s core-mantle boundary region”. In: *science* 315.5820 (2007), pp. 1813–1817. DOI: 10.1126/science.1137867.
- [573] Kei Hirose, John Brodholt, Thorne Lay, David A Yuen, et al. “An Introduction to Post-Perovskite: The Last Mantle Phase Transition”. In: *GEOPHYSICAL MONOGRAPH-AMERICAN GEOPHYSICAL UNION* 174 (2007), p. 1.
- [574] C.W Hirt and B.D Nichols. “Volume of fluid (VOF) method for the dynamics of free boundaries”. In: *J. Computational Physics* 39.1 (1981), pp. 201–225. DOI: 10.1016/0021-9991(81)90145-5.
- [575] C.W. Hirt, A. Amsden, and J.L. Cook. “An Arbitrary Lagrangian-Eulerian computing method for all flow speeds”. In: *J. Comp. Phys.* 14 (1974), pp. 227–253. DOI: 10.1016/0021-9991(74)90051-5.
- [576] Christian Hirt, SJ Claessens, Michael Kuhn, and WE Featherstone. “Kilometer-resolution gravity field of Mars: MGM2011”. In: *Planetary and Space Science* 67.1 (2012), pp. 147–154. DOI: 10.1016/j.pss.2012.02.006.
- [577] G. Hirth and D. L. Kohlstedt. “Water in the oceanic upper mantle: Implications for rheology, melt extraction and the evolution of the lithosphere”. In: *Earth Planet. Sci. Lett.* 144 (1996), pp. 93–108. DOI: 10.1016/0012-821X(96)00154-9.
- [578] G. Hirth and D.L. Kohlstedt. “Rheology of the upper mantle and the mantle wedge: A view from the experimentalists”. In: *in Inside the Subduction Factory, ed. J. Eiler, Geophysical Monograph American Geophysical Union, Washington, D.C.* 138 (2003), pp. 83–105.
- [579] BE Hobbs, H-B Mühlhaus, and A Ord. “Instability, softening and localization of deformation”. In: *Geological Society, London, Special Publications* 54.1 (1990), pp. 143–165. DOI: 10.1144/GSL.SP.1990.054.01.15.
- [580] BE Hobbs and A Ord. “Numerical simulation of shear band formation in a frictional-dilatational material”. In: *Ingenieur-Archiv* 59.3 (1989), pp. 209–220. DOI: 10.1007/BF00532251.
- [581] NRA Hoffman and DP McKenzie. “The destruction of geochemical heterogeneities by differential fluid motions during mantle convection”. In: *Geophysical Journal International* 82.2 (1985), pp. 163–206. DOI: 10.1111/j.1365-246X.1985.tb05134.x.



- [582] Klaus A Hoffmann and Steve T Chiang. “Computational fluid dynamics for engineers, Vol. 1”. In: *Wichita, Engineering Education System* (1993), pp. 124–137.
- [583] A.M. Hofmeister and D.A. Yuen. “Critical phenomena in thermal conductivity: Implications for lower mantle dynamics”. In: *Journal of Geodynamics* 44.3-5 (2007), pp. 186–199. DOI: 10.1016/j.jog.2007.03.002.
- [584] AJ Hogg and GP Matson. “Slumps of viscoplastic fluids on slopes”. In: *Journal of Non-Newtonian Fluid Mechanics* 158.1-3 (2009), pp. 101–112. DOI: 10.1016/j.jnnfm.2008.07.003.
- [585] John C Holden and Peter Vogt. “Graphic solutions to problems of plumacy”. In: *Eos, Transactions American Geophysical Union* 58.7 (1977), pp. 573–580. DOI: 10.1029/E0058i007p00573.
- [586] Bjørn Holmedal. “Spin and vorticity with vanishing rigid-body rotation during shear in continuum mechanics”. In: *Journal of the Mechanics and Physics of Solids* 137 (2020), p. 103835. DOI: 10.1016/j.jmps.2019.103835.
- [587] Satoru Honda. “Thermal structure beneath Tohoku, northeast Japan”. In: *Tectonophysics* 112.1-4 (1985), pp. 69–102. DOI: 10.1016/0040-1951(85)90173-8.
- [588] Seungwoo Hong, Kyungsoo Kim, and Sungyun Lee. “Modified cross-grid finite elements for the stokes problem”. In: *Applied mathematics letters* 16.1 (2003), pp. 59–64. DOI: 10.1016/S0893-9659(02)00144-1.
- [589] P. Hood and C. Taylor. “Navier-Stokes equations using mixed interpolation”. In: *Finite element methods in flow problems*. Ed. by J.T. Oden, R.H. Gallagher, O.C. Zienkiewicz, and C. Taylor. University of Alabama: Huntsville Press, 1974.
- [590] André Horbach, Hans-Peter Bunge, and Jens Oeser. “The adjoint method in geodynamics: derivation from a general operator formulation and application to the initial condition problem in a high resolution mantle circulation model”. In: *GEM-International Journal on Geomathematics* 5.2 (2014), pp. 163–194.
- [591] André Horbach, Marcus Mohr, and Hans-Peter Bunge. “A semi-analytic accuracy benchmark for Stokes flow in 3-D spherical mantle convection codes”. In: *GEM-International Journal on Geomathematics* 11.1 (2020), pp. 1–35. DOI: 10.1007/s13137-019-0137-3.
- [592] C.O. Horgan. “Korn’s inequalities and their applications in continuum mechanics”. In: *SIAM Review* 37.4 (1995), pp. 491–511. DOI: 10.1137/1037123.
- [593] GA Houseman. “Boundary conditions and efficient solution algorithms for the potential function formulation of the 3-D viscous flow equations”. In: *Geophysical Journal International* 100.1 (1990), pp. 33–38. DOI: 10.1111/j.1365-246X.1990.tb04565.x.
- [594] GA Houseman. “The thermal structure of mantle plumes: axisymmetric or triple-junction?” In: *Geophysical Journal International* 102.1 (1990), pp. 15–24. DOI: 10.1111/j.1365-246X.1990.tb00527.x.
- [595] Greg A Houseman, D Po McKenzie, and Peter Molnar. “Convective instability of a thickened boundary layer and its relevance for the thermal evolution of continental convergent belts”. In: *Journal of Geophysical Research: Solid Earth* 86.B7 (1981), pp. 6115–6132. DOI: 10.1029/JB086iB07p06115.
- [596] Gregory Houseman and Philip England. “Finite strain calculations of continental deformation: 1. Method and general results for convergent zones”. In: *Journal of Geophysical Research: Solid Earth* 91.B3 (1986), pp. 3651–3663. DOI: 10.1029/JB091iB03p03651.

- [597] Albert T Hsui. “Numerical simulation of finite-amplitude thermal convection with large viscosity variation in axisymmetric spherical geometry: effect of mechanical boundary conditions”. In: *Tectonophysics* 50.2-3 (1978), pp. 147–162. DOI: 10.1016/0040-1951(78)90132-4.
- [598] Albert T Hsui, Tang Xiao-Ming, and M Nafi Toksöz. “On the dip angle of subducting plates”. In: *Tectonophysics* 179.3-4 (1990), pp. 163–175.
- [599] Chongyu Hua. “An inverse transformation for quadrilateral isoparametric elements: analysis and application”. In: *Finite elements in analysis and design* 7.2 (1990), pp. 159–166. DOI: 10.1016/0168-874X(90)90007-2.
- [600] Jinshui Huang and Geoffrey F Davies. “Stirring in three-dimensional mantle convection models and implications for geochemistry: Passive tracers”. In: *Geochemistry, Geophysics, Geosystems* 8.3 (2007). DOI: 10.1029/2006GC001312.
- [601] Ch. Huber, A. Parmigiani, B. Chopard, M. Manga, and O. Bachmann. “Lattice Boltzmann model for melting with natural convection”. In: *International Journal of Heat and Fluid Flow* 29 (200), pp. 1469–1480. DOI: 10.1016/j.ijheatfluidflow.2008.05.002.
- [602] Aurélia Hubert-Ferrari, Geoffrey King, Isabelle Manighetti, Rolando Armijo, Bertrand Meyer, and Paul Tapponnier. “Long-term elasticity in the continental lithosphere; modelling the Aden Ridge propagation and the Anatolian extrusion process”. In: *Geophysical Journal International* 153.1 (2003), pp. 111–132. DOI: 10.1046/j.1365-246X.2003.01872.x.
- [603] A. Huerta and W.K. Liu. “Viscous flow with large free surface motion”. In: *Computer Methods in Applied Mechanics and Engineering* 69 (1988), pp. 277–324.
- [604] T.J.R. Hughes. *The Finite Element Method. Linear Static and Dynamic Finite Element Analysis*. Dover Publications, Inc., 2000. ISBN: 0-486-41181-8.
- [605] T.J.R. Hughes and A. Brooks. “A theoretical framework for Petrov-Galerkin methods with discontinuous weighting functions: application to the streamline-upwind procedure”. In: *Finite Elements in Fluids* 4 (1982), pp. 47–65.
- [606] T.J.R. Hughes, L.P. Franca, and M. Balestra. “A new finite element formulation for computational fluid dynamics: V. Circumventing the Babuška-Brezzi condition: A stable Petrov-Galerkin formulation of the Stokes problem accommodating equal-order interpolations”. In: *Computer Methods in Applied Mechanics and Engineering* 59.1 (1986), pp. 85–99.
- [607] T.J.R. Hughes, L.P. Franca, and M. Balestra. “A new finite element formulation for computational fluid dynamics: VII. The Stokes problem with various well-posed boundary conditions: symmetric formulations that converge for all velocity/pressure spaces”. In: *Computer Methods in Applied Mechanics and Engineering* 65 (1987), pp. 85–96.
- [608] T.J.R. Hughes, W.K. Liu, and A. Brooks. “Finite element analysis of Incompressible viscous flows by the penalty function formulation”. In: *J. Comp. Phys.* 30 (1979), pp. 1–60. DOI: 10.1016/0021-9991(79)90086-X.
- [609] Thomas JR Hughes, Wing Kam Liu, and Thomas K Zimmermann. “Lagrangian-Eulerian finite element formulation for incompressible viscous flows”. In: *Computer methods in applied mechanics and engineering* 29.3 (1981), pp. 329–349. DOI: 10.1016/0045-7825(81)90049-9.
- [610] Thomas JR Hughes, Michel Mallet, and Akira Mizukami. “A new finite element formulation for computational fluid dynamics: II. Beyond SUPG”. In: *Computer methods in applied mechanics and engineering* 54.3 (1986), pp. 341–355. DOI: 10.1016/0045-7825(86)90110-6.

- [611] Hoon Huh, Choong Ho Lee, and Wei H. Yang. “A general algorithm for plastic flow simulation by finite element limit analysis”. In: *International Journal of Solids and Structures* 36 (1999), pp. 1193–1207. DOI: 10.1016/S0020-7683(97)00347-8.
- [612] R. Huismans and C. Beaumont. “Depth-dependent extension, two-stage breakup and cratonic underplating at rifted margins”. In: *Nature* 473 (2011), pp. 74–79. DOI: 10.1038/nature09988.
- [613] R. S. Huismans and C. Beaumont. “Symmetric and asymmetric lithospheric extension: Relative effects of frictional-plastic and viscous strain softening”. In: *J. Geophys. Res.* 108 (B10).2496 (2003). DOI: 10.1029/2002JB002026.
- [614] R.S. Huismans and C. Beaumont. “Roles of lithospheric strain softening and heterogeneity in determining the geometry of rifts and continental margins”. In: *Imaging, Mapping and Modelling Continental Lithosphere Extension and Breakup*. Vol. 282. Geological Society, London, Special Publications, 2007, pp. 111–138. DOI: 10.1144/SP282.6.
- [615] R.S. Huismans, S.J.H. Buiter, and C. Beaumont. “Effect of plastic-viscous layering and strain softening on mode selection during lithospheric extension”. In: *J. Geophys. Res.* 110 (2005), B02406. DOI: 10.1029/2004JB003114.
- [616] J. van Hunen, A.P. van den Berg, and N.J. Vlaar. “On the role of subducting oceanic plateaus in the development of shallow flat subduction”. In: *Tectonophysics* 352.3-4 (2002), pp. 317–333. DOI: 10.1016/S0040-1951(02)00263-9.
- [617] C. Hüttig and K. Stemmer. “The spiral grid: A new approach to discretize the sphere and its application to mantle convection”. In: *Geochem. Geophys. Geosyst.* 9.2 (2008). DOI: 10.1029/2007GC001581.
- [618] G. Iaffaldano and H.-P. Bunge. “Relating rapid plate-motion variations to plate-boundary forces in global coupled models of the mantle/lithosphere system: Effects of topography and friction”. In: *Tectonophysics* 474.1-2 (2009), pp. 393–404. DOI: 10.1016/j.tecto.2008.10.035.
- [619] S. Idelsohn, M. Storti, and N. Nigro. “Stability analysis of mixed finite element formulations with special mention of equal-order interpolations”. In: *Int. J. Num. Meth. Fluids* 20 (1995), pp. 1003–1022. DOI: 10.1002/flid.1650200819.
- [620] F. Ilinca and D. Pelletier. “Computation of accurate nodal derivatives of finite element solutions: The finite node displacement method”. In: *Int. J. Num. Meth. Eng.* 71 (2007), pp. 1181–1207. DOI: 10.1002/nme.1979.
- [621] Wadi H Imseeh and Khalid A Alshibli. “3D finite element modelling of force transmission and particle fracture of sand”. In: *Computers and Geotechnics* 94 (2018), pp. 184–195. DOI: 10.1016/j.compgeo.2017.09.008.
- [622] Computational Inelasticity. 1998.
- [623] Tobin Isaac, Georg Stadler, and Omar Ghattas. “Solution of nonlinear Stokes equations discretized by high-order finite elements on nonconforming and anisotropic meshes, with application to ice sheet dynamics”. In: *SIAM Journal on Scientific Computing* 37.6 (2015), B804–B833. DOI: 10.1137/140974407.
- [624] Hideyuki Ishii, Yoshihisa Sakurai, and Takeo Maruyama. “Effect of soccer shoe upper on ball behaviour in curve kicks”. In: *Scientific reports* 4 (2014), p. 6067. DOI: 10.1038/srep06067.
- [625] A Ismail-Zadeh, A Korotkii, G Schubert, and I Tsepelev. “Quasi-reversibility method for data assimilation in models of mantle dynamics”. In: *Geophysical Journal International* 170.3 (2007), pp. 1381–1398. DOI: 10.1111/j.1365-246X.2007.03496.x.

- [626] Alik Ismail-Zadeh and Paul Tackley. *Computational Methods for Geodynamics*. Cambridge University Press, 2010.
- [627] J. Ita and S.D. King. “Sensitivity of convection with an endothermic phase change to the form of governing equations, initial conditions, boundary conditions, and equation of state”. In: *J. Geophys. Res.* 99.B8 (1994), pp. 15, 919–15, 938. DOI: 10.1029/94JB00852.
- [628] Y. Iwase. “Three-dimensional infinite Prandtl number convection in a spherical shell with temperature-dependent viscosity”. In: *J. Geomag. Geoelectr.* 48 (1996), pp. 1499–1514. DOI: 10.5636/jgg.48.1499.
- [629] Michel HG Jacobs and Bernard HWS de Jong. “Placing constraints on phase equilibria and thermophysical properties in the system MgO–SiO<sub>2</sub> by a thermodynamically consistent vibrational method”. In: *Geochimica et Cosmochimica Acta* 71.14 (2007), pp. 3630–3655. DOI: 10.1016/j.gca.2007.05.010.
- [630] Wolfgang R Jacoby and Harro Schmeling. “Convection experiments and the driving mechanism”. In: *Geologische Rundschau* 70.1 (1981), pp. 207–230. DOI: 10.1007/BF01764323.
- [631] J.C. Jaeger. *Elasticity, Fracture and Flow*. John Wiley and Sons, Inc., 1969.
- [632] S. Jammes and R.S. Huismans. “Structural styles of mountain building: Controls of lithospheric rheologic stratification and extensional inheritance”. In: *Journal of Geophysical Research: Solid Earth* 117.B10 (2012). DOI: 10.1029/2012JB009376.
- [633] Gang-Won Jang, Sangkeun Lee, Yoon Young Kim, and Dongwoo Sheen. “Topology optimization using non-conforming finite elements: three-dimensional case”. In: *International journal for numerical methods in engineering* 63.6 (2005), pp. 859–875. DOI: 10.1002/nme.1302.
- [634] O. Jaoul, J. Tullis, and A. Kronenberg. “The effect of varying water contents on the creep behavior of Heavtree quartzite”. In: *J. Geophys. Res.* 89.B6 (1984). DOI: 10.1029/JB089iB06p04298.
- [635] G.T. Jarvis. “Effects of curvature on two-dimensional models of mantle convection: cylindrical polar coordinates”. In: *J. Geophys. Res.* 98.B3 (1993), pp. 4477–4485. DOI: 10.1029/92JB02117.
- [636] Gary T Jarvis and WR Peltier. “Mantle convection as a boundary layer phenomenon”. In: *Geophysical Journal International* 68.2 (1982), pp. 389–427. DOI: 10.1111/j.1365-246X.1982.tb04907.x.
- [637] Gary T. Jarvis and Dan P. McKenzie. “Convection in a compressible fluid with infinite Prandtl number”. In: *Journal of Fluid Mechanics* 96.3 (1980), pp. 515–583. DOI: 10.1017/S002211208000225X.
- [638] GT Jarvis. “Two-dimensional numerical models of mantle convection”. In: *Advances in geophysics*. Vol. 33. 1991, pp. 1–80. DOI: 10.1016/S0065-2687(08)60440-9.
- [639] C. Jaupart and J.-C. Mareschal. *Heat Generation and Transport in the Earth*. Cambridge, 2011.
- [640] P. Jay, A. Magnin, and J.-M. Piau. “Viscoplastic Fluid Flow Through a Sudden Axisymmetric Expansion”. In: *AIChE Journal* 47.10 (2001), pp. 2155–2166. DOI: 10.1002/aic.690471004.
- [641] Raymond Jeanloz and S Morris. “Is the mantle geotherm subadiabatic?” In: *Geophysical research letters* 14.4 (1987), pp. 335–338. DOI: 10.1029/GL014i004p00335.
- [642] Raymond Jeanloz and S Morris. “Temperature distribution in the crust and mantle”. In: *Annual Review of Earth and Planetary Sciences* 14.1 (1986), pp. 377–415. DOI: xxxx.

- [643] Eleanor W Jenkins, Volker John, Alexander Linke, and Leo G Rebholz. “On the parameter choice in grad-div stabilization for the Stokes equations”. In: *Advances in Computational Mathematics* 40.2 (2014), pp. 491–516. DOI: 10.1007/s10444-013-9316-1.
- [644] P. Jenny, S.B. Pope, M. Muradoglu, and D.A. Caughey. “A Hybrid Algorithm for the Joint PDF Equation of Turbulent Reactive Flows”. In: *J. Comp. Phys.* 166 (2001), pp. 218–252. DOI: 10.1006/jcph.2000.6646.
- [645] Youngmok Jeon, Hyun Nam, Dongwoo Sheen, and Kwangshin Shim. “A class of nonparametric DSSY nonconforming quadrilateral elements”. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 47.6 (2013), pp. 1783–1796. DOI: 10.1051/m2an/2013088.
- [646] L. Jiang, J. Liu, J. Zhang, and Z. Feng. “Analytic Expressions for the Gravity Gradient Tensor of 3D Prisms with Depth-Dependent Density”. In: *Surv. Geophys.* 39 (2018), pp. 337–363. DOI: 10.1007/s10712-017-9455-x.
- [647] Liqing Jiao, Paul Tapponnier, Frédéric-victor Donzé, Luc Scholtès, Yves Gaudemer, and Xiwei Xu. “Discrete element modeling of southeast Asia’s 3D lithospheric deformation during the Indian collision”. In: *Journal of Geophysical Research: Solid Earth* 128.1 (2023). DOI: 10.1029/2022JB025578.
- [648] Stephen Jiménez, Ravindra Duddu, and Jeremy Bassis. “An updated-Lagrangian damage mechanics formulation for modeling the creeping flow and fracture of ice sheets”. In: *Computer Methods in Applied Mechanics and Engineering* 313 (2017), pp. 406–432. DOI: 10.1016/j.cma.2016.09.034.
- [649] Volker John. “A posteriori  $L^2$ -error estimates for the nonconforming  $P_1/P_0$ -finite element discretization of the Stokes equations”. In: *Journal of computational and applied mathematics* 96.2 (1998), pp. 99–116. DOI: 10.1016/S0377-0427(98)00095-8.
- [650] Volker John. *Finite Element Methods for Incompressible Flow Problems*. Springer, 2016. ISBN: 978-3-319-45749-9. DOI: 10.1007/978-3-319-45750-5.
- [651] Volker John. “Higher order finite element methods and multigrid solvers in a benchmark problem for the 3D Navier–Stokes equations”. In: *International Journal for Numerical Methods in Fluids* 40.6 (2002), pp. 775–798. DOI: 10.1002/d.377.
- [652] Volker John, Kristine Kaiser, and Julia Novo. “Finite element methods for the incompressible Stokes equations with variable viscosity”. In: *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik* 96.2 (2016), pp. 205–216. DOI: 10.1002/zamm.201400291.
- [653] Volker John, Songul Kaya, and Julia Novo. *Finite Element Error Analysis Of A Mantle Convection Model*. Vol. 15. 4–5. 2018, pp. 677–698. DOI: xxx.
- [654] Volker John and Petr Knobloch. “On discontinuity-capturing methods for convection-diffusion equations”. In: *Numerical mathematics and advanced applications*. 2006, pp. 336–344. DOI: 10.1007/978-3-540-34288-5\_27.
- [655] Volker John, Alexander Linke, Christian Merdon, Michael Neilan, and Leo G Rebholz. “On the divergence constraint in mixed finite element methods for incompressible flows”. In: *SIAM review* 59.3 (2017), pp. 492–544. DOI: 10.1137/15M1047696.
- [656] L. Jolivet, P. Davy, and P. Cobbold. “Right-lateral shear along the Northwest Pacific margin and the India-Eurasia collision”. In: *Tectonics* 9.6 (1990), pp. 1409–1419. DOI: 10.1029/TC009i006p01409.

- [657] MS Joun and MC Lee. “Quadrilateral finite-element generation and mesh quality control for metal forming simulation”. In: *International Journal for Numerical Methods in Engineering* 40.21 (1997), pp. 4059–4075. DOI: 10.1002/(SICI)1097-0207(19971115)40:21<4059::AID-NME249>3.0.CO;2-E.
- [658] Guillaume Jouvét and Jacques Rappaz. “Analysis and finite element approximation of a non-linear stationary Stokes problem arising in glaciology”. In: *Advances in Numerical Analysis* 2011 (2011). DOI: 10.1155/2011/164581.
- [659] M Jull and D McKenzie. “The effect of deglaciation on mantle melting beneath Iceland”. In: *J. Geophys. Res.* 101.B10 (1996), pp. 21815–21828. DOI: 10.1029/96JB01308.
- [660] Mikhail K Kaban, Alexey G Petrunin, Harro Schmeling, and Meysam Shahraki. “Effect of decoupling of lithospheric plates on the observed geoid”. In: *Surveys in Geophysics* 35.6 (2014), pp. 1361–1373. DOI: 10.1007/s10712-014-9281-3.
- [661] L.M. Kachanov. *Fundamentals of the Theory of Plasticity*. Dover Publications, Inc., 2004. ISBN: 0-486-43583-0.
- [662] B.J. Kadlec, G.A. Dorn, H.M. Tufo, and D.A. Yuen. “Interactive 3-D computation of fault surfaces using level sets”. In: *Visual Geosciences* 13.1 (2008), pp. 133–138. DOI: 10.1007/s10069-008-0016-9.
- [663] A. Kageyama and T. Sato. ““Yin-Yang grid”: An overset grid in spherical geometry”. In: *Geochem. Geophys. Geosyst.* 5.9 (2004), 10.1029/2004GC000734.
- [664] Juhani Kakkuri. “Fennoscandian Land Uplift: Past, Present and Future”. In: *From the Earth’s Core to Outer Space*. Ed. by Ilmari Haapala. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 127–136. ISBN: 978-3-642-25550-2. DOI: 10.1007/978-3-642-25550-2\_8.
- [665] M. Kameyama, A. Kageyamab, and T. Sato. “Multigrid-based simulation code for mantle convection in spherical shell using Yin-Yang grid”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 19–32.
- [666] M. Kameyama, D.A. Yuen, and S.-I. Karato. “Thermal-mechanical effects of low-temperature plasticity (the Peierls mechanism) on the deformation of a viscoelastic shear zone”. In: *Earth Planet. Sci. Lett.* 168 (1999), pp. 159–172. DOI: 10.1016/j.pepi.2006.02.005.
- [667] Masanori Kameyama. “ACuTEMan: A multigrid-based mantle convection simulation code and its optimization to the Earth Simulator”. In: *J. Earth Simulator* 4 (2005), pp. 2–10. DOI: xxxx.
- [668] Masanori Kameyama, Akira Kageyama, and Tetsuya Sato. “Multigrid iterative algorithm using pseudo-compressibility for three-dimensional mantle convection with strongly variable viscosity”. In: *Journal of Computational Physics* 206.1 (2005), pp. 162–181. DOI: 10.1016/j.jcp.2004.11.030.
- [669] Takeo Kaneko, Tomoeiki Nakakuki, and Hikaru Iwamori. “Mechanical coupling of the motion of the surface plate and the lower mantle slab: Effects of viscosity hill, yield strength, and depth-dependent thermal expansivity”. In: *Physics of the Earth and Planetary Interiors* 294 (2019), p. 106274. DOI: 10.1016/j.pepi.2019.106274.
- [670] Elias Karabelas, Gundolf Haase, Gernot Plank, and Christoph M Augustin. “Versatile stabilized finite element formulations for nearly and fully incompressible solid mechanics”. In: *Computational mechanics* 65.1 (2020), pp. 193–215. DOI: 10.1007/s00466-019-01760-w.

- [671] S. Karato, M.R. Riedel, and D.A. Yuen. “Rheological structure and deformation of subducted slabs in the mantle transition zone: implications for mantle circulation and deep earthquakes”. In: *Phys. Earth. Planet. Inter.* 127 (2001), pp. 83–108. DOI: 10.1016/S0031-9201(01)00223-0.
- [672] S.-I. Karato and H. Jung. “Effects of pressure on high-temperature dislocation creep in olivine”. In: *Philosophical Magazine* 83.3 (2003), pp. 401–414. DOI: 10.1080/0141861021000025829.
- [673] S.-I. Karato and P. Wu. “Rheology of the Upper Mantle: A synthesis”. In: *Science* 260 (1993), pp. 771–778. DOI: 10.1126/science.260.5109.771.
- [674] Shun-Ichiro Karato. *Deformation of Earth Materials*. Cambridge University Press, 2008. ISBN: 978-0-521-84404-8.
- [675] I. Katayama and S. Karato. “Low-temperature, high-stress deformation of olivine under water-saturated conditions”. In: *Phys. Earth. Planet. Inter.* 168.3-4 (2008), pp. 125–133.
- [676] Tomoo Katsura, Akira Yoneda, Daisuke Yamazaki, Takashi Yoshino, and Eiji Ito. “Adiabatic temperature profile in the mantle”. In: *Physics of the Earth and Planetary Interiors* 183.1-2 (2010), pp. 212–218. DOI: 10.1016/j.pepi.2010.07.001.
- [677] Rafael Katzman, Uri S ten Brink, and Jian Lin. “Three-dimensional modeling of pull-apart basins: Implications for the tectonics of the Dead Sea Basin”. In: *Journal of Geophysical Research: Solid Earth* 100.B4 (1995), pp. 6295–6312. DOI: 10.1029/94JB03101.
- [678] B. Kaus. “Modelling approaches to geodynamic processes”. PhD thesis. ETH Zurich, 2005.
- [679] B.J.P. Kaus. “Factors that control the angle of shear bands in geodynamic numerical models of brittle deformation”. In: *Tectonophysics* 484 (2010), pp. 36–47. DOI: 10.1016/j.tecto.2009.08.042.
- [680] B.J.P. Kaus and T.W. Becker. “Effects of elasticity on the Rayleigh-Taylor instability: implications for large-scale geodynamics”. In: *Geophy. J. Int.* 168.843–862 (2007). DOI: 10.1111/j.1365-246X.2006.03201.x.
- [681] B.J.P. Kaus, H. Mühlhaus, and D.A. May. “A stabilization algorithm for geodynamic numerical simulations with a free surface”. In: *Phys. Earth. Planet. Inter.* 181 (2010), pp. 12–20. DOI: 10.1016/j.pepi.2010.04.007.
- [682] B.J.P. Kaus and Y.Y. Podlachikov. “Initiation of localized shear zones in viscoelastoplastic rocks”. In: *J. Geophys. Res.* 111.B04412 (2006). DOI: 10.1029/2005JB003652.
- [683] B.J.P. Kaus et al. “Forward and Inverse Modelling of Lithospheric Deformation on Geological Timescales”. In: *NIC Symposium 2016* (2016), pp. 299–307.
- [684] M. Kawaguti. “Numerical solution of the Navier-Stokes equations for the flow in a two-dimensional cavity”. In: *Journal of the Physical Society of Japan* 16.12 (1961), pp. 2307–2315.
- [685] Takaaki Kawazoe, Shun-ichiro Karato, Kazuhiko Otsuka, Zhicheng Jing, and Mainak Mookherjee. “Shear deformation of dry polycrystalline olivine under deep upper mantle conditions using a rotational Drickamer apparatus (RDA)”. In: *Physics of the Earth and Planetary Interiors* 174.1-4 (2009), pp. 128–137. DOI: 10.1016/j.pepi.2008.06.027.
- [686] N. Kechkar and D. Silvester. “Analysis of locally stabilised mixed finite element methods for the Stokes problem”. In: *Mathematics of Computation* 58.197 (1992), pp. 1–10.
- [687] P. van Keken, D.A. Yuen, and A. van den Berg. “Pulsating diapiric flows: Consequences of vertical variations in mantle creep laws”. In: *Earth Planet. Sci. Lett.* 112 (1992), pp. 179–194. DOI: 10.1016/0012-821X(92)90015-N.

- [688] P. van Keken and S. Zhong. “Mixing in a 3D spherical model of present-day mantle convection”. In: *Earth Planet. Sci. Lett.* 171 (1999), pp. 533–547.
- [689] P.E. van Keken, D.A. Yuen, and A.P. van den Berg. “The effects of shallow rheological boundaries in the upper mantle on inducing shorter time scales of diapiric flows”. In: *Geophysical Research Letters* 20.18 (1993), pp. 1927–1930. DOI: 10.1029/93GL01768.
- [690] T. Keller, D.A. May, and B.J.P. Kaus. “Numerical modelling of magma dynamics coupled to tectonic deformation of lithosphere and crust”. In: *Geophy. J. Int.* 195.3 (2013), pp. 1406–1442. DOI: 10.1093/gji/ggt306.
- [691] L. H. Kellogg and S. D. King. “The effect of temperature dependent viscosity on the structure of new plumes in the mantle: Results of a finite element model in a spherical, axisymmetric shell”. In: *Earth and Planetary Science Letters* 148.1-2 (1997), pp. 13–26. DOI: 10.1016/S0012-821X(97)00025-3.
- [692] LH Kellogg and DL Turcotte. “Mixing and the distribution of heterogeneities in a chaotically convecting mantle”. In: *Journal of Geophysical Research: Solid Earth* 95.B1 (1990), pp. 421–432. DOI: 10.1029/JB095iB01p00421.
- [693] Louise H Kellogg, Bradford H Hager, and Rob D Van Der Hilst. “Compositional stratification in the deep mantle”. In: *Science* 283.5409 (1999), pp. 1881–1884. DOI: 10.1126/science.283.5409.1881.
- [694] Louise H Kellogg and Cheryl A Stewart. “Mixing by chaotic convection in an infinite Prandtl number fluid and implications for mantle convection”. In: *Physics of Fluids A: Fluid Dynamics* 3.5 (1991), pp. 1374–1378. DOI: 10.1063/1.858067.
- [695] B.L.N. Kennett. “On the density distribution within the Earth”. In: *Geophysical Journal International* 132.2 (1998), pp. 374–382. DOI: 10.1046/j.1365-246x.1998.00451.x.
- [696] B.L.N. Kennett, E.R. Engdahl, and R. Buland. “Travel times for global earthquake location and phase association”. In: *Geophy. J. Int.* 122 (1995), pp. 108–124. DOI: 10.1111/j.1365-246X.1991.tb06724.x.
- [697] D.F. Keppie, C.A. Currie, and C. Warren. “Subduction erosion modes: comparing finite element numerical models with the geological record”. In: *Earth Planet. Sci. Lett.* 287 (2009), pp. 241–254. DOI: 10.1016/j.epsl.2009.08.009.
- [698] A.S. Khan and S. Huang. *Continuum theory of plasticity*. Wiley and sons, 1995.
- [699] W.S. Kiefer and B. Hager. “Geoid anomalies and dynamic topography from convection in cylindrical geometry: applications to mantle plumes on Earth and Venus”. In: *Geophy. J. Int.* 108 (1992), pp. 198–214. DOI: 10.1111/j.1365-246X.1992.tb00850.x.
- [700] S. King et al. “A community benchmark for 2D Cartesian compressible convection in the Earth’s mantle”. In: *Geophy. J. Int.* 180 (2010), pp. 73–87.
- [701] S. D. King. “On topography and geoid from 2-D stagnant lid convection calculations”. In: *Geochemistry, Geophysics, Geosystems* 10.3 (2009), n/a–n/a. DOI: 10.1029/2008GC002250.
- [702] S.D. King and D.L. Anderson. “An alternative mechanism of flood basalt formation”. In: *Earth Planet. Sci. Lett.* 136 (1995), pp. 269–279.
- [703] S.D. King, D.J. Frost, and D.C. Rubie. “Why cold slabs stagnate in the transition zone”. In: *Geology* 43.3 (2015), pp. 231–234. DOI: 10.1130/G36320.1.
- [704] Scott D King. “Mantle convection, the asthenosphere, and Earth’s thermal history”. In: *Geological Society of America Special Papers* 514 (2015), SPE514–07. DOI: 10.1130/2015.2514(07).



- [705] Scott D King. “Reconciling laboratory and observational models of mantle rheology in geodynamic modelling”. In: *Journal of Geodynamics* 100 (2016), pp. 33–50. DOI: 10.1016/j.jog.2016.03.005.
- [706] Scott D King and Guy Masters. “An inversion for radial viscosity structure using seismic tomography”. In: *Geophysical Research Letters* 19.15 (1992), pp. 1551–1554. DOI: 10.1029/92GL01700.
- [707] Á Király, Clinton P Conrad, and LN Hansen. “Evolving viscous anisotropy in the upper mantle and its geodynamic implications”. In: *Geochem. Geophys. Geosyst.* 21 (2020), e2020GC009159. DOI: 10.1029/2020GC009159.
- [708] S.H. Kirby and A.K. Kronenberg. “Rheology of the lithosphere: Selected topics”. In: *Reviews of Geophysics* 25.6 (1987). DOI: 10.1029/RG025i006p01219.
- [709] Matthias Kirchhart, Sven Gross, and Arnold Reusken. “Analysis of an XFEM discretization for Stokes interface problems”. In: *SIAM Journal on Scientific Computing* 38.2 (2016), A1019–A1043. DOI: 10.1137/15M1011779.
- [710] P. Kloucek, B. Li, and M. Luskin. “Analysis of a class of nonconforming finite elements for crystalline microstructures”. In: *Mathematics of Computation* 65.215 (1996), pp. 1111–1135.
- [711] Erik A. Kneller, Markus Albertz, Garry D. Karner, and Christopher A. Johnson. “Testing inverse kinematic models of paleocrustal thickness in extensional systems with high-resolution forward thermo-mechanical models”. In: *Geochem. Geophys. Geosyst.* (2013).
- [712] Matthew G. Knepley. *Computational Science I*. Rice University - Department of Computational and Applied Mathematics, 2017.
- [713] P. Knobloch. “On Korn’s inequality for nonconforming Finite Elements”. In: *Technische Mechanik* 20.3 (2000), pp. 205–214.
- [714] P. Knobloch and L. Tobiska. “Stabilisation methods of bubble type for the  $Q_1/Q_1$  element applied to the incompressible Navier-Stokes equations”. In: *Mathematical Modelling and Numerical Analysis* 34.1 (2000), pp. 85–107. DOI: 10.1051/m2an:2000132.
- [715] Petr Knobloch. “On the definition of the SUPG parameter”. In: *Electronic Transactions on Numerical Analysis* 32 (2008), pp. 76–89. DOI: xxxx.
- [716] Dana A Knoll and David E Keyes. “Jacobian-free Newton–Krylov methods: a survey of approaches and applications”. In: *Journal of Computational Physics* 193.2 (2004), pp. 357–397. DOI: 10.1016/j.jcp.2003.08.010.
- [717] G.M. Kobelkov and M.A. Olshanskii. “Effective preconditioning of Uzawa type schemes for a generalized Stokes problem”. In: *Numer. Math.* 86 (2000), pp. 443–470.
- [718] D. E. Koglin Jr., S. R. Ghias, S. D. King, G. T. Jarvis, and J. P. Lowman. “Mantle convection with reversing mobile plates: A benchmark study”. In: *Geochemistry, Geophysics, Geosystems* 6.9 (2005). DOI: 10.1029/2005GC000924.
- [719] Dimitri Komatitsch, Seiji Tsuboi, Jeroen Tromp, A Levander, and G Nolet. “The spectral-element method in seismology”. In: *GEOPHYSICAL MONOGRAPH-AMERICAN GEOPHYSICAL UNION* 157 (2005), p. 205.
- [720] L. Komzsik and P. Poschmann. “Iterative solution techniques for finite element applications”. In: *Finite Elements in Analysis and Design* 14 (1993), pp. 373–379.
- [721] Henk Kooi and Christopher Beaumont. “Escarpment evolution on high-elevation rifted margins: Insights derived from a surface processes model that combines diffusion, advection, and reaction”. In: *Journal of Geophysical Research: Solid Earth* 99.B6 (1994), pp. 12191–12209. DOI: 10.1029/94JB00047.

- [722] U Kopitzke et al. “Finite element convection models: comparison of shallow and deep mantle convection, and temperatures in the mantle”. In: *Journal of Geophysics* 46.1 (1979), pp. 97–121. DOI: xxx.
- [723] Jun Korenaga. “Energetics of mantle convection and the fate of fossil heat”. In: *Geophysical Research Letters* 30.8 (2003). DOI: 10.1029/2003GL016982.
- [724] J.R. Koseff and R.L. Street. “The Lid-Driven Cavity Flow: A Synthesis of Qualitative and Quantitative Observations”. In: *J. Fluids Eng* 106 (1984), pp. 390–398.
- [725] D. Kourounis, A. Fuchs, and O. Schenk. “Towards the Next Generation of Multiperiod Optimal Power Flow Solvers”. In: *IEEE Transactions on Power Systems* PP.99 (2018), pp. 1–10. DOI: 10.1109/TPWRS.2017.2789187.
- [726] L.I.G. Kovaszny. “Laminar flow behind a two-dimensional grid”. In: *Mathematical Proceedings of the Cambridge Philosophical Society* 44.1 (1948), pp. 58–62. DOI: 10.1017/S0305004100023999.
- [727] Peter Kovesi. “Good Colour Maps: How to Design Them”. In: *CoRR* abs/1509.03700 (2015). arXiv: 1509.03700. URL: <http://arxiv.org/abs/1509.03700>.
- [728] M. Krabbendam. “Sliding of temperate basal ice on a rough, hard bed: creep mechanisms, pressure melting, and implications for ice streaming”. In: *The Cryosphere* 10 (2016), pp. 1915–1932. DOI: 10.5194/tc-10-1915-2016.
- [729] Rolf Krahel and Eberhard Bänsch. “COMPUTATIONAL COMPARISON BETWEEN THE TAYLOR–HOOD AND THE CONFORMING CROUZEIX–RAVIART ELEMENT”. In: *Proceedings of ALGORITHM*. 2005, pp. 369–379. DOI: xxx.
- [730] S.C. Kramer, D.R. Davies, and C.R. Wilson. “Analytical solutions for mantle flow in cylindrical and spherical shells”. In: *Geosci. Model Dev.* 14 (2021), pp. 1899–1919. DOI: 10.5194/gmd-14-1899-2021.
- [731] Stephan C Kramer, Cian R Wilson, and D Rhodri Davies. “An implicit free surface algorithm for geodynamical simulations”. In: *Physics of the Earth and Planetary Interiors* 194 (2012), pp. 25–37.
- [732] M. Kronbichler, T. Heister, and W. Bangerth. “High accuracy mantle convection simulation through modern numerical methods”. In: *Geophy. J. Int.* 191 (2012), pp. 12–29. DOI: 10.1111/j.1365-246X.2012.05609.x.
- [733] Marcin Krotkiewski and Marcin Dabrowski. “Parallel symmetric sparse matrix–vector product on scalar multi-core CPUs”. In: *Parallel Computing* 36.4 (2010), pp. 181–198. DOI: 10.1016/j.parco.2010.02.003.
- [734] E.J. Kubatko, B.A. Yeager, and A.L. Maggi. “New computationally efficient quadrature formulas for triangular prism elements”. In: *Computers and Fluids* 73 (2013), pp. 187–201. DOI: 10.1016/j.compfluid.2013.01.002.
- [735] Hsaio-Lan Kuo. “Solution of the non-linear equations of cellular convection and heat transport”. In: *Journal of Fluid Mechanics* 10.4 (1961), pp. 611–634. DOI: 10.1017/S0022112061000408.
- [736] NJ Kusznir and MHP Bott. “Stress concentration in the upper lithosphere caused by underlying visco-elastic creep”. In: *Tectonophysics* 43.3-4 (1977), pp. 247–256. DOI: 10.1016/0040-1951(77)90119-6.
- [737] NJ Kusznir and RG Park. “The strength of intraplate lithosphere”. In: *Physics of the Earth and Planetary Interiors* 36.3-4 (1984), pp. 224–235. DOI: 10.1016/0031-9201(84)90048-7.

- [738] Oh-In Kwon and C. Park. “A new quadrilateral MINI element for Stokes equations”. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 48 (2014), pp. 955–968. DOI: 10.1051/m2an/2013129.
- [739] P. Labbé and A. Garon. “A robust implementation of Zienkiewicz and Shu’s local patch recovery method”. In: *Communications in Numerical Methods in Engineering* 11 (1995), pp. 427–434. DOI: 10.1002/cnm.1640110507.
- [740] Olga Aleksandrovna Ladyzhenskaya. *The mathematical theory of viscous incompressible flow*. Vol. 2. Gordon and Breach New York, 1969.
- [741] B.P. Lamichhane. “A quadrilateral ’mini’ finite element for the Stokes problem using a single bubble function”. In: *International Journal of Numerical Analysis and Modeling* 14.6 (2017), pp. 869–878.
- [742] Bishnu P Lamichhane. “A mixed finite element method for nearly incompressible elasticity and Stokes equations using primal and dual meshes with quadrilateral and hexahedral grids”. In: *Journal of computational and applied mathematics* 260 (2014), pp. 356–363. DOI: 10.1016/j.cam.2013.09.056.
- [743] Bishnu P Lamichhane. “Inf-sup stable finite-element pairs based on dual meshes and bases for nearly incompressible elasticity”. In: *IMA Journal of Numerical Analysis* 29.2 (2009), pp. 404–420. DOI: 10.1093/imanum/drn013.
- [744] W. Landry, L. Hodkinson, and S. Kientz. *Gale User Manual*. Tech. rep. CIG, VPAC, 2011.
- [745] Ulrich Langer and Martin Neumüller. “Direct and iterative solvers”. In: *Computational Acoustics*. Springer, 2018, pp. 205–251. DOI: 10.1007/978-3-319-59038-7\_5.
- [746] Hans Petter Langtangen, Timothy J Barth, and Michael Griebel. *Python scripting for computational science*. Vol. 3. Springer, 2008.
- [747] BE Larock and LR Herrmann. “Improved Flux Prediction Using Low Order Finite Elements”. In: *International Conference on Finite Elements in Water Resources, Part II*. 1976.
- [748] T.B. Larsen and D.A. Yuen. “Fast plumeheads: Temperature-dependent versus non-Newtonian rheology”. In: *Geophysical Research Letters* 24.16 (1997), pp. 1995–1998. DOI: 10.1029/97GL01886.
- [749] T.B. Larsen and D.A. Yuen. “Ultrafast upwelling bursting through the upper mantle”. In: *Earth and Planetary Science Letters* 146.3-4 (1997), pp. 393–399. DOI: 10.1016/S0012-821X(96)00247-6.
- [750] T.B. Larsen, D.A. Yuen, A.V. Malevsky, and J.L. Smedsmo. “Dynamics of strongly time-dependent convection with non-Newtonian temperature-dependent viscosity”. In: *Physics of the Earth and Planetary Interiors* 94.1-2 (1996), pp. 75–103. DOI: 10.1016/0031-9201(95)03082-4.
- [751] Tine B Larsen, David A Yuen, Jiří Moser, and Bengt Fornberg. “A high-order finite-difference method applied to large Rayleigh number mantle convection”. In: *Geophysical & Astrophysical Fluid Dynamics* 84.1-2 (1997), pp. 53–83. DOI: 10.1080/03091929708208973.
- [752] Gautier Laurent, Guillaume Caumon, and Mark Jessell. “Interactive editing of 3D geological structures and tectonic history sketching via a rigid element method”. In: *Computers & geosciences* 74 (2015), pp. 71–86.
- [753] A. Lavecchia, C. Thieulot, F. Beekman, S. Cloetingh, and S. Clark. “Lithosphere erosion and continental breakup: Interaction of extension, plume upwelling and melting”. In: *Earth Planet. Sci. Lett.* 467 (2017), pp. 89–98. DOI: 10.1016/j.epsl.2017.03.028.

- [754] Luc L Lavier, W Roger Buck, and Alexei NB Poliakov. “Factors controlling normal fault offset in an ideal brittle layer”. In: *Journal of Geophysical Research: Solid Earth* 105.B10 (2000), pp. 23431–23442.
- [755] W. Layton. “Weak imposition of ‘no-slip’ conditions in finite element methods”. In: *Computers and Mathematics with Applications* 38 (1999), pp. 129–142. DOI: 10.1016/S0898-1221(99)00220-5.
- [756] Laetitia Le Pourhiet, Dave A May, Lucas Huille, Louise Watremez, and Sylvie Leroy. “A genetic link between transform and hyper-extended margins”. In: *Earth and Planetary Science Letters* 465 (2017), pp. 184–192.
- [757] E-S Lee, Charles Moulinec, Rui Xu, Damien Violeau, Dominique Laurence, and Peter Stansby. “Comparisons of weakly compressible and truly incompressible algorithms for the SPH mesh free particle method”. In: *Journal of computational Physics* 227.18 (2008), pp. 8417–8436. DOI: 10.1016/j.jcp.2008.06.005.
- [758] Junhwan Lee, Rodrigo Salgado, and Sooil Kim. “Bearing capacity of circular footings under surcharge using state-dependent finite element analysis”. In: *Computers and Geotechnics* 32 (2005), pp. 445–457. DOI: 10.1016/j.compgeo.2005.07.005.
- [759] R. Lee, P. Gresho, and R. Sani. “Smoothing techniques for certain primitive variable solutions of the Navier-Stokes equations”. In: *Int. J. Num. Meth. Eng.* 14 (1979), pp. 1785–1804.
- [760] Young-Ju Lee and Hengguang Li. “Axisymmetric Stokes equations in polygonal domains: regularity and finite element approximations”. In: *Computers & Mathematics with Applications* 64.11 (2012), pp. 3500–3521. DOI: 10.1016/j.camwa.2012.08.014.
- [761] R.S. Lehmann, M. Lukacova-Medvidova, B.J.P. Kaus, and A.A. Popov. “Comparison of continuous and discontinuous Galerkin approaches for variable-viscosity Stokes flow”. In: *Z. Angew. Math. Mech.* 96.6 (2015), pp. 733–746. DOI: 10.1002/zamm.201400274.
- [762] A.M. Leitch, V. Steinbach, and D.A. Yuen. “Centerline temperature of mantle plumes in various geometries: Incompressible flow”. In: *Journal of Geophysical Research B: Solid Earth* 101.B10 (1996), pp. 21829–21846. DOI: 10.1029/96JB01784.
- [763] A.M. Leitch, D.A. Yuen, and G. Sewell. “Mantle convection with internal heating and pressure-dependent thermal expansivity”. In: *Earth and Planetary Science Letters* 102.2 (1991), pp. 213–232. DOI: 10.1016/0012-821X(91)90009-7.
- [764] V. Lemiale, H.-B. Mühlhaus, L. Moresi, and J. Stafford. “Shear banding analysis of plastic models formulated for incompressible viscous flows”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 177–186. DOI: 10.1016/j.pepi.2008.07.038.
- [765] Vincent Lemiale, H-B Mühlhaus, Catherine Meriaux, L Moresi, and L Hodkinson. “Rate effects in dense granular materials: Linear stability analysis and the fall of granular columns”. In: *International Journal for Numerical and Analytical Methods in Geomechanics* 35.2 (2011), pp. 293–308. DOI: 10.1002/nag.895.
- [766] A. Lenardic and W. M. Kaula. “A numerical treatment of geodynamic viscous flow problems involving the advection of material interfaces”. In: *Journal of Geophysical Research: Solid Earth* 98.B5 (1993), pp. 8243–8260. DOI: 10.1029/92JB02858.
- [767] A. Lenardic and W. M. Kaula. “Near-surface thermal/chemical boundary layer convection at infinite Prandtl number: two-dimensional numerical experiments”. In: *Geophysical Journal International* 126.3 (1996), pp. 689–711. DOI: 10.1111/j.1365-246X.1996.tb04698.x.
- [768] Adrian Lenardic, CM Cooper, and L Moresi. “A note on continents and the Earth’s Urey ratio”. In: *Physics of the Earth and Planetary Interiors* 188.1-2 (2011), pp. 127–130. DOI: 10.1016/j.pepi.2011.06.008.

- [769] W. Leng and S. Zhong. “Implementation and application of adaptive mesh refinement for thermochemical mantle convection studies”. In: *Geochem. Geophys. Geosyst.* 12.4 (2011). DOI: 10.1029/2010GC003425.
- [770] W. Leng and S. Zhong. “Viscous heating, adiabatic heating and energetic consistency in compressible mantle convection”. In: *Geophys. J. Int.* 173 (2008), pp. 693–702. DOI: 10.1111/j.1365-246X.2008.03745.x.
- [771] Wei Leng, Lili Ju, Yan Xie, Tao Cui, and Max Gunzburger. “Finite element three-dimensional Stokes ice sheet dynamics model with enhanced local mass conservation”. In: *Journal of Computational Physics* 274 (2014), pp. 299–311.
- [772] Y Leroy and M Ortiz. “Finite element analysis of strain localization in frictional materials”. In: *International Journal for Numerical and Analytical Methods in Geomechanics* 13.1 (1989), pp. 53–74. DOI: 10.1002/nag.1610130106.
- [773] P. LeTallec. “Compatibility condition and existence results in discrete finite incompressible elasticity”. In: *Computer Methods in Applied Mechanics and Engineering* 27 (1981), pp. 239–259.
- [774] P. LeTallec and V. Ruas. “On the convergence of the bilinear-velocity constant-pressure finite element method in viscous flow”. In: *Computer Methods in Applied Mechanics and Engineering* 54 (1986), pp. 235–243.
- [775] Frank G Lether. “Computation of double integrals over a triangle”. In: *Journal of Computational and Applied Mathematics* 2.3 (1976), pp. 219–224. DOI: 10.1016/0771-050X(76)90008-5.
- [776] Einat Lev and Bradford H Hager. “Rayleigh–Taylor instabilities with anisotropic lithospheric viscosity”. In: *Geophysical Journal International* 173.3 (2008), pp. 806–814. DOI: 10.1111/j.1365-246X.2008.03731.x.
- [777] R.J. Leveque. “High-resolution conservative algorithms for advection in incompressible flow”. In: *SIAM J. Numer. Anal.* 33(2) (1996), pp. 627–665. DOI: 10.1137/0733033.
- [778] Randall J LeVeque. *Finite difference methods for ordinary and partial differential equations: steady-state and time-dependent problems*. SIAM, 2007. ISBN: 978-0-898716-29-0.
- [779] Ben Q Li. *Discontinuous finite elements in fluid dynamics and heat transfer*. Springer Science & Business Media, 2006.
- [780] D. Li, M. Gurnis, and G. Stadler. “Towards adjoint-based inversion of time-dependent mantle convection with nonlinear viscosity”. In: *Geophys. J. Int.* 209 (2017), pp. 86–105. DOI: 10.1093/gji/ggw493.
- [781] J. Li, Y. He, and Z. Chen. “Performance of several stabilized finite element methods for the Stokes equations based on the lowest equal-order pairs”. In: *Computing* 86 (2009), pp. 37–51. DOI: 10.1007/s00607-009-0064-5.
- [782] L. Li et al. “Deformation of olivine at mantle pressure using the D-DIA”. In: *Eur. J. Mineral.* 18 (2006), pp. 7–19. DOI: 10.1127/0935-1221/2006/0018-0007.
- [783] S. Li and W.K. Liu. *Meshfree Particle Methods*. Springer, 2004.
- [784] Zhenyu Li and Roger E Khayat. “Finite-amplitude Rayleigh–Bénard convection and pattern selection for viscoelastic fluids”. In: *Journal of Fluid Mechanics* 529 (2005), pp. 221–251. DOI: 10.1017/S0022112005003563.
- [785] Zhenyu Li and Roger E Khayat. “Three-dimensional thermal convection of viscoelastic fluids”. In: *Physical Review E* 71.6 (2005), p. 066305. DOI: 10.1103/PhysRevE.71.066305.

- [786] Jie Liao and Taras Gerya. “Partitioning of crustal shortening during continental collision: 2-D thermomechanical modeling”. In: *Journal of Geophysical Research: Solid Earth* 122.1 (2017), pp. 592–606. DOI: 10.1002/2016JB013398.
- [787] Q. Liao and D. Silvester. “Robust stabilized Stokes approximation methods for highly stretched grids”. In: *IMA Journal of Numerical Analysis* 33 (2013), pp. 413–431. DOI: 10.1093/imanum/drs012.
- [788] A. Limache, S. Idelsohn, R. Rossi, and E. Oñate. “The violation of objectivity in Laplace formulations of the Navier-Stokes equations”. In: *Int. J. Num. Meth. Fluids* 54 (2007), pp. 639–664. DOI: 0.1002/flid.1480.
- [789] S.-C. Lin and P.E. van Keken. “Dynamics of thermochemical plumes: 1. Plume formation and entrainment of a dense layer”. In: *Geochem. Geophys. Geosyst.* 7.2 (2006). DOI: 10.1029/2005GC001071.
- [790] S.-C. Lin and P.E. van Keken. “Dynamics of thermochemical plumes: 2. Complexity of plume structures and its implications for mapping mantle plumes”. In: *Geochem. Geophys. Geosyst.* 7.3 (2006). DOI: 10.1029/2005GC001072.
- [791] Shu-Chuan Lin and Peter E van Keken. “Multiple volcanic episodes of flood basalts caused by thermochemical mantle plumes”. In: *Nature* 436.7048 (2005), pp. 250–252. DOI: 10.1038/nature03697.
- [792] Johannes Linden, Guy Lonsdale, Barbara Steckel, and Klaus Stueben. “Multigrid for the steady-state incompressible navier-stokes equations: A survey”. In: *11th International Conference on Numerical Methods in Fluid Dynamics*. Springer. 1989, pp. 57–68.
- [793] Konstantin D Litasov and Eiji Ohtani. “Effect of water on the phase relations in Earth’s mantle and deep water cycle”. In: *Special Papers-Geological Society of America* 421 (2007), p. 115. DOI: 10.1130/2007.2421(08).
- [794] Benjamin T Liu, Susan J Muller, and Morton M Denn. “Convergence of a regularization method for creeping flow of a Bingham material about a rigid sphere”. In: *Journal of non-newtonian fluid mechanics* 102.2 (2002), pp. 179–191. DOI: 10.1016/S0377-0257(01)00177-X.
- [795] Chengjun Liu, Chengli Huang, and Mian Zhang. “The principal moments of inertia calculated with the hydrostatic equilibrium figure of the Earth”. In: *Geodesy and Geodynamics* 8.3 (2017), pp. 201–205. DOI: 10.1016/j.geog.2017.02.005.
- [796] G.R. Liu. *Mesh Free Methods*. CRC press, 2003.
- [797] G.R. Liu and Y.T. Gu. *An introduction to meshfree methods and their programming*. Springer, 2005.
- [798] G.R. Liu and M.B. Liu. *Smoothed Particle Hydrodynamics*. World Scientific, 2003.
- [799] GR Liu and ZH Tu. “An adaptive procedure based on background cells for meshless methods”. In: *Computer Methods in Applied Mechanics and Engineering* 191.17-18 (2002), pp. 1923–1943. DOI: 10.1016/S0045-7825(01)00360-7.
- [800] P. Ho-Liu, B.H. Hager, and A. Raefsky. “An improved method of Nusselt number calculation”. In: *Geophys. J. R. astr. Soc.* 88 (1987), pp. 205–215. DOI: 10.1111/j.1365-246X.1987.tb01375.x.
- [801] S. Liu and S.D. King. “A benchmark study of incompressible Stokes flow in a 3-D spherical shell using ASPECT”. In: *Geophy. J. Int.* 217 (2019), pp. 650–667. DOI: 10.1093/gji/ggz036.

- [802] X. Liu and S. Zhong. “Analyses of marginal stability, heat transfer and boundary layer properties for thermal convection in a compressible fluid with infinite Prandtl number”. In: *Geophy. J. Int.* 194 (2013), pp. 125–144. DOI: 10.1093/gji/ggt117.
- [803] Xi Liu and Shijie Zhong. “Constraining mantle viscosity structure for a thermochemical mantle using the geoid observation”. In: *Geochemistry, Geophysics, Geosystems* 17.3 (2016), pp. 895–913. DOI: 10.1002/2015GC006161.
- [804] IS Lobanov, I Yu Popov, AI Popov, and TV Gerya. “Numerical approach to the Stokes problem with high contrasts in viscosity”. In: *Applied Mathematics and Computation* 235 (2014), pp. 17–25. DOI: 10.1016/j.amc.2014.02.084.
- [805] J. Lof and A.H. van den Boogaard. “Adaptive return mapping algorithms for  $J_2$  elasto-viscoplastic flow”. In: *Int. J. Num. Meth. Eng.* 51 (2001), pp. 1283–1298. DOI: 10.1002/nme.203.
- [806] Anders Logg, Kent-Andre Mardal, and Garth Wells. *Automated solution of differential equations by the finite element method: The FEniCS book*. Vol. 84. Springer Science & Business Media, 2012.
- [807] C. Loiselet et al. “Subducting slabs: Jellyfishes in the Earth’s mantle”. In: *Geochem. Geophys. Geosyst.* 11.8 (2010). DOI: 10.1029/2010GC003172.
- [808] Benjamin Loret and Jean H Prevost. “Dynamic strain localization in elasto-(visco-) plastic solids, Part 1. General formulation and one-dimensional examples”. In: *Computer Methods in Applied Mechanics and Engineering* 83.3 (1990), pp. 247–273. DOI: 10.1016/0045-7825(90)90073-U.
- [809] Frank Losasso, Ronald Fedkiw, and Stanley Osher. “Spatially adaptive techniques for level set methods and incompressible flow”. In: *Computers & Fluids* 35.10 (2006), pp. 995–1010. DOI: 10.1016/j.compfluid.2005.01.006.
- [810] Aurélie Louis–Napoléon, Thomas Bonometti, Muriel Gerbault, Roland Martin, and Olivier Vanderhaeghe. “Models of convection and segregation in heterogeneous partially molten crustal roots with a VOF method—I: flow regimes”. In: *Geophysical Journal International* 229.3 (2022), pp. 2047–2080. DOI: 10.1093/gji/ggab510.
- [811] Aurélie Louis–Napoléon, Muriel Gerbault, Thomas Bonometti, Cedric Thieulot, Roland Martin, and Olivier Vanderhaeghe. “3D numerical modeling of crustal polydiapirs with Volume-Of-Fluid methods”. In: *Geophysical Journal International* 222 (2020), pp. 474–506. DOI: 10.1093/gji/ggaa141.
- [812] Diogo L Lourenço, Antoine B Rozel, Taras Gerya, and Paul J Tackley. “Efficient cooling of rocky planets by intrusive magmatism”. In: *Nature Geoscience* 11.5 (2018), p. 322.
- [813] Julian P Lowman and Gary T Jarvis. “Continental collisions in wide aspect ratio and high Rayleigh number two-dimensional mantle convection models”. In: *Journal of Geophysical Research: Solid Earth* 101.B11 (1996), pp. 25485–25497. DOI: 10.1029/96JB02568.
- [814] Julian P Lowman and Gary T Jarvis. “Mantle convection flow reversals due to continental collisions”. In: *Geophysical Research Letters* 20.19 (1993), pp. 2087–2090. DOI: 10.1029/93GL02047.
- [815] Julian P Lowman and Gary T Jarvis. “Mantle convection models of continental collision and breakup incorporating finite thickness plates”. In: *Physics of the Earth and Planetary Interiors* 88.1 (1995), pp. 53–68.
- [816] Jacob Lubliner. *Plasticity theory*. Courier Corporation, 2008. ISBN: 978-0-486-46290-5.
- [817] Richard A Lux, Geoffrey F Davies, and John H Thomas. “Moving lithospheric plates and mantle convection”. In: *Geophysical Journal International* 58.1 (1979), pp. 209–228.

- [818] F. Machado, F. Zinani, and S. Frey. “Herschel-Bulkley Fluid Flows Through a Sudden Axisymmetric Expansion via Galerkin Least-Squares Methodology”. In: ().
- [819] P. Machetel and D.A. Yuen. “Penetrative convective flows induced by internal heating and mantle compressibility”. In: *Journal of Geophysical Research* 94.B8 (1989), pp. 10, 609–10, 626. DOI: 10.1029/JB094iB08p10609.
- [820] S.J. Mackwell, M. E. Zimmerman, and D. L. Kohlstedt. “High-temperature deformation of dry diabase with application to tectonics on Venus”. In: *J. Geophys. Res.* 103 (1998), pp. 975–984.
- [821] Christopher W Macosko. “Rheology Principles”. In: *Measurements and Applications* (1994).
- [822] K Mader. “Das Newtonsche Raumpotential prismatischer Körper und seine Ableitungen bis zur dritten Ordnung, Osterr”. In: *Z. Vermess. Sonderheft* 11 (1951).
- [823] M. Maffione, C. Thieulot, D.J.J. van Hinsbergen, A. Morris, O. Plümper, and W. Spakman. “Dynamics of intraoceanic subduction initiation: 1. Oceanic detachment fault inversion and the formation of supra-subduction zone ophiolites”. In: *Geochem. Geophys. Geosyst.* 16 (2015), pp. 1753–1770. DOI: 10.1002/2015GC005746.
- [824] R. Mahmood, N. Kousar, M. Yaqub, and K. Jabeen. “Numerical Simulations of the Square Lid Driven Cavity Flow of Bingham Fluids Using Nonconforming Finite Elements Coupled with a Direct Solver”. In: *Advances in Mathematical Physics* (2017).
- [825] P. Maierová. “Evolution of the Bohemian Massif: Insights from numerical modeling”. PhD thesis. Charles University in Prague, 2012.
- [826] A.V. Malevsky and D.A. Yuen. “Characteristics-based methods applied to infinite Prandtl number thermal convection in the hard turbulent regime”. In: *Physics of Fluids A* 3.9 (1991), pp. 2105–2115. DOI: 10.1063/1.857893.
- [827] A.V. Malevsky and D.A. Yuen. “Plume structures in the hard-turbulent regime of three-dimensional infinite Prandtl number convection”. In: *Geophysical Research Letters* 20.5 (1993), pp. 383–386. DOI: 10.1029/93GL00293.
- [828] A.V. Malevsky and D.A. Yuen. “Strongly chaotic non-newtonian mantle convection”. In: *Geophysical & Astrophysical Fluid Dynamics* 65.1-4 (1992), pp. 149–171. DOI: 10.1080/03091929208225244.
- [829] Andrei V Malevsky, David A Yuen, and LM Weyer. “Viscosity and thermal fields associated with strongly chaotic non-Newtonian thermal convection”. In: *Geophysical research letters* 19.2 (1992), pp. 127–130.
- [830] D.S. Malkus and T.J.R. Hughes. “Mixed finite element methods - reduced and selective integration techniques: a unification of concepts”. In: *Comput. Meth. Appl. Mech. Eng.* 15 (1978), pp. 63–81. DOI: 10.1016/0045-7825(78)90005-1.
- [831] L.E. Malvern. *Introduction to the mechanics of a continuous medium*. Prentice-Hall, Inc., 1969.
- [832] N.S. Mancktelow. “Tectonic pressure: Theoretical concepts and modelled examples”. In: *Lithos* 103 (2008), pp. 149–177.
- [833] Michael Manga. “Low-viscosity mantle blobs are sampled preferentially at regions of surface divergence and stirred rapidly into the mantle”. In: *Physics of the Earth and Planetary Interiors* 180.1-2 (2010), pp. 104–107. DOI: 10.1016/j.pepi.2010.02.013.
- [834] Michael Manga. “Mixing of heterogeneities in the mantle: Effect of viscosity differences”. In: *Geophys. Res. Lett.* 23.4 (1996), pp. 403–406. DOI: 10.1029/96GL00242.



- [835] J.-P. Marcotte. “Methodes iteratives pour la resolution, par Elements Finis, du probleme de Stokes non lineaire”. PhD thesis. Ecole Polytechnique de Montreal, 2000.
- [836] FO Marques and YY Podladchikov. “A thin elastic core can control large-scale patterns of lithosphere shortening”. In: *Earth and Planetary Science Letters* 277.1-2 (2009), pp. 80–85. DOI: 10.1016/j.epsl.2008.10.009.
- [837] Robert S Marshall, Juan C Heinrich, and OC Zienkiewicz. “Natural convection in a square enclosure by a finite-element, penalty function method using primitive fluid variables”. In: *Numerical Heat Transfer, Part B: Fundamentals* 1.3 (1978), pp. 315–330.
- [838] J.G. Masek and C. Duncan. “Minimum-work mountain building”. In: *J. Geophys. Res.* 103.B1 (1998), pp. 907–917. DOI: 10.1029/97JB03213.
- [839] Paolo Massimi, Alfio Quarteroni, and G Scrofani. “An adaptive finite element method for modeling salt diapirism”. In: *Mathematical Models and Methods in Applied Sciences* 16.04 (2006), pp. 587–614. DOI: 10.1142/S0218202506001273.
- [840] A. Massmeyer, E. Di Giuseppe, A. Davaille, T. Rolf, and P.J. Tackley. “Numerical simulation of thermal plumes in a Herschel-Bulkley fluid”. In: *Journal of Non-Newtonian Rheology* 195 (2013), pp. 32–45.
- [841] Takeshi Matsumoto and Yoshiyumi Tomoda. “Numerical simulation of the initiation of subduction at the fracture zone”. In: *Journal of Physics of the Earth* 31.3 (1983), pp. 183–194. DOI: 10.4294/jpe1952.31.183.
- [842] Ctirad Matyska and David A Yuen. “Lower-mantle material properties and convection models of multiscale plumes”. In: *Special Papers – Geological Society of America* 430 (2007), p. 137.
- [843] Ctirad Matyska and David A Yuen. “Profiles of the Bullen parameter from mantle convection modelling”. In: *Earth and Planetary Science Letters* 178.1-2 (2000), pp. 39–46. DOI: 10.1016/S0012-821X(00)00060-1.
- [844] Ctirad Matyska, David A Yuen, Renata M Wentzcovitch, and Hana Čížková. “The impact of variability in the rheological activation parameters on lower-mantle viscosity stratification and its dynamics”. In: *Physics of the Earth and Planetary Interiors* 188.1-2 (2011), pp. 1–8. DOI: 10.1016/j.pepi.2011.05.012.
- [845] D.A. May, J. Brown, and L. Le Pourhiet. “A scalable, matrix-free multigrid preconditioner for finite element discretizations of heterogeneous Stokes flow”. In: *Computer Methods in Applied Mechanics and Engineering* 290 (2015), pp. 496–523. DOI: 10.1016/j.cma.2015.03.014.
- [846] D.A. May and L. Moresi. “Preconditioned iterative methods for Stokes flow problems arising in computational geodynamics”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 33–47. DOI: 10.1016/j.pepi.2008.07.036.
- [847] D.A. May, W.P. Schellart, and L. Moresi. “Overview of adaptive finite element analysis in computational geodynamics”. In: *Journal of Geodynamics* 70 (2013), pp. 1–20. DOI: 10.1016/j.jog.2013.04.002.
- [848] Dave A May, Jed Brown, and Laetitia Le Pourhiet. “pTatin3D: High-performance methods for long-term lithospheric dynamics”. In: *Proceedings of the international conference for high performance computing, networking, storage and analysis*. IEEE Press. 2014, pp. 274–284. DOI: 10.1109/SC.2014.28.
- [849] Neil McBride and Iain Gilmour. *The Solar System: Part 2*. Open University, 2003.
- [850] J.L. McGregor. “Semi-Lagrangian advection of conformal-cubic grids”. In: *Monthly Weather Review* 124 (1996), pp. 1311–1322.

- [851] S. McKee et al. “The MAC method”. In: *Computers and Fluids* 37 (2008), pp. 907–930.
- [852] D.P. McKenzie. “Speculations on the consequences and causes of plate motions”. In: *Geophys. J. R. astr. Soc.* 18 (1969), pp. 1–32. DOI: 10.1111/j.1365-246X.1969.tb00259.x.
- [853] Dan McKenzie. “Finite deformation during fluid flow”. In: *Geophysical Journal International* 58.3 (1979), pp. 689–715. DOI: 10.1111/j.1365-246X.1979.tb04803.x.
- [854] DAN Mckenzie and MJ Bickle. “The volume and composition of melt generated by extension of the lithosphere”. In: *Journal of petrology* 29.3 (1988), pp. 625–679. DOI: 10.1093/petrology/29.3.625.
- [855] Dan McKenzie, Jean Roberts, and Nigel Weiss. “Numerical models of convection in the earth’s mantle”. In: *Tectonophysics* 19.2 (1973), pp. 89–103. DOI: 10.1016/0040-1951(73)90034-6.
- [856] Dan P McKenzie. “Some remarks on heat flow and gravity anomalies”. In: *Journal of Geophysical Research* 72.24 (1967), pp. 6261–6273. DOI: 10.1029/JZ072i024p06261.
- [857] Dan P McKenzie, Jean M Roberts, and Nigel O Weiss. “Convection in the Earth’s mantle: towards a numerical simulation”. In: *Journal of Fluid Mechanics* 62.3 (1974), pp. 465–538. DOI: 10.1017/S0022112074000784.
- [858] Colin P McNally. “Divergence-free interpolation of vector fields from point values - exact  $\nabla \cdot B = 0$  in numerical simulations”. In: *Monthly Notices of the Royal Astronomical Society: Letters* 413.1 (2011), pp. L76–L80.
- [859] A. K. McNamara and S. Zhong. “Thermochemical structures within a spherical mantle: Superplumes or piles?” In: *Journal of Geophysical Research: Solid Earth* 109.B7 (2004). DOI: 10.1029/2003JB002847.
- [860] S. Mei, A.M. Suzuki, D.L. Kohlstedt, N.A. Dixon, and W.B. Durham. “Experimental constraints on the strength of the lithospheric mantle”. In: *J. Geophys. Res.* 115.B08204 (2010). DOI: 10.1029/2009JB006873.
- [861] Mark van der Meijde, Roland Pail, R Bingham, and R Floberghagen. “GOCE data, models, and applications: A review”. In: *International journal of applied earth observation and geoinformation* 35 (2015), pp. 4–15. DOI: 10.1016/j.jag.2013.10.001.
- [862] HJ Melosh and Arthur Raefsky. “The dynamical origin of subduction zone topography”. In: *Geophysical Journal International* 60.3 (1980), pp. 333–354.
- [863] HJ Melosh and CA Williams Jr. “Mechanics of graben formation in crustal rocks: A finite element analysis”. In: *Journal of Geophysical Research: Solid Earth* 94.B10 (1989), pp. 13961–13973.
- [864] C.A. Mériaux, A. May D, J. Mansour, Z. Chen, and O. Kaluza. “Benchmark of three-dimensional numerical models of subduction against a laboratory experiment”. In: *Phys. Earth. Planet. Inter.* 283 (2018), pp. 110–121.
- [865] Guy Metcalfe, Craig R Bina, and JM Ottino. “Kinematic considerations for mantle mixing”. In: *Geophysical research letters* 22.7 (1995), pp. 743–746. DOI: 10.1029/95GL00056.
- [866] AJF Metherell and TJ Quinn. “The gravitational field of a 111 tetrahedron”. In: *Metrologia* 22.2 (1986), p. 87. DOI: 10.1088/0026-1394/22/2/003.
- [867] M.J. van der Meulen, S.J.H. Buiter, J.E. Meulenkamp, and M.J.R. Wortel. “An early Pliocene uplift of the central Apenninic foredeep and its geodynamic significance”. In: *Tectonics* 19 (2000), pp. 300–313. DOI: 10.1029/1999TC900064.
- [868] D.W. Meyer and P. Jenny. “Conservative Velocity Interpolation for PDF Methods”. In: *Proc. Appl. Math. Mech.* 4 (2004), pp. 466–467.

- [869] Victor Benno Meyer-Rochow and Jozsef Gal. “Pressures produced when penguins pooh - calculations on avian defaecation”. In: *Polar Biology* 27.1 (2003), pp. 56–58.
- [870] R. L. Michalowski. “Upper-bound load estimates on square and rectangular footings”. In: *Géotechnique* 51.9 (2001), pp. 787–798.
- [871] K. Michibayashi and D. Mainprice. “The role of pre-existing mechanical anisotropy on shear zone development within oceanic mantle lithosphere: an example from the Oman ophiolite”. In: *J. Petrol.* 45(2) (2004), pp. 405–414.
- [872] Laurent Michon and Olivier Merle. “Crustal structures of the Rhinegraben and the Massif Central grabens: An experimental approach”. In: *Tectonics* 19.5 (2000), pp. 896–904. DOI: 10.1029/2000TC900015.
- [873] Chohong Min and Frédéric Gibou. “A second order accurate level set method on non-graded adaptive cartesian grids”. In: *Journal of Computational Physics* 225.1 (2007), pp. 300–321.
- [874] P. Ming, Z.-c. Shi, and Y. Xu. “A new superconvergence property of nonconforming rotated  $Q_1$  element in 3D”. In: *Computer Methods in Applied Mechanics and Engineering* 197 (2007), pp. 95–102.
- [875] P. Ming, Z.-c. Shi, and Y. Xu. “Superconvergence studies of quadrilateral nonconforming rotated  $Q_1$  elements”. In: *International Journal of Numerical analysis and modeling* 3.3 (2006), pp. 322–332.
- [876] R. von Mises. “Mechanik der festen Körper im plastisch deformablen Zustand”. In: *Nachrichten der Königlichen Gesellschaft der Wissenschaften* (1913), p. 582.
- [877] Y. Mishin. “Adaptive multiresolution methods for problems of computational geodynamics”. PhD thesis. ETH Zurich, 2011.
- [878] Y. Mishin. “Adaptive multiresolution methods for problems of computational geodynamics”. PhD thesis. ETH Zürich, 2011. DOI: 10.3929/ethz-a-007347901.
- [879] Y.A. Mishin, O.V. Vasilyev, and T.V. Gerya. “A Wavelet-Based Adaptive Finite Element Method for the Stokes Problems”. In: *Fluids* 7 (2022), p. 221. DOI: 10.3390/fluids7070221.
- [880] J.X. Mitrovica and A.M. Forte. “A new inference of mantle viscosity based upon joint inversion of convection and glacial isostatic adjustment data”. In: *Earth and Planetary Science Letters* 225.1-2 (2004), pp. 177–189. DOI: 10.1016/j.epsl.2004.06.005.
- [881] J.X. Mitrovica, R.N. Pysklywec, C. Beaumont, and A. Rutt. “The Devonian to Permian sedimentation of the Russian platform: An example of subduction-controlled long-wavelength tilting of continents”. In: *Journal of Geodynamics* 22.1-2 (1996), pp. 79–96. DOI: 10.1016/0264-3707(96)00008-7.
- [882] Jerry X Mitrovica. “Haskell [1935] revisited”. In: *Journal of Geophysical Research: Solid Earth* 101.B1 (1996), pp. 555–569. DOI: 10.1029/95JB03208.
- [883] E. Mitsoulis and S. Galazoulas. “Simulation of viscoplastic flow past cylinders in tubes”. In: *Journal of Non-Newtonian Fluid Mechanics* 158 (2009), pp. 132–141. DOI: 10.1016/j.jnnfm.2008.10.006.
- [884] E. Mitsoulis and Th. Zisis. “Flow of Bingham plastics in a lid-driven square cavity”. In: *Journal of Non-Newtonian Fluid Mechanics* 101 (2001), pp. 173–180.
- [885] Arata Miyauchi and Masanori Kameyama. “Influences of the depth-dependence of thermal conductivity and expansivity on thermal convection with temperature-dependent viscosity”. In: *Physics of the Earth and Planetary Interiors* 223 (2013), pp. 86–95.
- [886] A. Mizukami. “A mixed Finite Element method for boundary flux computation”. In: *Computer Methods in Applied Mechanics and Engineering* 57 (1986), pp. 239–243.

- [887] Arash Mohajeri, Yaron Finzi, Hans Mühlhaus, and Gideon Rosenbaum. “Melt and shear interactions in the lithosphere: Theory and numerical analysis of pure shear extension”. In: *Journal of Geophysical Research: Solid Earth* 118.5 (2013), pp. 2488–2499. DOI: 10.1111/j.1365-246X.2006.03225.x.
- [888] P.K. Mohapatra, V. Eswaran, and S. Murty Bhallamudi. “Two-dimensional analysis of dam-break flow in vertical plane”. In: *Journal of Hydraulic Engineering* 125.2 (1999), pp. 183–192.
- [889] P. Molnar. *Brace-Goetze strength profiles, the partitioning of strike-slip and thrust faulting at zones of oblique convergence, and the stress-heat flow paradox of the San Andreas Fault*. Academic Press Ltd, 1992.
- [890] P. Molnar and P. Tapponnier. “Relation of the tectonics of eastern China to the India-Eurasia collision: Application of the slip-line field theory to large-scale continental tectonics”. In: *Geology* 5 (1977), pp. 212–216.
- [891] JJ Monaghan. “Particle methods for hydrodynamics”. In: *Computer Physics Reports* 3.2 (1985), pp. 71–124. DOI: 10.1016/0167-7977(85)90010-3.
- [892] P. Mons and G. Rogé. “L’élément  $Q_1$ -bulle/ $Q_1$ ”. In: *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique* 26.4 (1992), pp. 507–521. DOI: xxxx.
- [893] A. Montlaur, S. Fernandez-Mendez, and A. Huerta. “Discontinuous Galerkin methods for the Stokes equations using divergence-free approximations”. In: *Int. J. Num. Meth. Fluids* 57.08 (2008), pp. 1071–1092.
- [894] A. Montlaur, S. Fernandez-Mendez, J. Peraire, and A. Huerta. “Discontinuous Galerkin methods for the Navier-Stokes equations using solenoidal approximations”. In: *Int. J. Num. Meth. Fluids* 64 (2010), pp. 549–564. DOI: 10.1002/flid.2161.
- [895] L Moresi and H-B Mühlhaus. “Anisotropic viscous models of large-deformation Mohr–Coulomb failure”. In: *Philosophical Magazine* 86.21-22 (2006), pp. 3287–3305. DOI: 10.1080/1478643050025541
- [896] L Moresi, H-B Mühlhaus, Vincent Lemiale, and D May. “Incompressible viscous formulations for deformation and yielding of the lithosphere”. In: *Geological Society, London, Special Publications* 282.1 (2007), pp. 457–472. DOI: 10.1144/SP282.19.
- [897] L-N Moresi and A. Lenardic. “Three-dimensional numerical simulations of crustal deformation and subcontinental mantle convection”. In: *Earth and Planetary Science Letters* 150.3-4 (1997), pp. 233–243. DOI: 10.1016/S0012-821X(97)00093-9.
- [898] L. Moresi, F. Dufour, and H.B. Mühlhaus. “A Lagrangian integration point finite element method for large deformation modeling of visco-elastic geomaterials”. In: *J. Comp. Phys.* 184.2 (2003), pp. 476–497. DOI: 10.1016/S0021-9991(02)00031-1.
- [899] L. Moresi, F. Dufour, and H.B. Mühlhaus. “Mantle Convection Modeling with Viscoelastic/Brittle Lithosphere: Numerical Methodology and Plate Tectonic Modeling”. In: *Pure and Applied Geophysics* 159 (2002), p. 159. DOI: 10.1007/s00024-002-8738-3.
- [900] L. Moresi and M. Gurnis. “Constraints on the lateral strength of slabs from three-dimensional dynamic flow models”. In: *Earth and Planetary Science Letters* 138.1 (1996), pp. 15–28. DOI: 10.1016/0012-821X(95)00221-W.
- [901] L. Moresi, S. Quenette, V. Lemiale, C. Mériaux, B. Appelbe, and H.-B. Mühlhaus. “Computational approaches to studying non-linear dynamics of the crust and mantle”. In: *Phys. Earth. Planet. Inter.* 163 (2007), pp. 69–82. DOI: 10.1016/j.pepi.2007.06.009.

- [902] L. Moresi, S. Zhong, and M. Gurnis. “The accuracy of finite element solutions of Stokes’ flow with strongly varying viscosity”. In: *Physics of the Earth and Planetary Interiors* 97.1-4 (1996), pp. 83–94. DOI: 10.1016/0031-9201(96)03163-9.
- [903] L.-N. Moresi and V.S. Solomatov. “Numerical investigation of 2D convection with extremely large viscosity variations”. In: *Physics of Fluids* 7.9 (1995), pp. 2154–2162.
- [904] Louis Moresi and Ben Mather. “Stripy: A Python module for (constrained) triangulation in Cartesian coordinates and on a sphere”. In: *Journal of Open Source Software* 4.38 (2019), p. 1410. DOI: 10.21105/joss.01410.
- [905] Isabelle Moretti and Claude Froidevaux. “Thermomechanical models of active rifting”. In: *Tectonics* 5.4 (1986), pp. 501–511.
- [906] Jason P Morgan, Jorge M Taramón, and Jörg Hasenclever. “Shape-preserving finite elements in cylindrical and spherical geometries: The double Jacobian approach”. In: *International Journal for Numerical Methods in Fluids* 92.6 (2020), pp. 635–668. DOI: 10.1002/flid.4799.
- [907] Gabriele Morra. “Pythonic Geodynamics”. In: *Lecture Notes in Earth System Sciences* (2018).
- [908] K.W. Morton and D.F. Mayers. *Numerical Solution of Partial Differential Equations: An Introduction*. Cambridge, 2005.
- [909] Jiří Moser, Ctirad Matyska, David A Yuen, Andrei V Malevsky, and Helmut Harder. “Mantle rheology, convection and rotational dynamics”. In: *Physics of the earth and planetary interiors* 79.3-4 (1993), pp. 367–381. DOI: 10.1016/0031-9201(93)90115-P.
- [910] P. Moulik and G. Ekström. “The relationships between large-scale variations in shear velocity, density, and compressional velocity in the Earth’s mantle”. In: *J. Geophys. Res.* 121 (2016). DOI: 10.1002/2015JB012679.
- [911] Lin Mu and Xiu Ye. “A simple finite element method for the Stokes equations”. In: *Advances in Computational Mathematics* 43.6 (2017), pp. 1305–1324. DOI: 10.1007/s10444-017-9526-z.
- [912] H-B Mühlhaus, L Moresi, and Miroslav Cada. “Emergent anisotropy and flow alignment in viscous rock”. In: *pure and applied geophysics* 161.11-12 (2004), pp. 2451–2463.
- [913] H-B Mühlhaus, L Moresi, and M Čada. “Anisotropy model for mantle convection”. In: *Computational Fluid and Solid Mechanics 2003*. 2003, pp. 1044–1046. DOI: 10.1016/B978-008044046-0.50255-4.
- [914] H.-B. Mühlhaus, L. Moresi, B. Hobbs, and F. Dufour. “Large amplitude folding in finely layered viscoelastic rock structures”. In: *Pure appl. Geophys.* 159 (2002), pp. 2311–2333.
- [915] H.B. Mühlhaus and K. Regenauer-Lieb. “Towards a self-consistent plate mantle model that includes elasticity: simple benchmarks and application to basic modes of convection”. In: *Geophy. J. Int.* 163 (2005), pp. 788–800. DOI: 10.1111/j.1365-246X.2005.02742.x.
- [916] Hans Mühlhaus, Louis Moresi, Lutz Gross, and Joseph Grotowski. “The influence of non-coaxiality on shear banding in viscous-plastic materials”. In: *Granular Matter* 12.3 (2010), pp. 229–238. DOI: 10.1007/s10035-010-0176-9.
- [917] Hans B Mühlhaus, Jingyu Shi, Louise Olsen-Kettle, and Louis Moresi. “Effects of a non-coaxial flow rule on shear bands in viscous-plastic materials”. In: *Granular Matter* 13.3 (2011), pp. 205–210.
- [918] E. Mulyukova, B. Steinberger, M. Dabrowski, and S.V. Sobolev. “Survival of LLSVPs for billions of years in a vigorously convecting mantle: Replenishment and destruction of chemical anomaly”. In: *J. Geophys. Res.* 120 (2015), pp. 3824–3847. DOI: 10.1002/2014JB011688.

- [919] Elvira Mulyukova. “Stability of the large low shear velocity provinces: numerical modeling of thermochemical mantle convection”. PhD thesis. Universität Potsdam, 2014.
- [920] K Nafa and RW Thatcher. “Low-order macroelements for two-and three-dimensional Stokes flow”. In: *Numerical Methods for Partial Differential Equations* 9.5 (1993), pp. 579–591. DOI: 10.1002/num.1690090507.
- [921] D. Nagy, G. Papp, and J. Benedek. “Corrections to ”The gravitational potential and its derivatives for the prism””. In: *Journal of Geodesy* 76.8 (2002), p. 475. DOI: xxxx.
- [922] D. Nagy, G. Papp, and J. Benedek. “The gravitational potential and its derivatives for the prism”. In: *Journal of Geodesy* 74.7–8 (2000), pp. 552–560. DOI: 10.1007/s001900000116.
- [923] Tomoeiki Nakakuki and Erika Mura. “Dynamics of slab rollback and induced back-arc basin formation”. In: *Earth Planet. Sci. Lett.* 361 (2013), pp. 287–297. DOI: 10.1016/j.epsl.2012.10.031.
- [924] Tomoeiki Nakakuki, Hiroki Sato, and Hiromi Fujimoto. “Interaction of the upwelling plume with the phase and chemical boundary at the 670 km discontinuity: Effects of temperature-dependent viscosity”. In: *Earth Planet. Sci. Lett.* 121.3–4 (1994), pp. 369–384. DOI: 10.1016/0012-821X(94)90078-7.
- [925] SM Nakiboglu and Kurt Lambeck. “A study of the Earth’s response to surface loading with application to Lake Bonneville”. In: *Geophysical Journal International* 70.3 (1982), pp. 577–620. DOI: 10.1111/j.1365-246X.1982.tb05975.x.
- [926] J. B. Naliboff and L. H. Kellogg. “Can large increases in viscosity and thermal conductivity preserve large-scale heterogeneity in the mantle?” In: *Physics of the Earth and Planetary Interiors* 161.1-2 (2007), pp. 86–102. DOI: 10.1016/j.pepi.2007.01.009.
- [927] J.B. Naliboff, S.J.H. Buiter, G. Péron-Pinvidic, P.T. Osmundsen, and J. Tetreault. “Complex fault interaction controls continental rifting”. In: *Nature Communications* 8 (2017), p. 1179. DOI: 10.1038/s41467-017-00904-x.
- [928] J.B. Naliboff, C. Lithgow-Bertelloni, L.J. Ruff, and N. de Koker. “The effects of lithospheric thickness and density structure on Earth’s stress field”. In: *Geophy. J. Int.* 188 (2012), pp. 1–17. DOI: 10.1111/j.1365-246X.2011.05248.x.
- [929] Michele Napolitano, Giuseppe Pascazio, and Luigi Quartapelle. “A review of vorticity conditions in the numerical solution of the  $\zeta$ – $\psi$  equations”. In: *Computers & Fluids* 28.2 (1999), pp. 139–185. DOI: 10.1016/S0045-7930(98)00024-3.
- [930] A. Napov and Y. Notay. “An algebraic multigrid method with guaranteed convergence rate”. In: *SIAM J. Sci. Comput.* 34 (2012), A1079–A1109. DOI: 10.1137/100818509.
- [931] Mauricio Nava-Flores et al. “3D gravity modeling of complex salt features in the southern gulf of Mexico”. In: *International Journal of Geophysics* 2016 (2016). DOI: 10.1155/2016/1702164.
- [932] A.M. Negredo, R. Sabadini, G. Bianco, and M. Fernandez. “Three-dimensional modelling of crustal motions caused by subduction and continental convergence in the central Mediterranean”. In: *Geophy. J. Int.* 136 (1999), pp. 261–274. DOI: 10.1046/j.1365-246X.1999.00726.x.
- [933] M. Nettesheim, T.A. Ehlers, D.M. Whipp, and A. Koptev. “The influence of upper-plate advance and erosion on overriding plate deformation in orogen syntaxes”. In: *Solid Earth* 9 (2018), pp. 1207–1224. DOI: 10.5194/se-9-1207-2018.
- [934] Derek Neuharth and Eric Mittelstaedt. “Temporal variations in plume flux: characterizing pulsations from tilted plume conduits in a rheologically complex mantle”. In: *Geophysical Journal International* 233.1 (2023), pp. 338–358. DOI: 10.1093/gji/ggac455.

- [935] William I Newman. *Continuum Mechanics in the Earth Sciences*. Cambridge University Press, 2012. ISBN: 978-0-521-56289-8.
- [936] N.C. Nguyen and J. Peraire. “Hybridizable discontinuous Galerkin methods for partial differential equations in continuum mechanics”. In: *J. Comp. Phys.* 231 (2012), pp. 5955–5988. DOI: 10.1016/j.jcp.2012.02.033.
- [937] N.C. Nguyen, J. Peraire, and B. Cockburn. “A hybridizable discontinuous Galerkin method for Stokes flow”. In: *Comput. Methods Appl. Mech. Engrg.* 199 (2010), pp. 582–597. DOI: 10.1016/j.cma.2009.10.007.
- [938] N.C. Nguyen, J. Peraire, and B. Cockburn. “An implicit high-order hybridizable discontinuous Galerkin method for the incompressible Navier-Stokes equations”. In: *J. Comp. Phys.* 230 (2011), pp. 1147–1170. DOI: 10.1016/j.jcp.2010.10.032.
- [939] G. Nicolas and T. Fouquet. “Adaptive mesh refinement for conformal hexahedral meshes”. In: *Finite Elements in Analysis and Design* 67 (2013), pp. 1–12. DOI: 10.1016/j.finel.2012.11.008.
- [940] G. Nicolas and T. Fouquet. “Conformal hexahedral meshes and adaptive mesh refinement”. In: *VI International Conference on Adaptive Modeling and Simulation* (2013). DOI: xxxx.
- [941] A.C. de Niet and F.W. Wubs. “Two preconditioners for the Saddle Point equation”. In: *European Congress on Computational Methods in Applied Sciences and Engineering ECCOMAS 2004* (2004).
- [942] Nicolai Nijholt and Rob Govers. “The role of passive margins on the evolution of Subduction-Transform Edge Propagators (STEPs)”. In: *Journal Of Geophysical Research* 120.10 (2015), pp. 7203–7230. DOI: 10.1002/2015JB012202.
- [943] Lena Noack and Doris Breuer. “First-and second-order Frank-Kamenetskii approximation applied to temperature-, pressure-and stress-dependent rheology”. In: *Geophysical Journal International* 195.1 (2013), pp. 27–46. DOI: 10.1093/gji/ggt248.
- [944] S. Norburn and D. Silvester. “Fourier analysis of stabilized Q1-Q1 mixed finite element approximation”. In: *SIAM J. Numer. Anal.* 39 (2001), pp. 817–833. DOI: 10.1137/S0036142999362274.
- [945] S. Norburn and D. Silvester. “Stabilised vs. stable mixed methods for incompressible flow”. In: *Computer Methods in Applied Mechanics and Engineering* 166 (1998), pp. 131–141. DOI: 10.1016/S0045-7825(98)00087-5.
- [946] Y. Notay. “Aggregation-based algebraic multigrid for convection-diffusion equations”. In: *SIAM J. Sci. Comput.* 34 (2012), A2288–A2316. DOI: 10.1137/110835347.
- [947] JL Nowinski. “On the Three-Dimensional Cerruti Problem for an Elastic Nonlocal Half-Space”. In: *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik* 72.7 (1992), pp. 243–249. DOI: 10.1002/zamm.19920720702.
- [948] K.A. O’Farrell, J.P. Lowman, and H.-P. Bunge. “Comparison of spherical-shell and plane-layer mantle convection thermal structure in viscously stratified models with mixed-mode heating: Implications for the incorporation of temperature-dependent parameters”. In: *Geophysical Journal International* 192.2 (2013), pp. 456–472. DOI: 10.1093/gji/ggs053.
- [949] C. O’Neill, L. Moresi, D. Müller, R. Albert, and F. Dufour. “Ellipsis 3D: a particle-in-cell finite element hybrid code for modelling mantle convection and lithospheric deformation”. In: *Computers and Geosciences* 32 (2006), pp. 1769–1779. DOI: 10.1016/j.cageo.2006.04.006.
- [950] J Tinsley Oden and Olivier-Pierre Jacquotte. “Stability of some mixed finite element methods for Stokesian flows”. In: *Computer methods in applied mechanics and engineering* 43.2 (1984), pp. 231–247.

- [951] J Tinsley Oden, Noboru Kikuchi, and Young Joon Song. “Penalty-finite element methods for the analysis of Stokesian flows”. In: *Computer Methods in Applied Mechanics and Engineering* 31.3 (1982), pp. 297–329.
- [952] J. Oeser, H.-P. Bunge, M. Mohr, and H. Igel. “Frontiers in computational geophysics: Simulations of mantle circulation, plate tectonics and seismic wave propagation”. In: *Notes on Numerical Fluid Mechanics and Multidisciplinary Design* 100 (2009), pp. 387–397. DOI: 10.1007/978-3-540-70805-6\_30.
- [953] Masaki Ogawa. “A Numerical Model of a Coupled Magmatism-Mantle Convection System in Venus and the Earth’s Mantle beneath Archean Continental Crusts”. In: *Icarus* 102.1 (1993), pp. 40–61.
- [954] Masaki Ogawa, Gerald Schubert, and Abdelfattah Zebib. “Numerical simulations of three-dimensional thermal convection in a fluid with strongly temperature-dependent viscosity”. In: *Journal of fluid mechanics* 233 (1991), pp. 299–328.
- [955] D. Olbertz, M.J.R. Wortel, and U. Hansen. “Trench migration and subduction zone geometry”. In: *Geophysical Research Letters* 24.3 (1997), pp. 221–224. DOI: 10.1029/96GL03971.
- [956] D. Oldham, J.H. Davies, and T.N. Phillips. “Generic polyhedron grid generation for solving partial differential equations on spherical surfaces”. In: *Computers and Geosciences* 39 (2012), pp. 11–17. DOI: 10.1016/j.cageo.2011.06.004.
- [957] Peter Olson, David A Yuen, and Derick Balsiger. “Mixing of passive heterogeneities by mantle convection”. In: *Journal of Geophysical Research: Solid Earth* 89.B1 (1984), pp. 425–436. DOI: 10.1029/JB089iB01p00425.
- [958] Peter Olson, David A. Yuen, and Derick Balsiger. “Convective mixing and the fine structure of mantle heterogeneity”. In: *Phys. Earth. Planet. Inter.* 36.3–4 (1984), pp. 291–304. DOI: 10.1016/0031-9201(84)90053-0.
- [959] E. Olsson and G. Kreiss. “A conservative level set method for two phase flow”. In: *J. Comp. Phys.* 210 (2005), pp. 225–246.
- [960] E. Olsson, G. Kreiss, and S. Zahedi. “A conservative level set method for two phase flow II”. In: *J. Comp. Phys.* 225 (2007), pp. 785–807.
- [961] S. Osher and R. Fedkiw. “Level set methods: an overview and some recent results”. In: *J. Comp. Phys.* 169 (2001), pp. 463–502.
- [962] S. Osher and C.-W. Shu. “High-order essentially non-oscillatory schemes for Hamilton-Jacobi equations”. In: *SIAM J. Numer. Anal* 28 (1991), pp. 907–922.
- [963] Tim A Osswald and Natalie Rudolph. *Polymer rheology: fundamentals and applications*. Carl Hanser Verlag GmbH Co KG, 2014.
- [964] J.M. Ottino. *The kinematics of mixing: stretching, chaos, and transport*. Vol. 3. Cambridge university press, 1989.
- [965] A. Ouazzi. “Finite Element Simulation of Nonlinear Fluids with Application to Granular Material and Powder”. PhD thesis. Universität Dortmund, 2005.
- [966] Guillaume Ovarlez and Sarah Hormozi. *Lectures on visco-plastic fluid mechanics*. Springer, 2019. ISBN: 978-3-319-89437-9.
- [967] D.R.J. Owen and E. Hinton. *Finite Elements in Plasticity*. Pineridge press, 1980.
- [968] M. OzBench et al. “A model comparison study of large-scale mantle-lithosphere dynamics driven by subduction”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 224–234. DOI: 10.1016/j.pepi.2008.08.011.



- [969] Karen Paczkowski, Laurent GJ Montési, Maureen D Long, and Christopher J Thissen. “Three-dimensional flow in the subslab mantle”. In: *Geochemistry, Geophysics, Geosystems* 15.10 (2014), pp. 3989–4008. DOI: 10.1002/2014GC005441.
- [970] C.C. Paige and M.A. Saunders. “Solution of sparse indefinite systems of linear equations”. In: *SIAM J. Numer. Anal* 12.4 (1975), pp. 617–629. DOI: 10.1137/0712047.
- [971] Enok Palm, Jan Erik Weber, and Oddmund Kvernfold. “On steady convection in a porous medium”. In: *Journal of Fluid Mechanics* 54.1 (1972), pp. 153–161. DOI: 10.1017/S0022112072000592.
- [972] Isabelle Panet, Gwendoline Pajot-Métivier, Marianne Greff-Lefftz, Laurent Métivier, Michel Diamant, and Mioara Manda. “Mapping the mass distribution of Earth’s mantle using satellite-derived gravity gradients”. In: *Nature Geoscience* 7.2 (2014), pp. 131–135. DOI: 10.1038/NGEO2063.
- [973] Tasos C Papanastasiou. “Flows of materials with yield”. In: *Journal of Rheology* 31.5 (1987), pp. 385–404.
- [974] Chunjae Park and Dongwoo Sheen. “ $P_1$ -nonconforming quadrilateral finite element methods for second-order elliptic problems”. In: *SIAM Journal on Numerical Analysis* 41.2 (2003), pp. 624–640. DOI: 10.1137/S0036142902404923.
- [975] Chunjae Park, Dongwoo Sheen, and Byeong-Chun Shin. “A subspace of the DSSY nonconforming quadrilateral finite element space for the Stokes equations”. In: *Journal of Computational and Applied Mathematics* 239 (2013), pp. 220–230. DOI: 10.1016/j.cam.2012.09.042.
- [976] EM Parmentier, C Sotin, and BJ Travis. “Turbulent 3-D thermal convection in an infinite Prandtl number, volumetrically heated fluid: implications for mantle dynamics”. In: *Geophysical Journal International* 116.2 (1994), pp. 241–251. DOI: 10.1111/j.1365-246X.1994.tb01795.x.
- [977] EM Parmentier, DL Turcotte, and KE Torrance. “Numerical experiments on the structure of mantle plumes”. In: *Journal of Geophysical Research* 80.32 (1975), pp. 4417–4424. DOI: 10.1029/JB080i032p04417.
- [978] EM Parmentier, DL Turcotte, and KE Torrance. “Studies of finite amplitude non-Newtonian thermal convection with application to convection in the Earth’s mantle”. In: *Journal of Geophysical Research* 81.11 (1976), pp. 1839–1846. DOI: 10.1029/JB081i011p01839.
- [979] M.H. Parrish. “A SELECTIVE APPROACH TO CONFORMAL REFINEMENT OF UNSTRUCTURED HEXAHEDRAL MESHES”. PhD thesis. Department of Civil and Environmental Engineering, Brigham Young University, 2007.
- [980] Barry Parsons and Stephen Daly. “The relationship between surface topography, gravity anomalies, and temperature structure of convection”. In: *Journal of Geophysical Research: Solid Earth* 88.B2 (1983), pp. 1129–1144.
- [981] S.V. Patankar and D.B. Spalding. “A calculation procedure for heat, mass and momentum transfer in three-dimensional parabolic flows”. In: *Int. J. Heat Mass Transfer* 15 (1972), pp. 1787–1806. DOI: 10.1016/B978-0-08-030937-8.50013-1.
- [982] WSB Paterson and WF Budd. “Flow parameters for ice sheet modeling”. In: *Cold Regions Science and Technology* 6.2 (1982), pp. 175–177. DOI: 10.1016/0165-232X(82)90010-6.
- [983] Archie Paulson, Shijie Zhong, and John Wahr. “Inference of mantle viscosity from GRACE and relative sea level data”. In: *Geophysical Journal International* 171.2 (2007), pp. 497–508. DOI: 10.1111/j.1365-246X.2007.03556.x.

- [984] D. Pelletier, A. Fortin, and R. Camarero. “Are FEM solutions of incompressible flows really incompressible ? (Or how simple flows can cause headaches!)” In: *International journal for numerical methods in fluids* 9 (1989), pp. 99–112. DOI: 10.1002/flid.1650090108.
- [985] Jon Pelletier. *Quantitative modelling of Earth surface processes*. Cambridge University Press, 2008.
- [986] WR Peltier. “Penetrative convection in the planetary mantle”. In: *Geophysical Fluid Dynamics* 5.1 (1972), pp. 47–88. DOI: 10.1080/03091927308236108.
- [987] WR Peltier. “The impulse response of a Maxwell Earth”. In: *Reviews of Geophysics* 12.4 (1974), pp. 649–669. DOI: 10.1029/RG012i004p00649.
- [988] WR Peltier, S Butler, and LP Solheim. “The influence of phase transformations on mantle mixing and plate tectonics”. In: *Earth’s Deep Interior. Gordon & Breach, Amsterdam* (1997), pp. 405–430.
- [989] WR Peltier and GT Jarvis. “Whole mantle convection and the thermal evolution of the Earth”. In: *Physics of the Earth and Planetary Interiors* 29.3-4 (1982), pp. 281–304. DOI: 10.1016/0031-9201(82)90018-8.
- [990] G. Peltzer and P. Tapponnier. “Formation and evolution of strike-slip faults, rifts, and basins during the india-asia collision: an experimental approach”. In: *J. Geophys. Res.* 93.B12 (1988), pp. 15085–15177. DOI: 10.1029/JB093iB12p15085.
- [991] D. Perić and C. Huang. “Analytical solutions for the three-invariant cam clay model: drained and undrained loading”. In: *16th ASCE Engineering Mechanics Conference* (2003).
- [992] J. Perry-Houts and L. Karlstrom. “Anisotropic viscosity and time-evolving lithospheric instabilities due to aligned igneous intrusions”. In: *Geophysical Journal International* 216.2 (2018), pp. 794–802. DOI: 10.1093/gji/ggy466.
- [993] P. Perzyna. “Constitutive Modelling of Dissipative Solids for Localization and Fracture”. In: *Localization and Fracture Phenomena in Inelastic Solids*. Ed. by Piotr Perzyna. Vienna: Springer Vienna, 1998, pp. 99–241. ISBN: 978-3-7091-2528-1. DOI: 10.1007/978-3-7091-2528-1\_3.
- [994] Piotr Perzyna. “Fundamental problems in viscoplasticity”. In: *Advances in applied mechanics*. Vol. 9. 1966, pp. 243–377. DOI: 10.1016/S0065-2156(08)70009-7.
- [995] A.G. Petrunin and S.V. Sobolev. “Three-dimensional numerical models of the evolution of pull-apart basins”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 387–399. DOI: 10.1016/j.pepi.2008.08.017.
- [996] Tim N Phillips, J Huw Davies, and D Oldham. “Towards global SEM mantle convection simulations on polyhedral-based grids”. In: *Journal of Computational and Applied Mathematics* 348 (2019), pp. 48–57. DOI: 10.1016/j.cam.2018.08.042.
- [997] E. Pichelin and T. Coupez. “Finite element solution of the 3D mold filling problem for viscous incompressible fluid”. In: *Computer Methods in Applied Mechanics and Engineering* 163 (1998), pp. 359–371. DOI: 10.1016/S0045-7825(98)00024-3.
- [998] RT Pierrehumbert. “Large-scale horizontal mixing in planetary atmospheres”. In: *Physics of Fluids A: Fluid Dynamics* 3.5 (1991), pp. 1250–1260. DOI: 10.1063/1.858053.
- [999] D.A. di Pietro and A. Ern. *Mathematical Aspects of Discontinuous Galerkin Methods*. Springer, 2012.
- [1000] D.A. Di Pietro, S. Lo Forte, and N. Parolini. “Mass preserving finite element implementations of the level set method”. In: *Applied Numerical Mathematics* 56 (2006), pp. 1179–1195. DOI: 10.1016/j.apnum.2006.03.003.

- [1001] J. Pitkäranta and T. Saarinen. “A Multigrid Version of a Simple Finite Element Method for the Stokes Problem”. In: *Mathematics of Computation* 45.171 (1985), pp. 1–14.
- [1002] A.-C. Plesa. “Mantle Convection in a 2D Spherical Shell”. In: *NFOCOMP 2011 : The First International Conference on Advanced Communications and Computation* (2011).
- [1003] Ana-Catalina Plesa, Nicola Tosi, and Christian Hüttig. “Thermo-chemical convection in planetary mantles: advection methods and magma ocean overturn simulations”. In: *Integrated Information and Computing Systems for Natural, Spatial, and Social Sciences*. IGI Global, 2013, pp. 302–323. DOI: 10.4018/978-1-4666-2190-9.ch015.
- [1004] Donald Plouff. “Gravity and magnetic fields of polygonal prisms and application to magnetic terrain corrections”. In: *Geophysics* 41.4 (1976), pp. 727–741. DOI: 10.1190/1.1440645.
- [1005] Meredith Plumley and Keith Julien. “Scaling laws in Rayleigh-Benard convection”. In: *Earth and Space Science* 6.9 (2019), pp. 1580–1592. DOI: 10.1029/2019EA000583.
- [1006] Jean-Paul Poirier. *Introduction to the Physics of the Earth’s Interior*. Cambridge University Press, 2000.
- [1007] A. Poliakov, P. Cundall, P. Podlachikov, and V. Lyakhovsky. “An explicit inertial method for the simulation of viscoelastic flow: an evaluation of elastic effects on diapiric flow in two- and three-layers models”. In: *Flow and creep in the solar system: Observations, Modeling and theory*. Kluwer Academic Publishers, 1993, pp. 175–195.
- [1008] A. Poliakov and Y. Podlachikov. “Diapirism and topography”. In: *Geophy. J. Int.* 109 (1992), pp. 553–564. DOI: 10.1111/j.1365-246X.1992.tb00117.x.
- [1009] A.N.B. Poliakov and H.J. Herrmann. “Self-organized criticality of plastic shear bands in rocks”. In: *Geophys. Res. Lett.* 21.19 (1994), pp. 2143–2146. DOI: 10.1029/94GL02005.
- [1010] J-P Ponthot and Ted Belytschko. “Arbitrary Lagrangian-Eulerian formulation for element-free Galerkin method”. In: *Computer methods in applied mechanics and engineering* 152.1-2 (1998), pp. 19–46. DOI: 10.1016/S0045-7825(97)00180-1.
- [1011] A.A. Popov and S.V. Sobolev. “SLIM3D: a tool for three-dimensional thermomechanical modelling of lithospheric deformation with elasto-visco-plastic rheology”. In: *Phys. Earth. Planet. Inter.* 171.1 (2008), pp. 55–75. DOI: 10.1016/j.pepi.2008.03.007.
- [1012] AI Popov, IS Lobanov, I Yu Popov, and TV Gerya. “Benchmark solutions for nanoflows”. In: *NANOSYSTEMS: PHYSICS, CHEMISTRY, MATHEMATICS* 5.3 (2014). DOI: xxxx.
- [1013] I.Yu. Popov, I.S. Lobanov, S.I. Popov, A.I. Popov, and T. Gerya. “Practical analytical solutions for benchmarking of 2D and 3D geodynamic Stokes problems with variable viscosity”. In: *Solid Earth* 5 (2014), pp. 461–476. DOI: 10.5194/se-5-461-2014.
- [1014] Igor Yu Popov and Ilya V Makeev. “A benchmark solution for 2D Stokes flow over cavity”. In: *Zeitschrift für angewandte Mathematik und Physik* 65.2 (2014), pp. 339–348.
- [1015] Jean H Prevost and Benjamin Loret. “Dynamic strain localization in elasto-(visco-) plastic solids, part 2. Plane strain examples”. In: *Computer Methods in Applied Mechanics and Engineering* 83.3 (1990), pp. 275–294. DOI: 10.1016/0045-7825(90)90074-V.
- [1016] Matthew Price. “Investigating the initial condition of mantle models using data assimilation”. PhD thesis. Cardiff University, 2016.
- [1017] TE Price. “Numerically exact integration of a family of axisymmetric finite elements”. In: *Communications in numerical methods in engineering* 19.4 (2003), pp. 253–261. DOI: 10.1002/cnm.583.

- [1018] P.J. Prince and J.R. Dormand. “High order embedded Runge-Kutta formulae”. In: *Journal of Computational and Applied Mathematics* 7.1 (1981), pp. 67–75. DOI: 10.1016/0771-050X(81)90010-3.
- [1019] A. Prosperetti. “Motion of two superposed viscous fluids”. In: *Phys. Fluids* 24.7 (1981), pp. 1217–1223.
- [1020] E.G. Puckett, D.L. Turcotte, Y. He, H. Lokavarapu, J.M. Robey, and L.H. Kellogg. “New numerical approaches for modeling thermochemical convection in a compositionally stratified fluid”. In: *Phys. Earth. Planet. Inter.* 276 (2018), pp. 10–35. DOI: 10.1016/j.pepi.2017.10.004.
- [1021] A.E. Pusok, B.J.P. Kaus, and A.A. Popov. “On the Quality of Velocity Interpolation Schemes for Marker-in-Cell Method and Staggered Grids”. In: *Pure and Applied Geophysics* (2016). DOI: 10.1007/s00024-016-1431-8.
- [1022] W.M. Putman and S.-J. Lin. “Finite-Volume transport on various cubed-sphere grids”. In: *J. Comp. Phys.* 227 (2007), pp. 55–78. DOI: 10.1016/j.jcp.2007.07.022.
- [1023] R.N. Pysklywec, C. Beaumont, and P. Fullsack. “Lithospheric deformation during the early stages of continental collision: Numerical experiments and comparison with South Island, New Zealand”. In: *J. Geophys. Res.* 107.B72133 (2002). DOI: 10.1029/2001JB000252.
- [1024] R.N. Pysklywec, C. Beaumont, and P. Fullsack. “Modeling the behavior of continental mantle lithosphere during plate convergence”. In: *Geology* 28.7 (2000), pp. 655–658. DOI: 10.1130/0091-7613(2000)28<655:MTBOTC>2.0.CO;2.
- [1025] J. Qin and S. Zhang. “On the selective local stabilization of the mixed Q1-P0 element”. In: *Int. J. Num. Meth. Eng.* 55.12 (2007), pp. 1121–1141. DOI: 10.1002/fld.1505.
- [1026] J. Qin and S. Zhang. “Stability and approximability of the P1-P0 element for stokes equations”. In: *Int. J. Num. Meth. Eng.* 54 (2007), pp. 497–515. DOI: 10.1002/fld.1407.
- [1027] M.E.T. Quinquis and S.J.H. Buiter. “Testing the effects of basic numerical implementations of water migration on models of subduction dynamics”. In: *Solid Earth* 5 (2014), pp. 537–555. DOI: 10.5194/se-5-537-2014.
- [1028] Matthieu E.T. Quinquis, Suzanne J.H. Buiter, and Susan Ellis. “The role of boundary conditions in numerical models of subduction zone dynamics”. In: *Tectonophysics* 497 (2011), pp. 57–70. DOI: 10.1016/j.tecto.2010.11.001.
- [1029] J. Quinteros, V.A. Ramos, and P.M. Jacovkis. “An elasto-visco-plastic model using the finite element method for crustal and lithospheric deformation”. In: *Journal of Geodynamics* 48 (2009), pp. 83–94. DOI: 10.1016/j.jog.2009.06.006.
- [1030] J. Quinteros, S.V. Sobolev, and A.A. Popov. “Viscosity in transition zone and lower mantle: Implications for slab penetration”. In: *Geophys. Res. Lett.* 37.L09307 (2010). DOI: 10.1029/2010GL043140.
- [1031] T. Rabczuk, P.M.A. Areias, and T. Belytschko. “A simplified mesh-free method for shear bands with cohesive surfaces”. In: *Int. J. Num. Meth. Eng.* 69 (2007), pp. 993–1021. DOI: 10.1002/nme.1797.
- [1032] A. Ramachandran. “Parallel adaptive finite element simulation using distributed quad-tree/octree forest”. MA thesis. Institute for Structural Mechanics, Ruhr University Bochum, 2016.
- [1033] Balasubramaniam Ramaswamy. “Numerical simulation of unsteady viscous free surface flow”. In: *Journal of Computational Physics* 90.2 (1990), pp. 396–430.

- [1034] Balasubramaniam Ramaswamy and Mutsuto Kawahara. “Arbitrary Lagrangian–Eulerian finite element method for unsteady, convective, incompressible viscous free surface fluid flow”. In: *International Journal for Numerical Methods in Fluids* 7.10 (1987), pp. 1053–1075.
- [1035] H. Ramberg. *Gravity, Deformation, and the Earth’s Crust: In Theory, Experiments and Geological Application*. Academic Press, London, 214pp., 1967.
- [1036] Hans Ramberg. “Folding of laterally compressed multilayers in the field of gravity, I”. In: *Physics of the Earth and Planetary Interiors* 2.4 (1970), pp. 203–232. DOI: 10.1016/0031-9201(70)90010-5.
- [1037] Hans Ramberg. *Gravity, deformation and the earth’s crust: in theory, experiments and geological application*. Academic press, 1981.
- [1038] Hans Ramberg. “Instability of layered systems in the field of gravity”. In: *Phys. Earth. Planet. Inter.* 1 (1968), pp. 427–447. DOI: 10.1016/0031-9201(68)90014-9.
- [1039] Ekkehard Ramm et al. *Error-controlled adaptive finite elements in solid mechanics*. 2003.
- [1040] G. Ranalli. “Rheology of the crust and its role in tectonic reactivation”. In: *Journal of Geodynamics* 30 (2000), pp. 3–15.
- [1041] G. Ranalli. *Rheology of the Earth*. Springer, 1995. ISBN: 0-412-54670-1.
- [1042] G. Ranalli. “Rheology of the lithosphere in space and time”. In: *Geological Society Special Publications* 121 (1997), pp. 19–37.
- [1043] Giorgio Ranalli. “Rheology and deep tectonics”. In: *Annals of Geophysics* 40.3 (1997).
- [1044] M. Rancic, R.J. Purser, and F. Mesinger. “A global shallow-water model using an expanded spherical cube: Gnomonic versus conformal coordinates”. In: *Q. J. R. Meteorol. Soc.* 122 (1996), pp. 959–982.
- [1045] R. Rannacher and S. Turek. “Simple Nonconforming Quadrilateral Stokes Element”. In: *Numerical Methods for Partial Differential Equations* 8 (1992), pp. 97–111. DOI: 10.1002/num.1690080202.
- [1046] J.T. Ratcliff, G. Schubert, and A. Zebib. “Steady tetrahedral and cubic patterns of spherical shell convection with temperature-dependent viscosity”. In: *J. Geophys. Res.* 101.B11 (1996), pp. 25, 473–25, 484. DOI: 10.1029/96JB02097.
- [1047] HT Rathod and Sridevi Kilari. “General complete Lagrange family for the cube in finite element interpolations”. In: *Computer methods in applied mechanics and engineering* 181.1-3 (2000), pp. 295–344. DOI: 10.1016/S0045-7825(99)00080-8.
- [1048] P-A Raviart and Jean-Marie Thomas. “Primal hybrid finite element methods for 2nd order elliptic equations”. In: *Mathematics of computation* 31.138 (1977), pp. 391–413. DOI: 10.1090/S0025-5718-1977-0431752-8.
- [1049] William H Raymond and Arthur Garder. “Selective damping in a Galerkin method for solving wave problems with variable grids”. In: *Monthly weather review* 104.12 (1976), pp. 1583–1590.
- [1050] J.N. Reddy. “On penalty function methods in the finite element analysis of flow problems”. In: *Int. J. Num. Meth. Fluids* 2 (1982), pp. 151–171. DOI: 10.1002/flid.1650020204.
- [1051] J.N. Reddy and D.K. Gartling. *The Finite Element Method in Heat Transfer and Fluid Dynamics*. CRC Press, 2010. ISBN: 978-1-4200-8598-3.
- [1052] Hannah L Redmond and Scott D King. “A numerical study of a mantle plume beneath the Tharsis Rise: Reconciling dynamic uplift and lithospheric support models”. In: *Journal of Geophysical Research: Planets* 109.E9 (2004). DOI: 10.1029/2003JE002228.

- [1053] K. Regenauer-Lieb and D.A. Yuen. “Rapid conversion of elastic energy into plastic shear heating during incipient necking of the lithosphere”. In: *Geophysical Research Letters* 25.14 (1998), pp. 2737–2740. DOI: 10.1029/98GL02056.
- [1054] Klaus Regenauer-Lieb and Jean-Pierre Petit. “Cutting of the European continental lithosphere: Plasticity theory applied to the present Alpine collision”. In: *J. Geophys. Res.* 102 (Apr. 1997), pp. 7731–7746.
- [1055] Klaus Regenauer-Lieb et al. “From point defects to plate tectonic faults”. In: *Philosophical Magazine* 86.21-22 (2006), pp. 3373–3392. DOI: 10.1080/14786430500375159.
- [1056] Mirko Reguzzoni and Daniele Sampietro. “GEMMA: An Earth crustal model based on GOCE satellite data”. In: *International Journal of Applied Earth Observation and Geoinformation* 35 (2015), pp. 31–43. DOI: 10.1016/j.jag.2014.04.002.
- [1057] Mirko Reguzzoni and Daniele Sampietro. “Moho estimation using GOCE data: a numerical simulation”. In: *Geodesy for Planet Earth*. Springer, 2012, pp. 205–214.
- [1058] Mirko Reguzzoni, Daniele Sampietro, and Fernando Sansò. “Global Moho from the combination of the CRUST2. 0 model and GOCE data”. In: *Geophysical Journal International* 195.1 (2013), pp. 222–237.
- [1059] M. ur Rehman, C. Vuik, and G. Segal. “SIMPLE-type preconditioners for the Oseen problem”. In: *International Journal for Numerical Methods in Fluids* 61 (2009), pp. 432–452.
- [1060] Robert J Renka. “Algorithm 751: TRIPACK: a constrained two-dimensional Delaunay triangulation package”. In: *ACM Transactions on Mathematical Software (TOMS)* 22.1 (1996), pp. 1–8.
- [1061] Robert J Renka. “Algorithm 772: STRIPACK: Delaunay triangulation and Voronoi diagram on the surface of a sphere”. In: *ACM Transactions on Mathematical Software (TOMS)* 23.3 (1997), pp. 416–434.
- [1062] J. Revenaugh and B. Parsons. “Dynamic topography and gravity anomalies for fluid layers whose viscosity varies exponentially with depth”. In: *Geophysical Journal of the Royal Astronomical Society* 90.2 (1987), pp. 349–368. DOI: 10.1111/j.1365-246X.1987.tb00731.x.
- [1063] Neil M Ribe. “Mantle flow induced by back arc spreading”. In: *Geophysical Journal International* 98.1 (1989), pp. 85–91.
- [1064] NM Ribe. “The dynamics of plume-ridge interaction - II. Off-ridge plumes”. In: *Journal of Geophysical Research: Solid Earth* 101.B7 (1996), pp. 16195–16204.
- [1065] NM Ribe and UR Christensen. “Three-dimensional modeling of plume-lithosphere interaction”. In: *Journal of Geophysical Research: Solid Earth* 99.B1 (1994), pp. 669–682.
- [1066] NM Ribe, UR Christensen, and J Theissing. “The dynamics of plume-ridge interaction, 1: Ridge-centered plumes”. In: *Earth and Planetary Science Letters* 134.1-2 (1995), pp. 155–168.
- [1067] FD Richards, MJ Hoggard, LR Cowton, and NJ White. “Reassessing the thermal structure of oceanic lithosphere with revised global inventories of basement depths and heat flow measurements”. In: *Journal of Geophysical Research: Solid Earth* 123.10 (2018), pp. 9136–9161. DOI: 10.1029/2018JB015998.
- [1068] F. Richter and D. McKenzie. “Simple plate models of mantle convection”. In: *J. Geophys.* 44 (1978), pp. 441–471. DOI: xxxx.
- [1069] F.M. Richter and S.F. Daly. “Convection models having a multiplicity of large horizontal scales”. In: *J. Geophys. Res.* 83 (1978), pp. 4951–4956.

- [1070] Frank M Richter. “Convection and the large-scale circulation of the mantle”. In: *Journal of Geophysical Research* 78.35 (1973), pp. 8735–8745. DOI: 10.1029/JB078i035p08735.
- [1071] Frank M Richter. “Dynamical models for sea floor spreading”. In: *Reviews of Geophysics* 11.2 (1973), pp. 223–287. DOI: 10.1029/RG011i002p00223.
- [1072] Frank M Richter, Stephen F Daly, and Henri-Claude Nataf. “A parameterized model for the evolution of isotopic heterogeneities in a convecting system”. In: *Earth and Planetary Science Letters* 60.2 (1982), pp. 178–194. DOI: 10.1016/0012-821X(82)90002-4.
- [1073] Frank M Richter and Dan P McKenzie. “Parameterizations for the horizontally averaged temperature of infinite Prandtl number convection”. In: *Journal of Geophysical Research: Solid Earth* 86.B3 (1981), pp. 1738–1744. DOI: 10.1029/JB086iB03p01738.
- [1074] A. E. Ringwood. *Composition and Petrology of the Earth’s Mantle*. McGraw-Hill, New York, 1975.
- [1075] B. Riviere. *Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations*. SIAM, 2008.
- [1076] Patrick J Roache. “Quantification of uncertainty in computational fluid dynamics”. In: *Annual review of fluid Mechanics* 29.1 (1997), pp. 123–160.
- [1077] GP Roberts, HA Barnes, and P Carew. “Modelling the flow behaviour of very shear-thinning liquids”. In: *Chemical Engineering Science* 56.19 (2001), pp. 5617–5623.
- [1078] J.M. Robey. “On the Design, Implementation, and Use of a Volume-of-Fluid Interface Tracking Algorithm For Modeling Convection and other Processes in the Earth’s Mantle”. PhD thesis. University of California Davis, 2019.
- [1079] Jonathan M Robey and Elbridge Gerry Puckett. “Implementation of a volume-of-fluid method in a finite element code with applications to thermochemical convection in a density stratified fluid in the earth’s mantle”. In: *Computers & Fluids* 190 (2019), pp. 217–253.
- [1080] M.P. Robichaud and P.A. Tanguy. “Comparison of 3-D finite elements for fluid flow”. In: *Communications in applied numerical methods* 3 (1987), pp. 223–233.
- [1081] Jörg Robl and Kurt Stüwe. “Continental collision with finite indenter strength: 1. Concept and model formulation”. In: *Tectonics* 24.4 (2005).
- [1082] T Rolf, Nicolas Coltice, and PJ Tackley. “Linking continental drift, plate tectonics and the thermal state of the Earth’s mantle”. In: *Earth and Planetary Science Letters* 351 (2012), pp. 134–146.
- [1083] B. Romanowicz. “Can we resolve 3D density heterogeneity in the lower mantle”. In: *Geophys. Res. Lett.* 28.6 (2001), pp. 1107–1110.
- [1084] C. Ronchi, R. Iacono, and P.S. Paolucci. “The ”Cubed Sphere”: A New Method for the Solution of Partial Differential Equations in Spherical Geometry”. In: *J. Comp. Phys.* 124 (1996), pp. 93–114.
- [1085] I. Rose, B. Buffet, and T. Heister. “Stability and accuracy of free surface time integration in viscous flows”. In: *Phys. Earth. Planet. Inter.* 262 (2017), pp. 90–100.
- [1086] Vitoriano Ruas. “Mixed finite element methods with discontinuous pressures for the axisymmetric Stokes problem”. In: *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik: Applied Mathematics and Mechanics* 83.4 (2003), pp. 249–264. DOI: 10.1002/zamm.200310032.
- [1087] Johann Rudi, Yu-hsuan Shih, and Georg Stadler. “Advanced Newton methods for geodynamical models of Stokes flow with viscoplastic rheologies”. In: *Geochemistry, Geophysics, Geosystems* 21.9 (2020), e2020GC009059. DOI: 10.1029/2020GC009059.

- [1088] Johann Rudi et al. “An extreme-scale implicit solver for complex PDEs: highly heterogeneous flow in earth’s mantle”. In: *Proceedings of the international conference for high performance computing, networking, storage and analysis*. ACM. 2015, p. 5. DOI: 10.1145/2807591.2807675.
- [1089] Maxwell L Rudolph, Vedran Lekić, and Carolina Lithgow-Bertelloni. “Viscosity jump in Earth’s mid-mantle”. In: *Science* 350.6266 (2015), pp. 1349–1352. DOI: 10.1126/science.aad4972.
- [1090] Reiner Rummel, Weiyong Yi, and Claudia Stummer. “GOCE gravitational gradiometry”. In: *Journal of Geodesy* 85.11 (2011), p. 777.
- [1091] Y. Saad. “A flexible inner-outer preconditioned GMRES algorithm”. In: *SIAM J. Sci. Comput.* 14.2 (1993), pp. 461–469.
- [1092] Y. Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.
- [1093] Y. Saad and M.H. Schultz. “GMRES: A Generalized Minimal Residual Algorithm for Solving Nonsymmetric Linear Systems”. In: *SIAM Journal on Scientific and Statistical Computing* 7.3 (1986), pp. 856–869. DOI: 10.1137/0907058.
- [1094] M.H. Sadd. *Elasticity. Theory, applications and numerics (Third edition)*. Elsevier, 2014.
- [1095] R. Sadourny. “Conservative Finite-Difference Approximations of the Primitive Equations on Quasi-Uniform Spherical Grids”. In: *Monthly Weather Review* 100.2 (1972), pp. 136–144.
- [1096] Edward B Saff and Amo BJ Kuijlaars. “Distributing many points on a sphere”. In: *The mathematical intelligencer* 19.1 (1997), pp. 5–11. DOI: xxxx.
- [1097] M. Saito and Y. Abe. “Consequences of anisotropic viscosity in the Earth’s mantle (in Japanese, with English Abstract)”. In: *Zisin* 37 (1984), pp. 237–245.
- [1098] A. Salih. *Streamfunction-Vorticity Formulation*. Tech. rep. Department of Aerospace Engineering Indian Institute of Space Science and Technology, Thiruvananthapuram, Mar. 2013.
- [1099] M. Sambridge, J. Braun, and H. McQueen. “Computational methods for natural neighbour interpolation in two and three dimensions”. In: *Computational Techniques and Applications*. Ed. by R. L. May and A. K. Easton (eds). World Scientific, 1996, pp. 685–692.
- [1100] M. Sambridge, J. Braun, and H. McQueen. “Geophysical parametrization and interpolation of irregular data using natural neighbours”. In: *Geophys. J. Int.* 122 (1995), pp. 837–857. DOI: 10.1111/j.1365-246X.1995.tb06841.x.
- [1101] Daniele Sampietro. “Geological units and Moho depth determination in the Western Balkans exploiting GOCE data”. In: *Geophysical Journal International* 202.2 (2015), pp. 1054–1063.
- [1102] H Samuel, V Aleksandrov, and B Deo. “The effect of continents on mantle convective stirring”. In: *Geophysical Research Letters* 38.4 (2011), p. L04307. DOI: 10.1029/2010GL046056.
- [1103] H. Samuel and M. Evonuk. “Modeling advection in geophysical flows with particle level sets”. In: *Geochem. Geophys. Geosyst.* 11.8 (2010). DOI: 10.1029/2010GC003081.
- [1104] Henri Samuel. “A deformable particle-in-cell method for advective transport in geodynamic modelling”. In: *Geophysical Journal International* 214.3 (2018), pp. 1744–1773.
- [1105] Henri Samuel and Nicola Tosi. “The influence of post-perovskite strength on the Earth’s mantle thermal and chemical evolution”. In: *Earth and Planetary Science Letters* 323 (2012), pp. 50–59.
- [1106] Patrick Sanan, Dave A May, Matthias Bollhöfer, and Olaf Schenk. “Pragmatic solvers for 3D Stokes and elasticity problems with heterogeneous coefficients: evaluating modern incomplete LDL T preconditioners”. In: *Solid Earth* 11.6 (2020), pp. 2031–2045. DOI: 10.5194/se-11-2031-2020.



- [1107] Patrick Sanan, Dave A May, Richard T Mills, et al. “DMStag: staggered, structured grids for PETSc”. In: *Journal of Open Source Software* 7.79 (2022), p. 4531. DOI: 10.21105/joss.04531.
- [1108] R.L. Sani, P.M. Gresho, R.L. Lee, and D.F. Griffiths. “The cause and cure (?) of the spurious pressures generated by certain FEM solutions of the incompressible Navier-Stokes equations: part 1”. In: *Int. J. Num. Meth. Fluids* 1 (1981), pp. 17–43. DOI: 10.1002/flid.1650010104.
- [1109] R.L. Sani, P.M. Gresho, R.L. Lee, D.F. Griffiths, and M. Engelman. “The cause and cure (?) of the spurious pressures generated by certain FEM solutions of the incompressible Navier-Stokes equations: part 2”. In: *Int. J. Num. Meth. Fluids* 1 (1981), pp. 171–204. DOI: 10.1002/flid.1650010206.
- [1110] Pierre Saramito. “A damped Newton algorithm for computing viscoplastic fluid flows”. In: *Journal of Non-Newtonian fluid mechanics* 238 (2016), pp. 6–15. DOI: 10.1016/j.jnnfm.2016.05.007.
- [1111] Pierre Saramito. “A new elastoviscoplastic model based on the Herschel–Bulkley viscoplastic model”. In: *Journal of Non-Newtonian Fluid Mechanics* 158.1-3 (2009), pp. 154–161. DOI: 10.1016/j.jnnfm.2008.12.001.
- [1112] Pierre Saramito. *Complex fluids*. Springer, 2016. ISBN: 978-3-319-44361-4.
- [1113] Pierre Saramito and Anthony Wachs. “Progress in numerical simulation of yield stress fluid flows”. In: *Rheologica Acta* 56.3 (2017), pp. 211–230. DOI: 10.1007/s00397-016-0985-9.
- [1114] Atuo Sato and Erik G Thompson. “Finite element models for creeping convection”. In: *Journal of Computational Physics* 22.2 (1976), pp. 229–244. DOI: 10.1016/0021-9991(76)90077-2.
- [1115] Michael Schäfer, Stefan Turek, Franz Durst, Egon Krause, and Rolf Rannacher. “Benchmark computations of laminar flow around a cylinder”. In: *Flow simulation with high-performance computers II*. 1996, pp. 547–566. DOI: 10.1007/978-3-322-89849-4\_39.
- [1116] S.M. Schmalholz. “A simple analytical solution for slab detachment”. In: *Earth Planet. Sci. Lett.* 304 (2011), pp. 45–54. DOI: 10.1016/j.epsl.2011.01.011.
- [1117] S.M. Schmalholz, Y.Yu. Podladchikov, and D.W. Schmid. “A spectral/finite difference method for simulating large deformations of heterogeneous, viscoelastic materials”. In: *Geophy. J. Int.* 145 (2001), pp. 199–208. DOI: 10.1046/j.0956-540x.2000.01371.x.
- [1118] Stefan M Schmalholz. “3D numerical modeling of forward folding and reverse unfolding of a viscous single-layer: Implications for the formation of folds and fold patterns”. In: *Tectonophysics* 446.1-4 (2008), pp. 31–41.
- [1119] J. Schmalzl and U. Hansen. “Mixing the Earth’s mantle by thermal convection: A scale dependent phenomenon”. In: *Geophysical Research Letters* 21.11 (1994), pp. 987–990. DOI: 10.1029/94GL00049.
- [1120] J. Schmalzl, G.A. Houseman, and U. Hansen. “Mixing in vigorous, time-dependent three-dimensional convection and application to Earth’s mantle”. In: *Journal of Geophysical Research B: Solid Earth* 101.B10 (1996), pp. 21847–21858.
- [1121] J. Schmalzl, G.A. Houseman, and U. Hansen. “Mixing properties of three-dimensional (3-D) stationary convection”. In: *Physics of Fluids* 7.5 (1995), pp. 1027–1033. DOI: 10.1063/1.868614.
- [1122] Jörg Schmalzl and Alexander Loddock. “Using subdivision surfaces and adaptive surface simplification algorithms for modeling chemical heterogeneities in geophysical flows”. In: *Geochemistry, Geophysics, Geosystems* 4.9 (2003).

- [1123] H. Schmeling and W.R. Jacoby. “On modelling the lithosphere in mantle convection with non-linear rheology”. In: *Journal of Geophysics* 50 (1981), pp. 89–100.
- [1124] H. Schmeling et al. “A benchmark comparison of spontaneous subduction models - Towards a free surface”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 198–223. DOI: 10.1016/j.pepi.2008.06.028.
- [1125] Harro Schmeling. “Compressible convection with constant and variable viscosity: The effect on slab formation, geoid, and topography”. In: *Journal of Geophysical Research: Solid Earth* 94.B9 (1989), pp. 12463–12481.
- [1126] Harro Schmeling. “On the relation between initial conditions and late stages of Rayleigh-Taylor instabilities”. In: *Tectonophysics* 133.1-2 (1987), pp. 65–80. DOI: 10.1016/0040-1951(87)90281-2.
- [1127] Harro Schmeling, Alexander R Cruden, and Gabriele Marquart. “Finite deformation in and around a fluid sphere moving through a viscous medium: implications for diapiric ascent”. In: *Tectonophysics* 149.1-2 (1988), pp. 17–34.
- [1128] D.W. Schmid and Y.Y. Podlachikov. “Analytical solutions for deformable elliptical inclusions in general shear”. In: *Geophy. J. Int.* 155 (2003), pp. 269–288. DOI: 10.1046/j.1365-246X.2003.02042.x.
- [1129] Max W Schmidt and Stefano Poli. “Experimentally based water budgets for dehydrating slabs and consequences for arc magma generation”. In: *Earth and Planetary Science Letters* 163.1-4 (1998), pp. 361–379. DOI: 10.1016/S0012-821X(98)00142-3.
- [1130] G.E. Schneider, G.D. Raithby, and M.M. Yovanovich. “Finite-element solution procedures for solving the incompressible Navier-Stokes equations using equal order variable interpolation”. In: *Numerical Heat Transfer* 1 (1978), pp. 433–451.
- [1131] Robert Schneiders. “A grid-based algorithm for the generation of hexahedral element meshes”. In: *Engineering with computers* 12.3-4 (1996), pp. 168–177.
- [1132] Robert Schneiders. “Algorithms for quadrilateral and hexahedral mesh generation”. In: *Proceedings of the VKI Lecture Series on Computational Fluid Dynamic, VKI-LS 4* (2000).
- [1133] Robert Schneiders. “Quadrilateral and Hexahedral Element Meshes”. In: *chapter 21?* (1999).
- [1134] Robert Schneiders. “Refining quadrilateral and hexahedral element meshes”. In: *transition* 2 (1996), p. 1.
- [1135] Robert Schneiders and Jürgen Debye. “Refining quadrilateral and brick element meshes”. In: *Modeling, Mesh Generation, and Adaptive Numerical Methods for Partial Differential Equations*. 1995, pp. 53–65.
- [1136] Martin PJ Schöpfer, Conrad Childs, and Tom Manzocchi. “Three-dimensional failure envelopes and the brittle-ductile transition”. In: *Journal of Geophysical Research: Solid Earth* 118.4 (2013), pp. 1378–1392. DOI: 10.1002/jgrb.50081.
- [1137] B Schott and H Schmeling. “Delamination and detachment of a lithospheric root”. In: *Tectonophysics* 296.3-4 (1998), pp. 225–247. DOI: 10.1016/S0040-1951(98)00154-1.
- [1138] P. Schroeder and G. Lube. “Stabilised dG-FEM for incompressible natural convection flows with boundary and moving interior layers on non-adapted meshes”. In: *J. Comp. Phys.* 335 (2017), pp. 760–779.
- [1139] G Schubert, DL Turcotte, and ER Oxburgh. “Stability of planetary interiors”. In: *Geophysical Journal International* 18.5 (1969), pp. 441–460. DOI: 10.1111/j.1365-246X.1969.tb03370.x.

- [1140] G. Schubert, D.L. Turcotte, and P. Olson. *Mantle Convection in the Earth and Planets*. Cambridge University Press, 2001. ISBN: 0-521-70000-0. DOI: 10.1017/CB09780511612879.
- [1141] Gerald Schubert and Charles A Anderson. “Finite element calculations of very high Rayleigh number thermal convection”. In: *Geophysical Journal International* 80.3 (1985), pp. 575–601. DOI: 10.1111/j.1365-246X.1985.tb05112.x.
- [1142] Gerald Schubert and David A Yuen. “Shear heating instability in the Earth’s upper mantle”. In: *Tectonophysics* 50.2-3 (1978), pp. 197–205. DOI: 10.1016/0040-1951(78)90135-X.
- [1143] Gerald Schubert, David A Yuen, and Donald L Turcotte. “Role of phase transitions in a dynamic mantle”. In: *Geophysical Journal International* 42.2 (1975), pp. 705–735. DOI: 10.1111/j.1365-246X.1975.tb05888.x.
- [1144] Melchior Schuh-Senlis, Cedric Thieulot, Paul Cupillard, and Guillaume Caumon. “Towards the application of Stokes flow equations to structural restoration simulations”. In: *Solid Earth* 11 (2020), pp. 1909–1930. DOI: 10.5194/se-11-1909-2020.
- [1145] P.R. Schunk, M.A. Heroux, R.R. Rao, T.A. Baer, S.R. Subia, and A.C. Sun. *Iterative solvers and preconditioners for fully-coupled finite element formulations of incompressible fluid mechanics and related transport problems*. Tech. rep. SAND2001-3512J. Sandia National Laboratories, 2001.
- [1146] Larkin Ridgway Scott and Michael Vogelius. “Conforming finite element methods for incompressible and nearly incompressible continua”. In: *Lectures in Applied Mathematics* 22.2 (1985). DOI: xxxx.
- [1147] A. Segal. *Finite element methods for the incompressible Navier-Stokes equations*. Delft University of Technology, 2012.
- [1148] A. Segal, M. ur Rehman, and C. Vuik. “Preconditioners for Incompressible Navier-Stokes Solvers”. In: *Numer. Math. Theor. Meth. Appl.* 3.3 (2010), pp. 245–275. DOI: 10.4208/nmtma.2010.33.1.
- [1149] C. Echevarria Serur. “Fast iterative methods for solving the incompressible Navier-Stokes equations”. PhD thesis. TU Delft, 2013.
- [1150] Ruben Sevilla and Thibault Duretz. “A face-centered finite volume method for high-contrast Stokes interface problems”. In: *International Journal for Numerical Methods in Engineering* 124 (2023), pp. 3709–3732. DOI: 10.1002/nme.7294.
- [1151] MH Shahnas and WR Peltier. “The impacts of mantle phase transitions and the iron spin crossover in ferropericlase on convective mixing - is the evidence for compositional convection definitive? New results from a Yin-Yang overset grid-based control volume model”. In: *Journal of Geophysical Research: Solid Earth* 120.8 (2015), pp. 5884–5910. DOI: 10.1002/2015JB012064.
- [1152] Farzin Shakib, Thomas JR Hughes, and Zdeněk Johan. “A new finite element formulation for computational fluid dynamics: X. The compressible Euler and Navier-Stokes equations”. In: *Computer Methods in Applied Mechanics and Engineering* 89.1-3 (1991), pp. 141–219.
- [1153] A. Shamekhi and A. Aliabadi. “Non-Newtonian lid-driven cavity flow simulation by mesh free method”. In: *ICCES* 11.3 (2009), pp. 67–72.
- [1154] W. Sharples, L.N. Moresi, M. Velic, M.A. Jadamec, and D.A. May. “Simulating faults and plate boundaries with a transversely isotropic plasticity model”. In: *Phys. Earth. Planet. Inter.* 252 (2016), pp. 77–90. DOI: 10.1016/j.pepi.2015.11.007.

- [1155] Dongwoo Sheen. “P1–Nonconforming Polyhedral Finite Elements in High Dimensions”. In: *2018 MATRIX Annals*. Ed. by Jan de Gier, Cheryl E. Praeger, and Terence Tao. Cham: Springer International Publishing, 2020, pp. 121–133. ISBN: 978-3-030-38230-8. DOI: 10.1007/978-3-030-38230-8\_9.
- [1156] J.R. Shewchuk. “An Introduction to the Conjugate Gradient Method Without the Agonizing Pain”. In: (1994).
- [1157] J.R. Shewchuk. “Triangle: Engineering a 2D Quality Mesh Generator and Delaunay Triangulator”. In: *Applied Computational Geometry: Towards Geometric Engineering*. Ed. by Ming C. Lin and Dinesh Manocha. Vol. 1148. Lecture Notes in Computer Science. From the First ACM Workshop on Applied Computational Geometry. 1996, pp. 203–222. DOI: 10.1007/BFb0014497.
- [1158] Jonathan Richard Shewchuk. “Delaunay refinement algorithms for triangular mesh generation”. In: *Computational geometry* 22.1-3 (2002), pp. 21–74.
- [1159] Jonathan Richard Shewchuk. “Reprint of: Delaunay refinement algorithms for triangular mesh generation”. In: *Computational Geometry* 47.7 (2014), pp. 741–778. DOI: 10.1016/j.comgeo.2014.02.005.
- [1160] Yanan Shi, Dongping Wei, Zhong-Hai Li, Ming-Qi Liu, and Mengxue Liu. “Subduction mode selection during slab and mantle transition zone interaction: Numerical modeling”. In: *Pure Appl. Geophys.* (2017), pp. 5–24. DOI: 10.1007/s00024-017-1762-0.
- [1161] J. S. Shiau, A.V. Lyamin, and S. W. Sloan. “Bearing capacity of a sand layer on clay by finite element limit analysis”. In: *Can. Geotech. J.* 40 (2003), pp. 900–915.
- [1162] M. Shipeng and S. Zhongci. “Nonconforming rotated  $Q_1$  element on non-tensor product anisotropic meshes”. In: *Science in China Series A: Mathematics* 49.10 (2006), pp. 1363–1375. DOI: 10.1007/s11425-006-1363-3.
- [1163] Hang Si. “TetGen, a Delaunay-based quality tetrahedral mesh generator”. In: *ACM Transactions on Mathematical Software (TOMS)* 41.2 (2015), pp. 1–36. DOI: 10.1145/2629697.
- [1164] Karin Sigloch and Mitchell G Mihalynuk. “Intra-oceanic subduction shaped the assembly of Cordilleran North America”. In: *Nature* 496.7443 (2013), p. 50.
- [1165] Amit Kaur Sihota. “Conjugate gradient methods using MPI for distributed systems”. MA thesis. McGill University, 2004.
- [1166] Jesse L Silverberg, Matthew Bierbaum, James P Sethna, and Itai Cohen. “Collective motion of humans in mosh and circle pits at heavy metal concerts”. In: *Physical review letters* 110.22 (2013), p. 228701. DOI: 10.1103/PhysRevLett.110.228701.
- [1167] D.J. Silvester and N. Kechkar. “Stabilised bilinear-constant velocity-pressure finite elements for the conjugate gradient solution of the stokes problem”. In: *Computer Methods in Applied Mechanics and Engineering* 79.1 (1990), pp. 71–86. DOI: 10.1016/0045-7825(90)90095-4.
- [1168] Nathan Sime, Jakob M Maljaars, Cian R Wilson, and Peter E van Keken. “An exactly mass conserving and pointwise divergence free velocity method: Application to compositional buoyancy driven flow problems in geodynamics”. In: *Geochemistry, Geophysics, Geosystems* 22.4 (2021), e2020GC009349. DOI: 10.1029/2020GC009349.
- [1169] Nathan Sime and Cian R Wilson. “Automatic weak imposition of free slip boundary conditions via Nitsche’s method: application to nonlinear problems in geodynamics”. In: *arXiv preprint arXiv:2001.10639* (2020).
- [1170] Nathan Sime, Cian R Wilson, and Peter E van Keken. “A pointwise conservative method for thermochemical convection under the compressible anelastic liquid approximation”. In: *Geochemistry, Geophysics, Geosystems* 23.2 (2022), e2021GC009922. DOI: 10.1029/2021GC009922.

- [1171] N.A. Simmons, S.C. Myers, G. Johannesson, and E. Matzel. “LLNL-G3Dv3: Global P wave tomography model for improved regional and teleseismic travel time prediction”. In: *J. Geophys. Res.* 117.B10302 (2012). DOI: 10.1029/2012JB009525.
- [1172] Guy Simpson. *Practical Finite Element Modelin in Earth Science Using Matlab*. Wiley-Blackwell, 2017. ISBN: 978-1-119-24862-0.
- [1173] N.H. Sleep, S. Stein, R.J. Geller, and R.G. Gordon. “Comment on ”The use of the minimum-dissipation principle in tectonophysics””. In: *Earth Planet. Sci. Lett.* 45 (1979), pp. 218–220.
- [1174] S.W. Sloan. “A fortran program for profile and wavefront reduction”. In: *Int. J. Num. Meth. Eng.* 28 (1989), pp. 2651–2679.
- [1175] S.W. Sloan. “An algorithm for profile and wavefront reduction of sparse matrices”. In: *International Journal for Numerical Methods in Engineering* 23.2 (1986), pp. 239–251.
- [1176] F. Soboutia, A. Ghodsb, and J. Arkani-Hamed. “On the advection of sharp material interfaces in geodynamic problems: entrainment of the D” layer”. In: *Journal of Geodynamics* 31 (2001), pp. 459–479.
- [1177] Larry P Solheim and WR Peltier. “Avalanche effects in phase transition modulated thermal convection: A model of Earth’s mantle”. In: *Journal of Geophysical Research: Solid Earth* 99.B4 (1994), pp. 6997–7018.
- [1178] LP Solheim and WR Peltier. “Heat transfer and the onset of chaos in a spherical, axisymmetric, anelastic model of whole mantle convection”. In: *Geophysical & Astrophysical Fluid Dynamics* 53.4 (1990), pp. 205–255.
- [1179] VS Solomatov. “Scaling of temperature-and stress-dependent viscosity convection”. In: *Physics of Fluids* 7.2 (1995), pp. 266–274.
- [1180] C Sotin and S Labrosse. “Three-dimensional thermal convection in an iso-viscous, infinite Prandtl number fluid heated from within and from below: applications to the transfer of heat through planetary mantles”. In: *Physics of the Earth and Planetary Interiors* 112.3-4 (1999), pp. 171–190. DOI: 10.1016/S0031-9201(99)00004-7.
- [1181] A. Soulaïmani, M. Fortin, Y. Ouellet, G. Dhatt, and F. Bertrand. “Simple continuous pressure elements for two- and three-dimensional incompressible flows”. In: *Computer Methods in Applied Mechanics and Engineering* 62 (1987), pp. 47–69. DOI: 10.1016/0045-7825(87)90089-2.
- [1182] M. Souli and J.P. Zolesio. “Arbitrary Lagrangian-Eulerian and free surface methods in fluid mechanics”. In: *Comput. Methods Appl. Mech. Engrg* 191 (2001), pp. 451–466. DOI: 10.1016/S0045-7825(01)00313-9.
- [1183] E.A. de Souza Neto, D. Peric, and D.R.J. Owen. *Computational methods for plasticity*. Wiley, 2008.
- [1184] Frank J Spera, David A Yuen, and Stephen J Kirschvink. “Thermal boundary layer convection in silicic magma chambers: Effects of temperature-dependent rheology and implications for thermogravitational chemical fractionation”. In: *Journal of Geophysical Research: Solid Earth* 87.B10 (1982), pp. 8755–8767. DOI: 10.1029/JB087iB10p08755.
- [1185] Charles G Speziale. “On the advantages of the vorticity-velocity formulation of the equations of fluid dynamics”. In: *J. Comp. Phys.* 73 (1987), pp. 476–480.
- [1186] Edward A Spiegel and G Veronis. “On the Boussinesq approximation for a compressible fluid.” In: *The Astrophysical Journal* 131 (1960), p. 442.

- [1187] M. Spiegelman, D.A. May, and C. Wilson. “On the solvability of incompressible Stokes with viscoplastic rheologies in geodynamics”. In: *Geochem. Geophys. Geosyst.* 17 (2016), pp. 2213–2238. DOI: 10.1002/2015GC006228.
- [1188] Marc Spiegelman and Richard F Katz. “A semi-Lagrangian Crank-Nicolson algorithm for the numerical solution of advection-diffusion problems”. In: *Geochemistry, Geophysics, Geosystems* 7.4 (2006), Q04014. DOI: 10.1029/2005GC001073.
- [1189] Marc Spiegelman and Dan McKenzie. “Simple 2-D models for melt extraction at mid-ocean ridges and island arcs”. In: *Earth and Planetary Science Letters* 83.1-4 (1987), pp. 137–152.
- [1190] Frank D Stacey. “A thermal model of the Earth”. In: *Physics of the Earth and Planetary Interiors* 15.4 (1977), pp. 341–348.
- [1191] Frank D Stacey and Paul M Davis. *Physics of the Earth*. 2008.
- [1192] G. Stadler, M. Gurnis, C. Burstedde, L.C. Wilcox, L. Alisic, and O. Ghattas. “The dynamics of plate tectonics and mantle flow: from local to global scales”. In: *Science* 329 (2010), pp. 1033–1038. DOI: 10.1126/science.1191223.
- [1193] Andrew Staniforth and Jean Côté. “Semi-Lagrangian integration schemes for atmospheric models - A review”. In: *Monthly weather review* 119.9 (1991), pp. 2206–2223.
- [1194] Matthew Staten and Scott A Canann. “Post refinement element shape improvement for quadrilateral meshes”. In: *220 Trends in Unstructured Mesh Generation, ASME*. Citeseer. 1997.
- [1195] Ph. Steer, R. Cattin, J. Lavé, and V. Godard. “Surface Lagrangian Remeshing: A new tool for studying long term evolution of continental lithosphere from 2D numerical modelling”. In: *Computers and Geosciences* 37.8 (2011), pp. 1067–1074. DOI: 10.1016/j.cageo.2010.05.023.
- [1196] C Stein and U Hansen. “Arrhenius rheology versus Frank-Kamenetskii rheology - Implications for mantle dynamics”. In: *Geochemistry, Geophysics, Geosystems* 14.8 (2013), pp. 2757–2770. DOI: 10.1002/ggge.20158.
- [1197] C. Stein, J. Lowman, and U. Hansen. “A comparison of mantle convection models featuring plates”. In: *Geochem. Geophys. Geosyst.* 15 (2014), pp. 2689–2698. DOI: 10.1002/2013GC005211.
- [1198] S. Stein. “A model for the relation between spreading rate and oblique spreading”. In: *Earth Planet. Sci. Lett.* 39 (1978), pp. 313–318.
- [1199] Volker Steinbach, Ulrich Hansen, and Adolf Ebel. “Compressible convection in the earth’s mantle: a comparison of different approaches”. In: *Geophysical Research Letters* 16.7 (1989), pp. 633–636. DOI: 10.1029/GL016i007p00633.
- [1200] Volker Steinbach and David A Yuen. “The non-adiabatic nature of mantle convection as revealed by passive tracers”. In: *Earth and Planetary Science Letters* 136.3-4 (1995), pp. 241–250. DOI: 10.1016/0012-821X(95)00166-A.
- [1201] Volker Steinbach, David A Yuen, and Wuling Zhao. “Instabilities from phase transitions and the timescales of mantle thermal evolution”. In: *Geophysical research letters* 20.12 (1993), pp. 1119–1122. DOI: 10.1029/93GL01243.
- [1202] B Steinberger and R Holme. “Mantle flow models with core-mantle boundary constraints and chemical heterogeneities in the lowermost mantle”. In: *Journal of Geophysical Research: Solid Earth* 113.B5 (2008).

- [1203] B. Steinberger and A.R. Calderwood. “Models of large-scale viscous flow in the Earth’s mantle with constraints from mineral physics and surface observations”. In: *Geophys. J. Int.* 167 (2006), pp. 1461–1481. DOI: 10.1111/j.1365-246X.2006.03131.x.
- [1204] Bernhard Steinberger. “Topography caused by mantle density variations: observation-based estimates and models derived from tomography and lithosphere thickness”. In: *Geophysical Supplements to the Monthly Notices of the Royal Astronomical Society* 205.1 (2016), pp. 604–621. DOI: 10.1093/gji/ggw040.
- [1205] K. Stemmer, H. Harder, and U. Hansen. “A new method to simulate convection with strongly temperature- and pressure-dependent viscosity in a spherical shell: Applications to the Earth’s mantle”. In: *Phys. Earth. Planet. Inter.* 157 (2006), pp. 223–249. DOI: 10.1016/j.pepi.2006.04.007.
- [1206] R. Stenberg. “Analysis of Mixed Finite Element Methods for the Stokes Problem: A unified approach”. In: *Mathematics of Computation* 42.165 (1984), pp. 9–23.
- [1207] R. Stenberg. “Error analysis of some finite element methods for the Stokes problem”. In: *Mathematics of Computation* 54.190 (1990), pp. 495–508.
- [1208] Ove Stephansson and Harald Berner. “The finite element method in tectonic processes”. In: *Physics of the Earth and Planetary Interiors* 4.4 (1971), pp. 301–321.
- [1209] Pietro Sternai. “Surface processes forcing on extensional rock melting”. In: *Scientific Reports* 10.1 (2020), pp. 1–13. DOI: 10.1038/s41598-020-63920-w.
- [1210] Eli Sternberg and F Rosenthal. “The elastic sphere under concentrated loads”. In: *J. Appl. Mech.* 19(4) (1952), pp. 413–421. DOI: 10.1115/1.4010536.
- [1211] D. Sterpi. “An analysis of geotechnical problems involving strain softening effects”. In: *International Journal for Numerical and Analytical Methods in Geomechanics* 23 (1999), pp. 1427–1454. DOI: 10.1002/(SICI)1096-9853(199911)23:13<1427::AID-NAG6>3.0.CO;2-B.
- [1212] WNR Stevens. “Finite element, stream function–vorticity solution of steady laminar natural convection”. In: *International Journal for Numerical Methods in Fluids* 2.4 (1982), pp. 349–366. DOI: 10.1002/flid.1650020404.
- [1213] Lars Stixrude and Carolina Lithgow-Bertelloni. “Mineralogy and elasticity of the oceanic upper mantle: Origin of the low-velocity zone”. In: *Journal of Geophysical Research: Solid Earth* 110.B3 (2005).
- [1214] Lars Stixrude and Carolina Lithgow-Bertelloni. “Thermodynamics of mantle minerals – I. Physical properties”. In: *Geophysical Journal International* 162.2 (2005), pp. 610–632. DOI: 10.1111/j.1365-246X.2005.02642.x.
- [1215] Glen S Stockmal, Christopher Beaumont, and Ross Boutilier. “Geodynamic models of convergent margin tectonics: transition from rifted margin to overthrust belt and consequences for foreland-basin development”. In: *AAPG Bulletin* 70.2 (1986), pp. 181–190. DOI: 10.1306/94885656-1704-11D7-8645000102C1865D.
- [1216] Yvonne Marie Stokes. “Very Viscous Flows Driven by Gravity with particular application to Slumping of Molten Glass”. PhD thesis. University of Adelaide, 1998.
- [1217] J. Suckale, B.H. Hager, L.T. Elkins-Tanton, and J.-Ch. Nave. “It takes three to tango: 2. Bubble dynamics in basaltic volcanoes and ramifications for modeling normal Strombolian activity”. In: *J. Geophys. Res.* 115.B7 (2010). DOI: 10.1029/2009JB006917.
- [1218] J. Suckale, J.-C. Nave, and B.H. Hager. “It takes three to tango: 1. Simulating buoyancy-driven flow in the presence of large viscosity contrasts”. In: *J. Geophys. Res.* 115.B07409 (2010). DOI: 10.1029/2009JB006916.

- [1219] J. Sung, H.G. Choi, and J.Y. Yoo. “Time-accurate computation of unsteady free surface flows using an ALE-segregated equal-order FEM”. In: *Computer Methods in Applied Mechanics and Engineering* 190 (2000), pp. 1425–1440.
- [1220] M. Sussman and E.G. Puckett. “A Coupled Level Set and Volume-of-Fluid Method for Computing 3D and Axisymmetric Incompressible Two-Phase Flows”. In: *J. Comp. Phys.* 162 (2000), pp. 301–337.
- [1221] K. Sverdrup, N. Nikiforakis, and A. Almgren. “Highly parallelisable simulations of time-dependent viscoplastic fluid flow simulations with structured adaptive mesh refinement”. In: *Physics of Fluids* 30 (2018), p. 093102. DOI: 10.1063/1.5049202.
- [1222] R. Swinbank and R.J. Purser. “Fibonacci grids: A novel approach to global modelling”. In: *Quarterly Journal of the Royal Meteorological Society* 132 (2006), pp. 1769–1793.
- [1223] A. Syrakos, G.C. Georgiou, and A.N. Alexandrou. “Performance of the finite volume method in solving regularised Bingham flows: Inertia effects in the lid-driven cavity flow”. In: *Journal of Non-Newtonian Fluid Mechanics* 208–209 (2014), pp. 88–107. DOI: 10.1016/j.jnnfm.2014.03.004.
- [1224] M. Tabata and A. Suzuki. “A stabilized finite element method for the Rayleigh-Bénard equations with infinite Prandtl number in a spherical shell”. In: *Computer Methods in Applied Mechanics and Engineering* 190 (2000), pp. 387–402.
- [1225] Masahisa Tabata. “Finite element analysis of axisymmetric flow problems”. In: *Zeitschrift für angewandte Mathematik und Mechanik* 76 (1996), pp. 171–174.
- [1226] M. Tachibana and Y. Iemoto. “Steady flow around, and drag on a circular cylinder moving at low speeds in a viscous liquid between two parallel planes”. In: *Fluid Dynamics Research* 2 (1987), pp. 125–137.
- [1227] P. Tackley. “Three-dimensional models of mantle convection: Influence of phase transitions and temperature-dependent viscosity”. PhD thesis. California Institute of Technology, 1994, p290.
- [1228] P.J. Tackley. “Modelling compressible mantle convection with large viscosity contrasts in a three-dimensional spherical shell using the yin-yang grid”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 7–18.
- [1229] P.J. Tackley and S.D. King. “Testing the tracer ratio method for modeling active compositional fields in mantle convection simulations”. In: *Geochem. Geophys. Geosyst.* 4.4 (2003). DOI: 10.1029/2001GC000214.
- [1230] Paul J Tackley and Shunxing Xie. “The thermochemical structure and evolution of Earth’s mantle: constraints and numerical models”. In: *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences* 360.1800 (2002), pp. 2593–2609. DOI: 10.1098/rsta.2002.1082.
- [1231] E. Taliadorou, G.C. Georgiou, and I. Moulitsas. “Weakly compressible Poiseuille flows of a Herschel-Bulkley fluid”. In: *Journal of Non-Newtonian Fluid Mechanics* 158 (2009), pp. 162–169. DOI: 10.1016/j.jnnfm.2008.11.010.
- [1232] M. Talwani and M. Ewing. “Rapid Computation of Gravitational Attraction of Three-Dimensional Bodies of Arbitrary Shape”. In: *Geophysics* 25.1 (1960), pp. 203–225.
- [1233] Manik Talwani, J Lamar Worzel, and Mark Landisman. “Rapid gravity computations for two-dimensional bodies with application to the Mendocino submarine fracture zone”. In: *Journal of geophysical research* 64.1 (1959), pp. 49–59.
- [1234] E. Tan and M. Gurnis. “Compressible thermochemical convection and application to lower mantle structures”. In: *J. Geophys. Res.* 112.B06304 (2007).



- [1235] A. Tanaka, Y. Sanada, Y. Yamada, T. Matsuoka, and Y. Ashida. “discrete element simulations of continental collision in asia”. In: *Exploration Geophysics* 36.1 (2005), pp. 1–6.
- [1236] R.I. Tanner and E. Tanner. “Heinrich Hencky: a rheological pioneer”. In: *Rheol. Acta* 42 (2003), pp. 93–101.
- [1237] P. Tapponnier, G. Peltzer, A.Y. Le Dain, R. Armijo, and P. Cobbold. “Propagating extrusion tectonics in Asia: new insights from simple experiments with plasticine”. In: *Geology* 10 (1982), pp. 611–616.
- [1238] Paul Tapponnier and Peter Molnar. “Slip-line field theory and large-scale continental tectonics”. In: *Nature* 264 (Nov. 1976), pp. 319–324.
- [1239] J. Tarduno, H.P. Bunge, N. Sleep, and U. Hansen. “The bent hawaiian-emperor hotspot track: inheriting the mantle wind”. In: *Science* 324.5923 (2009), pp. 50–53. DOI: 10.1126/science.1161256.
- [1240] C. Tayloor and P. Hood. “A numerical solution of the Navier-Stokes equations using the finite element technique”. In: *Comput. Fluids* 1 (1973), pp. 73–100. DOI: 10.1016/0045-7930(73)90027-3.
- [1241] Alexandra J. Taylor and Simon D. R. Wilson. “Conduit flow of an incompressible, yield-stress fluid”. In: *Journal of Rheology* 41 (1997), pp. 93–101. DOI: 10.1122/1.550802.
- [1242] Arkady Ten, David A Yuen, Tine B Larsen, and Andrei V Malevsky. “The evolution of material surfaces in convection with variable viscosity as monitored by a characteristics-based method”. In: *Geophysical research letters* 23.16 (1996), pp. 2001–2004. DOI: 10.1029/96GL02182.
- [1243] Arkady Ten, David A Yuen, Yu Yu Podladchikov, Tine B Larsen, Elizaveta Pachepsky, and Andrei V Malevsky. “Fractal features in mixing of non-Newtonian and Newtonian mantle convection”. In: *Earth and planetary science letters* 146.3-4 (1997), pp. 401–414. DOI: 10.1016/S0012-821X(96)00244-0.
- [1244] Arkady A Ten, Yuri Yu Podladchikov, David A Yuen, Tine B Larsen, and Andrei V Malevsky. “Comparison of mixing properties in convection with the Particle-Line Method”. In: *Geophysical research letters* 25.16 (1998), pp. 3205–3208. DOI: <https://doi.org/10.1029/98GL51991>.
- [1245] J.L. Tetreault and S.J.H. Buiter. “Geodynamic models of terrane accretion: Testing the fate of island arcs, oceanic plateaus, and continental fragments in subduction zones”. In: *J. Geophys. Res.* 117 (2012), B08403. DOI: 10.1029/2012JB009316.
- [1246] Jannis Teunissen and Ute Ebert. “Afino: A framework for quadtree/octree AMR with shared-memory parallelization and geometric multigrid methods”. In: *Computer Physics Communications* 233 (2018), pp. 156–166.
- [1247] T.E. Tezduyar. “Stabilized Finite Element Formulations for Incompressible Flow Computations”. In: *Advances in Applied Mechanics* 28 (1992), pp. 1–44.
- [1248] Tayfun E Tezduyar, Mittal Behr, S Mittal, and J Liou. “A new strategy for finite element computations involving moving boundaries and interfaces - the deforming-spatial-domain/space-time procedure: II. Computation of free-surface flows, two-liquid flows, and flows with drifting cylinders”. In: *Computer methods in applied mechanics and engineering* 94.3 (1992), pp. 353–371.
- [1249] Tayfun E Tezduyar, Sanjay Mittal, SE Ray, and R Shih. “Incompressible flow computations with stabilized bilinear and linear equal-order-interpolation velocity-pressure elements”. In: *Computer Methods in Applied Mechanics and Engineering* 95.2 (1992), pp. 221–242.

- [1250] Tayfun E. Tezduyar and Yasuo Osawa. “Finite element stabilization parameters computed from element matrices and vectors”. In: *Computer Methods in Applied Mechanics and Engineering* 190 (2000), pp. 411–430.
- [1251] TE Tezduyar and YJ Park. “Discontinuity-capturing finite element formulations for nonlinear convection-diffusion-reaction equations”. In: *Computer Methods in Applied Mechanics and Engineering* 59.3 (1986), pp. 307–325.
- [1252] Thomas M Tharp. “Numerical models of subduction and forearc deformation”. In: *Geophysical Journal International* 80.2 (1985), pp. 419–437. DOI: 10.1111/j.1365-246X.1985.tb05102.x.
- [1253] M. Thielmann, B.J.P. Kaus, and A.A. Popov. “Lithospheric stresses in Rayleigh-Benard convection: effects of a free surface and a viscoelastic Maxwell rheology”. In: *Geophy. J. Int.* 203 (2015), pp. 2200–2219. DOI: 10.1093/gji/ggv436.
- [1254] M. Thielmann, D.A. May, and B.J.P. Kaus. “Discretization errors in the Hybrid Finite Element Particle-In-Cell Method”. In: *Pure and Applied Geophysics* 171.9 (2014), pp. 2164–2184. DOI: 10.1007/s00024-014-0808-9.
- [1255] P. van Thienen. “Convective vigour and heat flow in chemically differentiated systems”. In: *Geophy. J. Int.* 169.2 (2007), pp. 747–766. DOI: 10.1111/j.1365-246X.2007.03377.x.
- [1256] C. Thieulot. “Analytical solution for viscous incompressible Stokes flow in a spherical shell”. In: *Solid Earth* 2017 (2017), pp. 1–19. DOI: 10.5194/se-2017-71.
- [1257] C. Thieulot. “ELEFANT: a user-friendly multipurpose geodynamics code”. In: *Solid Earth Discussions* 6 (2014), pp. 1949–2096. DOI: 10.5194/sed-6-1949-2014.
- [1258] C. Thieulot. “FANTOM: two- and three-dimensional numerical modelling of creeping flows for the solution of geological problems”. In: *Phys. Earth. Planet. Inter.* 188.1 (2011), pp. 47–68. DOI: 10.1016/j.pepi.2011.06.011.
- [1259] C. Thieulot. “GHOST: Geoscientific Hollow Sphere Tessellation”. In: *Solid Earth* 9.1–9 (2018). DOI: 10.5194/se-9-1169-2018.
- [1260] C. Thieulot and W. Bangerth. “On the choice of finite element for applications in geodynamics”. In: *Solid Earth* 13 (2022), pp. 229–249. DOI: 10.5194/se-13-1-2022.
- [1261] C. Thieulot, P. Fullsack, and J. Braun. “Adaptive octree-based finite element analysis of two- and three-dimensional indentation problems”. In: *J. Geophys. Res.* 113 (2008), B12207. DOI: 10.1029/2008JB005591.
- [1262] C. Thieulot, P. Steer, and R.S. Huismans. “Three-dimensional numerical simulations of crustal systems undergoing orogeny and subjected to surface processes”. In: *Geochem. Geophys. Geosyst.* 15 (2014). DOI: 10.1002/2014GC005490.
- [1263] Cedric Thieulot and Pep Español. “Non-isothermal diffusion in a binary mixture with smoothed particle hydrodynamics”. In: *Computer physics communications* 169.1-3 (2005), pp. 172–176. DOI: 10.1016/j.cpc.2005.03.039.
- [1264] Cedric Thieulot, LPBM Janssen, and Pep Español. “Smoothed particle hydrodynamics model for phase separating fluid mixtures. I. General equations”. In: *Physical Review E* 72.1 (2005), p. 016713. DOI: 10.1103/PhysRevE.72.016713.
- [1265] Cedric Thieulot, LPBM Janssen, and Pep Español. “Smoothed particle hydrodynamics model for phase separating fluid mixtures. II. Diffusion in a binary mixture”. In: *Physical Review E* 72.1 (2005), p. 016714. DOI: 10.1103/PhysRevE.72.016714.

- [1266] B Thomas, H Samuel, CG Farnetani, J Aubert, and C Chauvel. “Mixing time of heterogeneities in a buoyancy-dominated magma ocean”. In: *Geophysical Journal International* 236.2 (2024), pp. 764–777. DOI: 10.1093/gji/ggad452.
- [1267] J.F. Thompson, B.K. Soni, and N.P. Weatherill. *Handbook of grid generation*. CRC press, 1998.
- [1268] F. Tin-Loi and N.S. Ngo. “Performance of the p-version finite element method for limit analysis”. In: *International Journal of Mechanical Sciences* 45 (2003), pp. 1149–1166. DOI: 10.1016/j.ijmecsci.2003.08.004.
- [1269] C. Tirel, J.-P. Brun, and E. Burov. “Dynamics and structural development of metamorphic core complexes”. In: *J. Geophys. Res.* 113.B04403 (2008).
- [1270] Max Tirone. “On the thermal gradient in the Earth’s deep interior”. In: *Solid Earth* 7.1 (2016), p. 229. DOI: 10.5194/se-7-229-2016.
- [1271] Andréa Tommasi, Mickael Knoll, Alain Vauchez, Javier W Signorelli, Catherine Thoraval, and Roland Logé. “Structural reactivation in plate tectonics controlled by olivine crystal anisotropy”. In: *Nature Geoscience* 2.6 (2009), pp. 423–427. DOI: 10.1038/ngeo528.
- [1272] Rosaria Tondi, Maddalena Gilardoni, and Mirko Reguzzoni. “The combined inversion of seismological and GOCE gravity data: New insights into the current state of the Pacific lithosphere and upper mantle”. In: *Tectonophysics* 705 (2017), pp. 101–115.
- [1273] KE Torrance and DL Turcotte. “Structure of convection cells in the mantle”. In: *Journal of Geophysical Research* 76.5 (1971), pp. 1154–1161. DOI: 10.1029/JB076i005p01154.
- [1274] KE Torrance and DL Turcotte. “Thermal convection with large viscosity variations”. In: *Journal of Fluid Mechanics* 47.1 (1971), pp. 113–125. DOI: 10.1017/S002211207100096X.
- [1275] A. Toselli and O. Widlund. *Domain decomposition methods - Algorithms and Theory*. Springer, 2005.
- [1276] N. Tosi et al. “A community benchmark for viscoplastic thermal convection in a 2-D square box”. In: *Geochem. Geophys. Geosyst.* 16.7 (2015), pp. 2175–2196. DOI: 10.1002/2015GC005807.
- [1277] Aaron Tovish, Gerald Schubert, and Bruce P Luyendyk. “Mantle flow pressure and the angle of subduction: Non-Newtonian corner flows”. In: *Journal of Geophysical Research: Solid Earth* 83.B12 (1978), pp. 5892–5898.
- [1278] B.J. Travis et al. “A benchmark comparison of numerical methods for infinite Prandtl number thermal convection in two-dimensional Cartesian geometry”. In: *Geophysical & Astrophysical Fluid Dynamics* 55.3-4 (1990), pp. 137–160.
- [1279] Bryan Travis and Peter Olson. “Convection with internal heat sources and thermal turbulence in the Earth’s mantle”. In: *Geophysical Journal International* 118.1 (1994), pp. 1–19. DOI: 10.1111/j.1365-246X.1994.tb04671.x.
- [1280] Bryan Travis, Peter Olson, and Gerald Schubert. “The transition from two-dimensional to three-dimensional planforms in infinite-Prandtl-number thermal convection”. In: *Journal of Fluid Mechanics* 216 (1990), pp. 71–91. DOI: 10.1017/S0022112090000349.
- [1281] S.J. Trim, J. P. Lowman, and S.L. Butler. “Improving Mass Conservation With the Tracer Ratio Method: Application to Thermochemical Mantle Flows”. In: *Geochem. Geophys. Geosyst.* 22 (2020), e2019GC008799. DOI: 10.1029/2019GC008799.
- [1282] Sean J Trim, Samuel L Butler, and Raymond J Spiteri. “Benchmarking multiphysics software for mantle convection”. In: *Computers & Geosciences* 154 (2021), p. 104797. DOI: 10.1016/j.cageo.2021.104797.

- [1283] R.A. Trompert and U. Hansen. “On the Rayleigh number dependence of convection with a strongly temperature-dependent viscosity”. In: *Physics of Fluids* 10.2 (1998), pp. 351–360.
- [1284] R.A. Trompert and U. Hansen. “The application of a finite volume multigrid method to three-dimensional flow problems in a highly viscous fluid with a variable viscosity”. In: *Geophysical and Astrophysical Fluid Dynamics* 83.3-4 (1996), pp. 261–291. DOI: 10.1080/03091929608208968.
- [1285] Ulrich Trottenberg, Cornelius W. Oosterlee, and Anton Schuller. *Multigrid*. Elsevier, 2001.
- [1286] VP Trubitsyn, IE Rogozhina, and MK Kaban. “On a spectral method of solving the Stokes equation”. In: *Izvestiya, Physics of the Solid Earth* 44.1 (2008), pp. 18–25.
- [1287] Terry E Tullis, Franklin G Horowitz, and Jan Tullis. “Flow laws of polyphase aggregates from end-member flow laws”. In: *Journal of Geophysical Research: Solid Earth* 96.B5 (1991), pp. 8081–8096. DOI: 10.1029/90JB02491.
- [1288] D.L. Turcotte and G. Schubert. *Geodynamics, 2nd edition*. Cambridge, 2012. ISBN: 9780511807442. DOI: 10.1017/CB09780511807442.
- [1289] D.L. Turcotte and G. Schubert. *Geodynamics, 3rd edition*. Cambridge University Press, 2014. ISBN: 9780521186230,
- [1290] DL Turcotte and ER Oxburgh. “Finite amplitude convective cells and continental drift”. In: *Journal of Fluid Mechanics* 28.1 (1967), pp. 29–42. DOI: 10.1017/S0022112067001880.
- [1291] S. Turek. “A comparative study of time-stepping techniques for the incompressible Navier-Stokes equations: from fully implicit non-linear schemes to semi-implicit projection methods”. In: *Int. J. Num. Meth. Fluids* 22 (1996), pp. 987–1011.
- [1292] S. Turek. *Efficient Solvers for Incompressible Flow Problems*. Springer, Berlin, Heidelberg, 1999. ISBN: 978-3-540-65433-9. DOI: 10.1007/978-3-642-58393-3.
- [1293] S. Turek. “Tools for simulating non-stationary incompressible flow via discretely divergence-free finite element models”. In: *Int. J. Num. Meth. Fluids* 18 (1994), pp. 71–105.
- [1294] S. Turek and A. Ouazzi. “Unified edge-oriented stabilization of nonconforming finite element methods for incompressible flow problems”. In: *Journal of Numerical Mathematics* 15.4 (2007).
- [1295] Stefan Turek, Abderrahim Ouazzi, and Rainer Schmachtel. “Multigrid methods for stabilized nonconforming finite elements for incompressible flow involving the deformation tensor formulation”. In: *Journal of Numerical Mathematics* 10.3 (2002), pp. 235–248.
- [1296] L. Uieda, V.C.F. Barbosa, and C. Braitenberg. “Tesseroids: Forward-modeling gravitational fields in spherical coordinates”. In: *Geophysics* 81.5 (2015), pp. 41–48.
- [1297] Leonardo Uieda. “Forward modeling and inversion of gravitational fields in spherical coordinates”. PhD thesis. Observatorio Nacional, 2016.
- [1298] Leonardo Uieda, Everton P Bomfim, Carla Braitenberg, and Eder Molina. “Optimal forward calculation method of the Marussi tensor due to a geologic structure at GOCE height”. In: *Proceedings of the 4th International GOCE User Workshop*. Munich Germany. 2011.
- [1299] H.C. Upadhyaya, O.P. Sharma, R. Mittal, and H. Fatima. “Icosahedral-hexagonal grids on a sphere for CFD applications”. In: *Asia-Pac. J. Chem. Eng.* 6 (2011), pp. 110–119.
- [1300] G de Vahl Davis and IP Jones. “Natural convection in a square cavity: a comparison exercise”. In: *International Journal for numerical methods in fluids* 3.3 (1983), pp. 227–248. DOI: 10.1002/flid.1650030304.

- [1301] JL Valera, Ana M Negredo, and Ivone Jiménez-Munt. “Deep and near-surface consequences of root removal by asymmetric continental delamination”. In: *Tectonophysics* 502.1-2 (2011), pp. 257–265. DOI: 10.1016/j.tecto.2010.04.002.
- [1302] Juan-Luis Valera, Ana-María Negredo, and Antonio Villaseñor. “Asymmetric delamination and convective removal numerical modeling: comparison with evolutionary models for the Alboran Sea region”. In: *Pure appl. geophys.* 165 (2008), pp. 1683–1706. DOI: 10.1007/s00024-008-0395-8.
- [1303] FN Van de Vosse et al. “Finite-element-based computational methods for cardiovascular fluid-structure interaction”. In: *Journal of engineering mathematics* 47.3-4 (2003), pp. 335–368. DOI: 10.1023/B:ENGI.00000007985.17625.43.
- [1304] A.P. Van Den Berg and D.A. Yuen. “The role of shear heating in lubricating mantle flow”. In: *Earth and Planetary Science Letters* 151.1-2 (1997), pp. 33–42. DOI: 10.1016/S0012-821X(97)00110-6.
- [1305] A.P. van den Berg and D.A. Yuen. “Modelling planetary dynamics by using the temperature at the core-mantle boundary as a control variable: effects of rheological layering on mantle heat transport”. In: *Physics of the Earth and Planetary Interiors* 108.3 (1998), pp. 219–234. DOI: 10.1016/S0031-9201(98)00101-0.
- [1306] SP Van Der Pijl, A Segal, and C Vuik. “Modelling of three-dimensional bubbly flows with a mass-conserving Level-Set method”. In: *Proceedings of the 4th European Congress on Computational Methods in Applied Sciences and Engineering, ECCOMAS*. 2004.
- [1307] SP Van der Pijl, A Segal, C Vuik, and P Wesseling. “A mass-conserving level-set method for modelling of multi-phase flows”. In: *International journal for numerical methods in fluids* 47.4 (2005), pp. 339–361. DOI: 10.1002/d.817.
- [1308] SP Van der Pijl, A Segal, C Vuik, and P Wesseling. “Computing three-dimensional two-phase flows with a mass-conserving level set method”. In: *Computing and Visualization in Science* 11.4-6 (2008), pp. 221–235. DOI: 10.1007/s00791-008-0106-0.
- [1309] P.E. van Keken, S.D. King, H. Schmeling, U.R. Christensen, D. Neumeister, and M.-P. Doin. “A comparison of methods for the modeling of thermochemical convection”. In: *J. Geophys. Res.* 102.B10 (1997), pp. 22, 477–22, 495.
- [1310] P.E. van Keken, C.J. Spiers, A.P. van den Berg, and E.J. Muzert. “The effective viscosity of rocksalt: implementation of steady-state creep laws in numerical models of salt diapirism”. In: *Tectonophysics* 225 (1993), pp. 457–476.
- [1311] P.E. van Keken et al. “A community benchmark for subduction zone modelling”. In: *Phys. Earth. Planet. Inter.* 171 (2008), pp. 187–197. DOI: 10.1016/j.pepi.2008.04.015.
- [1312] R.S. Varga. *Matrix Iterative Analysis*. Prentice-Hall, Inc., 1963.
- [1313] Oleg V. Vasilyev, Yuri Yu. Podladchikov, and David A. Yuen. “Modelling of viscoelastic plume-lithosphere interaction using the adaptive multilevel wavelet collocation method”. In: *Geophy. J. Int.* 147 (2001), pp. 579–589.
- [1314] J. Vatteville, P.E. van Keken, A. Limare, and A. Davaille. “Starting laminar plumes: Comparison of laboratory and numerical modeling”. In: *Geochem. Geophys. Geosyst.* 10.12 (2009). DOI: 10.1029/2009GC002739.
- [1315] A. Vauchez, A. Tomassi, and G. Barroul. “Rheological heterogeneity, mechanical anisotropy and deformation of the continental lithosphere”. In: *Tectonophysics* 296 (1998), pp. 61–86. DOI: 10.1016/S0040-1951(98)00137-1.
- [1316] M Velić, L. Moresi, D. May, and M. Knepley. “A Family of Numerically Stable Analytic Solutions for Geodynamic Code Verification”. In: ().

- [1317] Mirko Velić, Dave May, and Louis Moresi. “A fast robust algorithm for computing discrete voronoi diagrams”. In: *Journal of Mathematical Modelling and Algorithms* 8.3 (2009), pp. 343–355. DOI: 10.1007/s10852-008-9097-6.
- [1318] Fabio Verbosio, Arne De Coninck, Drosos Kourounis, and Olaf Schenk. “Enhancing the scalability of selected inversion factorization algorithms in genomic prediction”. In: *Journal of Computational Science* 22.Supplement C (2017), pp. 99–108. ISSN: 1877-7503.
- [1319] Jean Verhoogen. “The adiabatic gradient in the mantle”. In: *Eos, Transactions American Geophysical Union* 32.1 (1951), pp. 41–43.
- [1320] A Verruijt. “Deformations of an elastic half plane with a circular cavity”. In: *International Journal of Solids and Structures* 35.21 (1998), pp. 2795–2804. DOI: 10.1016/S0020-7683(97)00194-7.
- [1321] Valérie Vidal. “Interaction des Différentes échelles de Convection dans le Manteau Terrestre”. PhD thesis. Institut de physique du globe de paris-IPGP, 2004.
- [1322] Jean-Pierre Vilotte, M Daignieres, and Raul Madariaga. “Numerical modeling of intraplate deformation: Simple mechanical models of continental collision”. In: *Journal of Geophysical Research: Solid Earth* 87.B13 (1982), pp. 10709–10728. DOI: 10.1029/JB087iB13p10709.
- [1323] Jean-Pierre Vilotte, Marc Daignieres, R Madariaga, and OC Zienkiewicz. “The role of a heterogeneous inclusion during continental collision”. In: *Physics of the earth and planetary interiors* 36.3-4 (1984), pp. 236–259. DOI: 10.1016/0031-9201(84)90049-9.
- [1324] Jean-Pierre Vilotte, Raul Madariaga, Marc Daignieres, and O Zienkiewicz. “Numerical study of continental collision: influence of buoyancy forces and an initial stiff inclusion”. In: *Geophysical Journal International* 84.2 (1986), pp. 279–310. DOI: 10.1111/j.1365-246X.1986.tb04357.x.
- [1325] A.P. Vincent and D.A. Yuen. “Thermal attractor in chaotic convection with high-Prandtl-number fluids”. In: *Physical Review A* 38.1 (1988), pp. 328–334. DOI: 10.1103/PhysRevA.38.328.
- [1326] C. Vincent and R. Boyer. “A preconditioned conjugate gradient Uzawa-type method for the solution of the Stokes problem by mixed Q1-P0 stabilised finite elements”. In: *Int. J. Num. Meth. Fluids* 14 (1992), pp. 289–298.
- [1327] N.J. Vlaar, P.E. van Keken, and A.P. van den Berg. “Cooling of the Earth in the Archean: Consequences of pressure-release melting in a hotter mantle”. In: *Earth Planet. Sci. Lett.* 121 (1994), pp. 1–18. DOI: 10.1016/0012-821X(94)90028-0.
- [1328] M. von Tscharn and S. M. Schmalholz. “A 3-D Lagrangian finite element algorithm with remeshing for simulating large-strain hydrodynamic instabilities in power law viscoelastic fluids”. In: *Geochem. Geophys. Geosyst.* 16.1 (2015), pp. 215–245. DOI: 10.1002/2014GC005628.
- [1329] H.A. van der Vorst and C. Vuik. “GMRESR: a family of nested GMRES methods”. In: *Num. Lin. Alg. Appl.* 1 (1994), pp. 369–386.
- [1330] F.N. van de Vosse, A.A. van Steenhoven, A. Segal, and J.D. Janssen. “A finite element analysis of the steady laminar entrance flow in a 90° curved tube”. In: *Int. J. Num. Meth. Fluids* 9 (1989), pp. 275–287.
- [1331] C Vuik, A Saghir, and GP Boerstol. “The Krylov accelerated SIMPLE (R) method for flow problems in industrial furnaces”. In: *International Journal for Numerical methods in fluids* 33.7 (2000), pp. 1027–1040.
- [1332] Kees Vuik. “Een kwart eeuw iteratieve methoden”. In: (2009).

- [1333] L. Vynnytska, M.E. Rognes, and S.R. Clark. “Benchmarking FEniCS for mantle convection simulations”. In: *Computers & Geosciences* 50 (2013), pp. 95–105. DOI: 10.1016/j.cageo.2012.05.012.
- [1334] Uwe Walzer and Roland Hendel. “A new convection–fractionation model for the evolution of the principal geochemical reservoirs of the Earth’s mantle”. In: *Physics of the Earth and Planetary Interiors* 112.3-4 (1999), pp. 211–256. DOI: 10.1016/S0031-9201(99)00035-7.
- [1335] Uwe Walzer and Roland Hendel. “Tectonic episodicity and convective feedback mechanisms”. In: *Physics of the earth and planetary interiors* 100.1-4 (1997), pp. 167–188. DOI: 10.1016/S0031-9201(96)03238-4.
- [1336] Bo Wang and BC Khoo. “Hybridizable discontinuous Galerkin method (HDG) for Stokes interface flow”. In: *Journal of Computational Physics* 247 (2013), pp. 262–278. DOI: 10.1016/j.jcp.2013.03.064.
- [1337] H. Wang, R. Agrusta, and J. van Hunen. “Advantages of a conservative velocity interpolation (CVI) scheme for particle-in-cell methods with application in geodynamic modeling”. In: *Geochem. Geophys. Geosyst.* 16 (2015). DOI: 10.1002/2015GC005824.
- [1338] Li Wang and Dimitri J Mavriplis. “Adjoint-based h–p adaptive discontinuous Galerkin methods for the 2D compressible Euler equations”. In: *Journal of Computational Physics* 228.20 (2009), pp. 7643–7661.
- [1339] W.M. Wang, L.J. Sluys, and R. de Borst. “Viscoplasticity for instabilities due to strain softening and strain-rate softening”. In: *Int. J. Num. Meth. Eng.* 40 (1997), pp. 3839–3864.
- [1340] X. Wang and W.K. Liu. “Extended immersed boundary method using FEM and RKPM”. In: *Comput. Methods Appl. Mech. Engrg.* 193 (2004), pp. 1305–1321. DOI: 10.1016/j.cma.2003.12.024.
- [1341] Xinguo Wang, William E Holt, and Attreyee Ghosh. “Joint modeling of lithosphere and mantle dynamics: Evaluation of constraints from global tomography models”. In: *Journal of Geophysical Research: Solid Earth* 120.12 (2015), pp. 8633–8655. DOI: 10.1002/2015JB012188.
- [1342] G.H. Wannier. “A contribution to the hydrodynamics of lubrication”. In: *Quarterly of Applied Mathematics* VIII (1950), pp. 1–32.
- [1343] C.J. Warren, C. Beaumont, and R.A. Jamieson. “Formation and exhumation of ultra-high-pressure rocks during continental collision: Role of detachment in the subduction channel”. In: *Geochem. Geophys. Geosyst.* 9 (2008). DOI: 10.1029/2007GC001839.
- [1344] C.J. Warren, C. Beaumont, and R.A. Jamieson. “Modelling tectonic styles and ultra-high pressure (UHP) rock exhumation during the transition from oceanic subduction to continental collision”. In: *Earth Planet. Sci. Lett.* 267 (2008), pp. 129–145.
- [1345] Clare J Warren, Christopher Beaumont, and Rebecca A Jamieson. “Deep subduction and rapid exhumation: Role of crustal strength and strain weakening in continental subduction and ultrahigh-pressure rock exhumation”. In: *Tectonics* 27.6 (2008).
- [1346] R.F. Weinberg and H. Schmeling. “Polydiapirs: multiwavelength gravity structures”. In: *Journal of Structural Geology* 14.4 (1992), pp. 425–436. DOI: 10.1016/0191-8141(92)90103-4.
- [1347] S.A. Weinstein, P.L. Olson, and D.A. Yuen. “Time-dependent large aspect-ratio thermal convection in the earth’s mantle”. In: *Geophysical & Astrophysical Fluid Dynamics* 47.1-4 (1989), pp. 157–197. DOI: 10.1080/03091928908221820.

- [1348] M. B. Weller and A. Lenardic. “The Energetics and Convective Vigor of Mixed-mode Heating: Velocity Scalings and Implications for the Tectonics of Exoplanets: The Energetics of Mixed-mode convection”. In: *Geophysical Research Letters* 43 (2016), pp. 9469–9474. DOI: 10.1002/2016GL069927.
- [1349] M.B. Weller, A. Lenardic, and W.B. Moore. “Scaling relationships and physics for mixed heating convection in planetary interiors: Isoviscous spherical shells”. In: *J. Geophys. Res.* 121 (2016), pp. 7598–7617. DOI: 10.1002/2016JB013247.
- [1350] R.A. Werner and D.J. Scheeres. “Exterior gravitation of a polyhedron derived and compared with harmonic and mascon gravitation representations of asteroid 4769 Castalia”. In: *Celestial Mechanics and Dynamical Astronomy* 65 (1997), pp. 313–344. DOI: 10.1007/BF00053511.
- [1351] Pieter Wesseling and Cornelis W Oosterlee. “Geometric multigrid with applications to computational fluid dynamics”. In: *Journal of Computational and Applied Mathematics* 128.1-2 (2001), pp. 311–334. DOI: 10.1016/S0377-0427(00)00517-3.
- [1352] D.M. Whipp, C. Beaumont, and J. Braun. “Feeding the ”aneurysm”: Orogen-parallel mass transport into Nanga Parbat and the western Himalayan syntaxis”. In: *J. Geophys. Res.* 119 (2014). DOI: 10.1002/2013JB010929.
- [1353] A Whittaker, MHP Bott, and GD Waghorn. “Stresses and plate boundary forces associated with subduction plate margins”. In: *Journal of Geophysical Research: Solid Earth* 97.B8 (1992), pp. 11933–11944.
- [1354] Mark A Wieczorek and Matthias Meschede. “Shtools: Tools for working with spherical harmonics”. In: *Geochemistry, Geophysics, Geosystems* 19.8 (2018), pp. 2574–2592. DOI: 10.1029/2018GC007529.
- [1355] E. van der Wiel, D.J.J. van Hinsbergen, C. Thieulot, and W. Spakman. “Linking rates of slab sinking to long-term lower mantle flow and mixing”. In: *Earth Planet. Sci. Lett.* 625 (2024), p. 118471. DOI: 10.1016/j.epsl.2023.118471.
- [1356] William SD Wilcock and JA Whitehead. “The Rayleigh-Taylor instability of an embedded layer of low-viscosity fluid”. In: *Journal of Geophysical Research: Solid Earth* 96.B7 (1991), pp. 12193–12200. DOI: 10.1029/91JB00339.
- [1357] M.L. Wilkins. *Computer simulation of dynamic phenomena*. Springer, 1999.
- [1358] K.R. Wilks and N.L. Carter. “Rheology of some continental lower crustal rocks”. In: *Tectonophysics* 182.1-2 (1990), pp. 57–77. DOI: 10.1016/0040-1951(90)90342-6.
- [1359] S.D. Willett. “Dynamic and kinematic growth and change of a Coulomb wedge”. In: *Thrust Tectonics*. Ed. by K.R. McClay. Chapman and Hall, 1992, pp. 19–31.
- [1360] S.D. Willett. “Orogeny and orography: The effects of erosion on the structure of mountain belts”. In: *J. Geophys. Res.* 104.B12 (1999), p. 28957.
- [1361] S. Williams, L. Oliker, R. Vuduc, J. Shalf, K. Yelick, and J. Demmel. “Optimization of Sparse Matrix-Vector Multiplication on Emerging Multicore Platforms”. In: *SC’07 proceedings of the 2007 ACM/IEEE conference on supercomputing*. ACM NY, 2007.
- [1362] Samuel Williams et al. “PERI-auto-tuning memory-intensive kernels for multicore”. In: *Journal of Physics: Conference Series*. Vol. 125. 1. IOP Publishing, 2008, p. 012038. DOI: 10.1088/1742-6596/125/1/012038.
- [1363] Erskine D Williamson and Leason H Adams. “Density distribution in the Earth”. In: *Journal of the Washington Academy of Sciences* 13.19 (1923), pp. 413–428. DOI: xxxx.



- [1364] J Tuzo Wilson. “Did the Atlantic close and then re-open?” In: *Nature* 211.5050 (1966), pp. 676–681. DOI: 10.1038/211676a0.
- [1365] HH Winter. “Viscous dissipation term in energy equations”. In: *Calculation and Measurement Techniques for Momentum, Energy and Mass Transfer* 7 (1987), pp. 27–34. DOI: xxxx.
- [1366] W.-D. Woidt. “Finite element calculations applied to salt dome analysis”. In: *Tectonophysics* 50 (1978), pp. 369–386. DOI: 10.1016/0040-1951(78)90143-9.
- [1367] Marek Wojciechowski. “A note on the differences between Drucker-Prager and Mohr-Coulomb shear strength criteria”. In: *Studia Geotechnica et Mechanica* (2018). DOI: 10.2478/sgem-2018-0016.
- [1368] M. Wolstencroft, J.H. Davies, and D.R. Davies. “Nusselt-Rayleigh number scaling for spherical shell Earth mantle simulation up to a Rayleigh number of  $10^9$ ”. In: *Phys. Earth. Planet. Inter.* 176 (2009), pp. 132–141. DOI: 10.1016/j.pepi.2009.05.002.
- [1369] Jennifer Worthen, Georg Stadler, Noemi Petra, Michael Gurnis, and Omar Ghattas. “Towards adjoint-based inversion for rheological parameters in nonlinear viscous mantle flow”. In: *Physics of the Earth and Planetary Interiors* 234 (2014), pp. 23–34.
- [1370] G.B. Wright, N. Flyer, and D.A. Yuen. “A hybrid radial basis function-pseudospectral method for thermal convection in a 3-D spherical shell”. In: *Geochem. Geophys. Geosyst.* 11.7 (2010). DOI: 10.1029/2009GC002985.
- [1371] B. Wu, Z. Xu, and Z. Li. “A note on imposing displacement boundary conditions in finite element analysis”. In: *Commun. Numer. Meth. Engng* 24 (2008), pp. 777–784.
- [1372] Patrick Wu. “Deformation of an incompressible viscoelastic flat earth with powerlaw creep: a finite element approach”. In: *Geophysical Journal International* 108.1 (1992), pp. 35–51. DOI: 10.1111/j.1365-246X.1992.tb00837.x.
- [1373] Patrick Wu and WR Peltier. “Viscous gravitational relaxation”. In: *Geophysical Journal International* 70.2 (1982), pp. 435–485. DOI: 10.1111/j.1365-246X.1982.tb04976.x.
- [1374] H. Xing, W. Yu, and J. Zhang. “3D Mesh Generation in Geocomputing”. In: *Advances in Geocomputing, Lecture Notes in Earth Sciences*. Berlin Heidelberg: Springer-Verlag, 2009. DOI: 10.1007/978-3-540-85879-9\_2.
- [1375] Dongbin Xiu and George Em Karniadakis. “A semi-Lagrangian high-order method for Navier-Stokes equations”. In: *Journal of computational physics* 172.2 (2001), pp. 658–684. DOI: 10.1006/jcph.2001.6847.
- [1376] Y. Lina nd Y. Cao. “A new nonlinear Uzawa algorithm for generalised saddle point problems”. In: *Applied Mathematics and Computation* 175 (2006), pp. 1432–1454.
- [1377] P. Yamato, L. Husson, J. Braun, C. Loiselet, and C. Thieulot. “Influence of surrounding plates on 3D subduction dynamics”. In: *Geophys. Res. Lett.* 36.L07303 (2009). DOI: 10.1029/2008GL036942.
- [1378] Haibin Yang, Louis Moresi, and John Mansour. “Stress Recovery for the Particle-in-cell Finite Element Method”. In: *Phys. Earth. Planet. Inter.* 311 (2021), p. 106637. DOI: 10.1016/j.pepi.2020.106637.
- [1379] S. Yang and Y. Shi. “Three-dimensional numerical simulation of glacial trough forming process”. In: *Science China: Earth Sciences* (2015). DOI: 10.1007/s11430-015-5120-8.
- [1380] Ting Yang, Louis Moresi, Dapeng Zhao, Dan Sandiford, and Joanne Whittaker. “Cenozoic lithospheric deformation in Northeast Asia and the rapidly-aging Pacific Plate”. In: *Earth Planet. Sci. Lett.* 492 (2018), pp. 1–11. DOI: 10.1016/j.epsl.2018.03.057.

- [1381] Woo-Sun Yang and John R Baumgardner. “A matrix-dependent transfer multigrid method for strongly variable viscosity infinite Prandtl number thermal convection”. In: *Geophysical & Astrophysical Fluid Dynamics* 92.3-4 (2000), pp. 151–195. DOI: 10.1080/03091920008203715.
- [1382] K. Yasuda, R.C. Armstrong, and R.E. Cohen. “Shear flow properties of concentrated solutions of linear and star branched polystyrenes”. In: *Rheol. Acta* 20 (1981), pp. 163–178. DOI: 10.1007/BF01513059.
- [1383] Irad Yavneh. “Why multigrid methods are so efficient”. In: *Computing in science & engineering* 8.6 (2006), p. 12. DOI: 10.1109/MCSE.2006.125.
- [1384] Tao Ye, Rajat Mittal, HS Udaykumar, and Wei Shyy. “An accurate Cartesian grid method for viscous incompressible flows with complex immersed boundaries”. In: *Journal of computational physics* 156.2 (1999), pp. 209–240. DOI: 10.1006/jcph.1999.6356.
- [1385] Liang Yin, Chao Yang, Shi-Zhuang Ma, and Ke-Ke Zhang. “Parallel and fully implicit simulations of the thermal convection in the Earth’s outer core”. In: *Computers & Fluids* 193 (2019), p. 104278. DOI: 10.1016/j.compfluid.2019.104278.
- [1386] He Yiqian and Yang Haitian. “Solving inverse couple-stress problems via an element-free Galerkin (EFG) method and Gauss–Newton algorithm”. In: *Finite Elements in Analysis and Design* 46.3 (2010), pp. 257–264. DOI: 10.1016/j.finel.2009.09.009.
- [1387] M. Yoshida and A. Kageyama. “Application of the Yin-Yang grid to a thermal convection of a Boussinesq fluid with infinite Prandtl number in a three-dimensional spherical shell”. In: *Geophys. Res. Lett.* 31.L12609 (2004). DOI: 10.1029/2004GL019970.
- [1388] Masaki Yoshida, Satoru Honda, Motoyuki Kido, and Yasuyuki Iwase. “Numerical simulation for the prediction of the plate motions”. In: *Earth, planets and space* 53.7 (2001), pp. 709–721. DOI: 10.1186/BF03352399.
- [1389] Masaki Yoshida and Akira Kageyama. “Low-degree mantle convection with strongly temperature- and depth-dependent viscosity in a three-dimensional spherical shell”. In: *Journal of Geophysical Research: Solid Earth* 111.B3 (2006). DOI: 10.1029/2005JB003905.
- [1390] S. Yoshioka and M.J.R. Wortel. “Three-dimensional numerical modeling of detachment of subducted lithosphere”. In: *J. Geophys. Res.* 100.B10 (1995), pp. 20, 223–20, 244.
- [1391] DL Youngs. *Numerical methods for fluid dynamics*, ed. by KW Morton and MJ Baines. Academic Press, Massachusetts, 1982.
- [1392] Hongzheng Yu and Shimin Wang. “Unified linear stability analysis for thermal convections in spherical shells under different boundary conditions and heating modes”. In: *Earth and Space Science* (2019).
- [1393] X. Yu and F. Tin-Loi. “A simple mixed finite element for static limit analysis”. In: *Computers and Structures* 84 (2006), pp. 1906–1917.
- [1394] K.Y. Yuan, Y.S. Huang, H.T. Yang, and T.H.H. Pian. “The inverse mapping and distortion measures for 8-node hexahedral isoparametric elements”. In: *Computational Mechanics* 14 (1994), pp. 189–199.
- [1395] D.A. Yuen and W.R. Peltier. “Mantle plumes and the thermal stability of the D” layer”. In: *Geophysical Research Letters* 7.9 (1980), pp. 625–628. DOI: 10.1029/GL007i009p00625.
- [1396] D.A. Yuen, W.R. Peltier, and G. Schubert. “On the existence of a second scale of convection in the upper mantle”. In: *Geophysical Journal of the Royal Astronomical Society* 65.1 (1981), pp. 171–190. DOI: 10.1111/j.1365-246X.1981.tb02707.x.

- [1397] David A Yuen, Ctirad Matyska, Ondrej Cadek, and Masanori Kameyama. “The dynamical influences from physical properties in the lower mantle and post-perovskite phase transition”. In: *Geophysical Monograph - American Geophysical Union* 174 (2007), p. 249. DOI: xxxx.
- [1398] David A Yuen, DM Reuteler, S Balachandar, V Steinbach, AV Malevsky, and JJ Smedsmo. “Various influences on three-dimensional mantle convection with phase transitions”. In: *Physics of the Earth and Planetary interiors* 86.1-3 (1994), pp. 185–203. DOI: 10.1016/0031-9201(94)05068-6.
- [1399] S.T. Zalesak. “Fully Multidimensional Flux-Corrected Transport Algorithms for Fluids”. In: *J. Comp. Phys.* 31 (1979), pp. 335–362. DOI: 10.1016/0021-9991(79)90051-2.
- [1400] S. Zaleski and P. Julien. “Numerical simulation of Rayleigh-Taylor instability for single and multiple salt diapirs”. In: *Tectonophysics* 206 (1992), pp. 55–69. DOI: 10.1016/0040-1951(92)90367-F.
- [1401] A. Zebib. “Linear and Weakly Nonlinear Variable Viscosity Convection in Spherical Shells”. In: *Theoret. Comput. Fluid Dynamics* 4 (1993), pp. 241–253.
- [1402] Abdelfattah Zebib, Gerald Schubert, and Joe M Straus. “Infinite Prandtl number thermal convection in a spherical shell”. In: *Journal of Fluid Mechanics* 97.2 (1980), pp. 257–277. DOI: 10.1017/S0022112080002558.
- [1403] Iris van Zelst, Fabio Cramer, Adina E. Pusok, Anne Glerum, Julianne Dannberg, and Cedric Thieulot. “101 Geodynamic modelling: How to design, interpret, and communicate numerical studies of the solid Earth”. In: 13 (2022), pp. 583–637. DOI: 10.5194/se-13-583-2022.
- [1404] Guoxiang Zhang and Junyu Xiang. “Eight-node conforming straight-side quadrilateral element with high-order completeness (QH8-C1)”. In: *International Journal for Numerical Methods in Engineering* (2020). DOI: 10.1002/nme.6360.
- [1405] Huai Zhang, Lili Ju, Max Gunzburger, Todd Ringler, and Stephen Price. “Coupled models and parallel simulations for three-dimensional full-Stokes ice sheet modeling”. In: *Numerical Mathematics: Theory, Methods and Applications* 4.3 (2011), pp. 396–418. DOI: 10.1017/S1004897900000416.
- [1406] N. Zhang, S. Zhong, and A.K. McNamara. “Supercontinent formation from stochastic collision and mantle convection models”. In: *Gondwana Research* 15 (2009), pp. 267–275. DOI: 10.1016/j.gr.2008.10.002.
- [1407] L. Zhao, X. Bai, T. Li, and J.J.R. Williams. “Improved conservative level set method”. In: *Int. J. Num. Meth. Fluids* (2014). DOI: 10.1002/flid.3907.
- [1408] S. Zhong and M. Gurnis. “Role of plates and temperature-dependent viscosity in phase change dynamics”. In: *Journal of Geophysical Research* 99.B8 (1994), p. 15903. DOI: 10.1029/94JB00545.
- [1409] S. Zhong, M. Gurnis, and G. Hulbert. “Accurate determination of surface normal stress in viscous flow from a consistent boundary flux method”. In: *Physics of the Earth and Planetary Interiors* 78.1-2 (1993), pp. 1–8. DOI: 10.1016/0031-9201(93)90078-N.
- [1410] S. Zhong, M. Gurnis, and L. Moresi. “Free-surface formulation of mantle convection-I. Basic theory and application to plumes”. In: *Geophysical Journal International* 127.3 (1996), pp. 708–718. DOI: 10.1111/j.1365-246X.1996.tb04049.x.
- [1411] S. Zhong, M. Gurnis, and L. Moresi. “Role of faults, nonlinear rheology, and viscosity structure in generating plates from instantaneous mantle flow models”. In: *J. Geophys. Res.* 103.B7 (1998), pp. 15, 255–15, 268. DOI: 10.1029/98JB00605.

- [1412] S. Zhong, A. McNamara, E. Tan, L. Moresi, and M. Gurnis. “A benchmark study on mantle convection in a 3-D spherical shell using CitcomS”. In: *Geochem. Geophys. Geosyst.* 9.10 (2008). DOI: 10.1029/2008GC002048.
- [1413] S. Zhong and A.B. Watts. “Lithospheric deformation induced by loading of the Hawaiian Islands and its implications for mantle rheology”. In: *J. Geophys. Res.* 118 (2013), pp. 6025–6048. DOI: 10.1002/2013JB010408.
- [1414] S. Zhong, M.T. Zuber, L.N. Moresi, and M. Gurnis. “The role of temperature-dependent viscosity and surface plates in spherical shell models of mantle convection”. In: *J. Geophys. Res.* 105.B5 (2000), pp. 11, 063–11, 082. DOI: 10.1029/2000JB900003.
- [1415] S.J. Zhong, D.A. Yuen, L.N. Moresi, and M.G. Knepley. “7.05 - Numerical Methods for Mantle Convection”. In: *Treatise on Geophysics (Second Edition)*. Ed. by Gerald Schubert. Second Edition. Oxford, 2015, pp. 197–222. DOI: 10.1016/B978-0-444-53802-4.00130-5.
- [1416] Shijie Zhong. “Analytic solutions for Stokes’ flow with lateral variations in viscosity”. In: *Geophys. J. Int.* 124.1 (1996), pp. 18–28. DOI: 10.1111/j.1365-246X.1996.tb06349.x.
- [1417] S-S Zhou and X-L Gao. “Solutions of the generalized half-plane and half-space Cerruti problems with surface effects”. In: *Zeitschrift für angewandte Mathematik und Physik* 66.3 (2015), pp. 1125–1142. DOI: 10.1007/s00033-014-0419-4.
- [1418] D.Y. Zhu, C.F. Lee, and K.T. Law. “Determination of bearing capacity of shallow foundations without using superposition approximation”. In: *Can. Geotech. J.* 40 (2003), pp. 450–459. DOI: 10.1139/t02-105.
- [1419] G. Zhu, T.V. Gerya, P.J. Tackley, and E. Kissling. “Four-dimensional numerical modeling of crustal growth at active continental margins”. In: *J. Geophys. Res.* 118 (2013), pp. 4682–4698. DOI: 10.1002/jgrb.50357.
- [1420] Yongning Zhu and Robert Bridson. “Animating sand as a fluid”. In: *ACM Transactions on Graphics (TOG)* 24.3 (2005), pp. 965–972. DOI: 10.1145/1073204.1073298.
- [1421] O. Zienkiewicz and S. Nakazawa. “The penalty function method and its application to the numerical solution of boundary value problems”. In: *The American Society of Mechanical Engineers* 51 (1982). DOI: xxxx.
- [1422] O.C. Zienkiewicz. “Visco-plasticity, plasticity, creep and visco-plastic flow”. In: *Lecture Notes in Mathematics* 461 (1975), pp. 297–328.
- [1423] O.C. Zienkiewicz and I.C. Corneau. “Visco-plasticity and creep in elastic solids - A unified numerical solution approach”. In: *Int. J. Num. Meth. Eng.* 8 (1974), pp. 821–845. DOI: 10.1002/nme.1620080411.
- [1424] O.C. Zienkiewicz and I.C. Corneau. “Visco-plasticity and plasticity. An alternative for Finite Element solution of material nonlinearities”. In: *Lecture notes in Computer Science* 10 (1974), pp. 259–287. DOI: 10.1007/BFb0015179.
- [1425] O.C. Zienkiewicz and P.N. Godbole. “Flow of plastic and visco-plastic solids with special reference to extrusion and forming processes”. In: *Int. J. Num. Meth. Eng.* 8 (1974), pp. 3–16. DOI: 10.1002/nme.1620080102.
- [1426] O.C. Zienkiewicz, M. Huang, and M. Pastor. “Localization problems in plasticity using Finite Elements with adaptive remeshing”. In: *International Journal for Numerical and Analytical Methods in Geomechanics* 19 (1995), pp. 127–148. DOI: 10.1002/nag.1610190205.
- [1427] O.C. Zienkiewicz, C. Humpheson, and R.W. Lewis. “Associated and non-associated visco-plasticity and plasticity in soil mechanics”. In: *Géotechnique* 25.4 (1975), pp. 671–689. DOI: 10.1680/geot.1975.25.4.671.

- [1428] O.C. Zienkiewicz, P.C. Jain, and E. Oñate. “Flow of solids during forming and extrusion: some aspects of numerical solutions”. In: *Int. J. Solids Structures* 14 (1978), pp. 15–38. DOI: 10.1016/0020-7683(78)90062-8.
- [1429] O.C. Zienkiewicz, M. Pastor, and Maosong Huang. “Softening, localisation and adaptive remeshing. Capture of discontinuous solutions”. In: *Computational Mechanics* 17 (1995), pp. 98–106. DOI: 10.1007/BF00356482.
- [1430] O.C. Zienkiewicz and R.L. Taylor. *The Finite Element Method. Vol. 1: The basis*. Butterworth and Heinemann, 2002. ISBN: 0-7506-5049-4.
- [1431] O.C. Zienkiewicz and R.L. Taylor. *The Finite Element Method. Vol. 2: Solid Mechanics*. Butterworth and Heinemann, 2002.
- [1432] O.C. Zienkiewicz and R.L. Taylor. *The Finite Element Method. Vol. 3: Fluid Dynamics*. Butterworth and Heinemann, 2002.
- [1433] O.C. Zienkiewicz, R.L. Taylor, and D.D. Fox. *The Finite Element Method for solid and structural mechanics*. Elsevier B.H., 2014.
- [1434] O.C. Zienkiewicz, J.P. Vilotte, and S. Toyoshima. “Iterative method for constrained and mixed approximation. An inexpensive improvement of FEM performance”. In: *Computer Methods in Applied Mechanics and Engineering* 51 (1985), pp. 3–29. DOI: 10.1016/0045-7825(85)90025-8.
- [1435] O.C. Zienkiewicz and J. Wu. “Incompressibility without tears - How to avoid restrictions of mixed formulations”. In: *Int. J. Num. Meth. Eng.* 32 (1991), pp. 1189–1203. DOI: 10.1002/nme.1620320603.
- [1436] O.C. Zienkiewicz and J.Z. Zhu. “The superconvergent patch recovery and a posteriori error estimates. Part 1: the recovery technique”. In: *Int. J. Num. Meth. Eng.* 33 (1992), pp. 1331–1364.
- [1437] O.C. Zienkiewicz and J.Z. Zhu. “The superconvergent patch recovery and a posteriori error estimates. Part 2: error estimates and adaptativity”. In: *Int. J. Num. Meth. Eng.* 33 (1992), pp. 1365–1382.
- [1438] OC Zienkiewicz, B Boroomand, and Jian Zhong Zhu. “Recovery procedures in error estimation and adaptivity: adaptivity in linear problems”. In: *Advances in Adaptive Computational Methods in Mechanics*. 1998, pp. 3–23.
- [1439] F. Zinani and S. Frey. “Galerkin Least-Squares Solutions for Purely Viscous Flows of Shear-Thinning Fluids and Regularized Yield Stress Fluids”. In: *J. of the Braz. Soc. of Mech. Sci. & Eng.* XXIX (2007), pp. 432–443.
- [1440] S. Zlotnik, P. Diez, M. Fernandez, and J. Verges. “Numerical modelling of tectonic plates subduction using X-FEM”. In: *Comput. Methods Appl. Mech. Engrg* 196 (2007), pp. 4283–4293. DOI: 10.1016/j.cma.2007.04.006.
- [1441] S. Zlotnik, M. Fernandez, P. Diez, and J. Verges. “Modelling gravitational instabilities: slab break-off and Rayleigh-Taylor diapirism”. In: *Pure appl. geophys.* 165 (2008), pp. 1491–1510. DOI: 10.1007/s00024-004-0386-9.
- [1442] M.T. Zuber and E.M. Parmentier. “Lithospheric necking: a dynamic model for rift morphology”. In: *Earth Planet. Sci. Lett.* 77 (1986), pp. 373–383. DOI: 10.1016/0012-821X(86)90147-0.
- [1443] M.T. Zuber, E.M. Parmentier, and R.C. Fletcher. “Extension of Continental Lithosphere: A Model for Two Scales of Basin and Range Deformation”. In: *J. Geophys. Res.* 91.B5 (1986), pp. 4826–4838. DOI: 10.1029/JB091iB05p04826.

- [1444] T. Zwinger, R. Greve, O. Gagliardini, T. Shiraiwa, and M. Lyly. “A full Stokes-flow thermo-mechanical model for firn and ice applied to the Gorshkov crater glacier, Kamchatka”. In: *Annals of Glaciology* 45 (2007), pp. 29–37.

# General index

- $H^1$  norm, 561
- $H^1$  semi-norm, 561
- $H^1(\Omega)$  space, 561
- $L_1$  norm, 561
- $L_2$  norm, 561
- $P_1$ , 229, 232, 260, 264
- $P_1 \times P_0$ , 341
- $P_1^+$ , 233, 261
- $P_1^{NC}$ , 258
- $P_2$ , 235
- $P_2^+$ , 237, 263
- $P_3$ , 239
- $P_3 \times P_2$  element, 346
- $P_4$ , 243
- $P_m \times P_n$ , 203
- $P_m \times P_{-n}$ , 203
- $QH8 - C1$ , 228
- $Q_1$ , 205, 220
- $Q_1^+$ , 248
- $Q_2 \times Q_1$ , 203
- $Q_2$ , 206, 222, 262
- $Q_2^{(20)}$ , 266
- $Q_2^{(8)}$ , 226
- $Q_3$ , 207, 224
- $Q_4$ , 210, 229
- $Q_5$ , 212
- $Q_6$ , 213
- $Q_m \times P_{-n}$ , 204
- $Q_m \times Q_n$ , 203
- $Q_m \times Q_{-n}$ , 203
- $\bar{Q}_1$ , 254, 266
- $\bar{Q}_1 \times P_0$ , 340
- $Q_1 \times P_0$ , 326
- $\bar{Q}_2 \times Q_0$ , 342
- Adaptive Mesh Refinement, 521
- Advection-Diffusion Equation, 301
- ALA, 47
- ALE, 482
- AMR, 521
- Anelastic Liquid Approximation, 47
- Angular Momentum, 444
- Angular Velocity, 444
- Arbitrary Lagrangian Eulerian, 482
- Arrhenius law, 118
- Associated Legendre polynomials, 724
- Augmented Lagrangian, 421
- Average Operator, 470
- Backward Euler, 307
- Backward FD Derivative, 154
- Bandwidth Reduction, 666
- Barycentric Coordinates, 233
- basis functions, 219
- BDF-2, 308
- BDM element, 347
- Bernstein Polynomials, 616
- BiCG solver, 162, 660
- Biharmonic Operator, 595, 956
- Bingham model, 92
- Bird-Carreau model, 91
- Bird-Carreau-Yasuda model, 92
- Boussinesq Approximation, 60
- Bow-tied element, 479
- Bubble Function, 204, 233, 315
- Bulk Modulus, 62, 990
- Bulk Viscosity, 45
- Cam-clay Failure Criterion, 114
- Capacitance Matrix, 296
- Carreau model, 91
- CG, 415
- CG solver, 162, 659
- Chain Rule, 569
- Checkerboard mode, 535
- Cholesky decomposition, 162
- Colatitude, 37
- Compliance Matrix, 61
- Compositional Field, 582
- Compositional Rayleigh Number, 68
- Compressed Sparse Column, 499
- Compressed Sparse Row, 499
- Conformal Mesh Refinement, 525
- Conforming Element, 204
- Conjugate Gradient, 415
- Connectivity Array, 515
- Continuity Equation, 44
- Convex Polygon, 514
- Courant-Friedrichs-Lewy Condition, 545
- Crank-Nicolson, 164, 307
- Crank-Nicolson Method, 176
- Critical Rayleigh Number, 68
- Crouzeix-Raviart, 237
- CSC, 499
- CSR, 499

Defect Correction Formulation, 590  
 Differential Curvature Magnitude, 741  
 Diffusion creep, 95  
 Dirichlet Boundary Condition, 66  
 Dislocation creep, 95  
 Dissipation Number, 48  
 Distributive Iterative Method, 423  
 Divergence Operator in Cartesian Coordinates, 33  
 Divergence Operator in Cylindrical Coordinates, 35  
 Divergence Theorem, 569  
 Divergence-free, 44  
 Divergence-free Flow, 322  
 Domain Decomposition, 592  
 Drucker-Prager, 107  
 Drunken Sailor, 489  
 DSSY element, 347  
 Dynamic Viscosity, 45  
  
 EBA, 52, 60  
 ENO, 581  
 Equation of State, 46  
 Essential Boundary Conditions, 430  
 Extended Boussinesq Approximation, 52, 60  
 Extrapolation, 562  
  
 Flow Index, 93  
 Forward Euler, 307  
 Forward FD Derivative, 152  
 Frank-Kamenetskii, 118  
 Free Surface, 479  
 Free Surface Stabilisation Algorithm, 488  
 FSSA, 488  
  
 Gauss-Legendre Quadrature, 186  
 Gauss-Lobatto, 200  
 Gauss-Seidel Iterative Method, 414  
 Gauss-Seidel solver, 162, 659  
 Generalized Newtonian Fluid, 91  
 Geoid, 603  
 Geometric Multigrid, 599  
 GMRES solver, 162, 659  
 Gradient Operator in Cartesian Coordinates, 33  
 Gradient Operator in Cylindrical Coordinates, 35  
 Gradient Tensor, 477  
 Gradient-Based Formulation (DG-FEM), 477  
 Griffith-Murrell, 114  
 Haigh-Westergaard Coordinates, 76  
 Han element, 350  
 Heat Diffusivity, 157  
 Herschel-Bulkley model, 93  
 Heun's emthod, 552  
 Horizontal Gradient Magnitude, 741  
 hyperbolic PDE, 545  
  
 Isoparametric, 377  
 Isotropic, 310  
  
 Jacobi Iterative Method, 414  
 Jacobi solver, 162, 659  
 Jump Operator, 470  
  
 Korn's inequality, 340  
  
 Lamé Parameter, 61  
 Laplace Operator, 54, 57  
 Laplace Operator in Cartesian Coordinates, 33  
 Laplace Operator in Cylindrical Coordinates, 35  
 Laplacian, 37, 54, 57  
 Lax-Friedrichs flux, 458  
 Lax-Friedrichs Method, 176  
 Lax-Friedrichs method, 166  
 Lax-Wendroff method, 166  
 LBB, 324  
 Leapfrog method, 166  
 Legendre Polynomial, 186  
 Level-set Function, 581  
 Level-set Method, 581  
 Load Balancing, 593  
 Lode Angle, 75, 87  
 Lode Coordinates, 76  
 Lode Parameter, 76  
 LSF, 581  
 LSM, 581  
 LU decomposition, 162  
 Lyapunov Time, 606  
  
 MAC, 575  
 Marker Chain method, 584  
 Marker-and-Cell, 575  
 Marussi Tensor, 738  
 Mass Conservation Equation, 44  
 Mass Conservation Equation (Cylindrical Coordinates), 44  
 Mass Conservation Equation (Spherical Coordinates), 44  
 Mass Matrix, 296  
 Maximum Shear Stress, 71  
 Maxwell Time, 612



Meshless, 514  
 Method of Manufactured Solutions, 798  
 Midpoint Method, 552  
 Midpoint Rule, 184  
 MINI element, 342  
 Mixed Formulation, 360  
 MMS, 798  
 Mohr-Coulomb, 104  
 Moment Invariant, 78, 1187  
 Moment of Inertia, 444  
  
 Natural Boundary Conditions, 430  
 Neumann Boundary Condition, 66  
 Newton's method, 558  
 Newton-Cotes, 186  
 Newtonian fluid, 90  
 Non-Conforming Element, 204  
 Nonconforming element, 340, 347, 350  
 Nonlinear PDE, 589  
 Normal Stress, 40  
 Nullspace, 443  
 Nusselt Number, 68  
  
 Optimal Rate, 324  
 optimal rate, 562  
 Orthotropic, 310  
  
 P.R.E.M., 674, 1023  
 Particle-in-Cell, 575  
 Path Increment in Cartesian Coordinates, 33  
 Path Increment in Cylindrical Coordinates, 35  
 Path Increment in Polar Coordinates, 34  
 Peclet Number, 181, 315  
 Peierls creep, 117  
 Penalty Formulation, 355, 368  
 Periodic Boundary Conditions, 431  
 Perzyna Model, 135  
 PIC, 575  
 Picard Iterations, 589  
 Piecewise, 203  
 Plain Strain, 81  
 Plastic Hardening, 495  
 Poiseuille flow, 807  
 Poisson Ratio, 62, 991  
 Power Law Rheology, 91  
 Prandtl Number, 68  
 Prandtl number, 48  
 Preconditioned Conjugate Gradient, 418  
 Pressure Mass Matrix, 368  
 pressure nullspace, 808  
 pressure scaling, 370  
  
 Pressure Smoothing, 535  
 Principal Invariant, 1187  
 Principal Invariants, 74  
 Principal Stress, 71  
 Principal Stress Ratio, 86  
 Prism, 739  
  
 Quadrature, 186  
  
 Rannacher-Turek element, 340  
 Rayleigh Number, 67  
 Rectangle Rule, 184  
 reference element, 379  
 Relaxation, 589  
 Richardson Iterations, 413  
 RK2, 552  
 RK3, 552  
 RK4, 552  
 RK45, 553  
 Runge-Kutta-Fehlberg method, 553  
  
 Schur Complement, 413  
 Second Viscosity, 45  
 Serendipity element, 226, 228, 266  
 Shear Heating, 807  
 Shear Modulus, 61  
 Shear Stress, 40  
 SIMPLE, 423  
 Solenoidal Field, 44, 322  
 SOR Iterative Method, 414  
 SOR iterative method, 659  
 Space Filling Curve, 593  
 Sparse Matrix-Vector Multiplication, 509  
 SPD, 413  
 Spin Tensor, 41  
 SpMV, 509  
 SSOR Iterative Method, 414  
 SSOR iterative method, 659  
 SSOR solver, 162  
 Static Condensation, 342, 587  
 Sticky Air, 480  
 Strain Rate, 41  
 Strain rate partitioning, 120  
 Strain Rate Tensor, 67  
 Strain Tensor, 62, 67, 988  
 Stream Function, 595, 953  
 Stress Tensor, 40, 67  
 Stress Tensor (Cylindrical Coordinates), 41  
 Stress Tensor (Spherical Coordinates), 41  
 Stress Vector, 40  
 Strong Form, 293

- strong form, 323
- Structured Grid, 514
- Subparametric, 377
- Superparametric, 377
- Symmetric Positive Definite, 1046
  
- Taylor-Hood, 325
- Tensor Invariant, 78
- Tesseract, 745
- Total Gradient Magnitude, 741
- Traction, 40
- Trapezoidal Rule, 184
- Tresca, 102
- Truncation Error, 152
  
- Unstructured Grid, 514
- Uzawa algorithm, 412
  
- Validation, 148
- Velocity Gradient, 41
- Verification, 148
- VOF, 584
- Volume of a Sphere, 37
- Volume of a Spherical shell, 38
- Volume-of-Fluid Method, 584
- von Mises, 99
  
- Weak Form, 357
- weak form, 323
  
- X-FEM, 343
  
- Young's Modulus, 62, 991

# Contributors

A. Hendrickx, 293, 833

B. Myhill, 412  
B. Root, 497, 599, 618

D. Bonté, 507  
D. Duck, 955

E. Hoogen, 65  
E. van der Wiel, 286, 830, 888  
E.G.P. Puckett, 47, 98

F. Garel, 789  
F. Gueydan, 767, 880

I. van Zelst, 368, 484

J. Austermann, 497  
J. Jansen, 19, 26, 295, 456, 707  
J. Mos, 19  
J. Naliboff, 261  
J. Veenhof, 882  
J. Wolbers, 817

L. van de Wiel, 188, 195, 300, 717, 772  
L. Verbrugh, 591

M. Blasweiler, 599, 817

N. Ribe, 286, 823, 826

P. Pitard, 658

R. Elbertsen, 147, 183, 388, 497  
R. Hassani, 103, 511  
R. Maguire, 497  
R. Meghezi, 839, 841

S. Hassing, 489

T. Broerse, 550, 797  
T. Shinohara, 306, 328  
T. Weir, 110, 113  
Th. Sanders, 136

W. Klessens, 707

Z. Erdős, 641, 653