

The Finite Element Method in Geodynamics

C. Thieulot

July 10, 2019

Contents

1	Introduction	8
1.1	Philosophy	8
1.2	Acknowledgments	8
1.3	Essential literature	8
1.4	Installation	8
1.5	What is a fieldstone?	9
1.6	Why the Finite Element method?	9
1.7	Notations	9
1.8	Colour maps for visualisation	9
2	List of tutorials	10
3	The physical equations	12
3.1	Strain rate and spin tensor	12
3.2	The heat transport equation - energy conservation equation	12
3.3	The momentum conservation equations	13
3.4	The mass conservation equations	13
3.5	The equations in ASPECT manual	14
3.6	The Boussinesq approximation: an Incompressible flow	15
3.7	Stokes equation for elastic medium	16
3.8	The strain rate tensor in all coordinate systems	17
3.8.1	Cartesian coordinates	17
3.8.2	Polar coordinates	17
3.8.3	Cylindrical coordinates	17
3.8.4	Spherical coordinates	17
3.9	Boundary conditions	18
3.9.1	The Stokes equations	18
3.9.2	The heat transport equation	18
3.10	Meaningful physical quantities	19
3.11	The need for numerical modelling	21
4	The building blocks of the Finite Element Method	22
4.1	Numerical integration	22
4.1.1	in 1D - theory	22
4.1.2	in 1D - examples	24
4.1.3	in 2D/3D - theory	25
4.1.4	quadrature on triangles	25
4.1.5	quadrature on tetrahedra	26
4.2	The mesh	27
4.3	A bit of FE terminology	27
4.4	Elements and basis functions in 1D	28
4.4.1	Linear basis functions (Q_1)	28
4.4.2	Quadratic basis functions (Q_2)	29

4.4.3	Cubic basis functions (Q_3)	30
4.5	Elements and basis functions in 2D	32
4.5.1	Bilinear basis functions in 2D (Q_1)	33
4.5.2	Biquadratic basis functions in 2D (Q_2)	35
4.5.3	Eight node serendipity basis functions in 2D ($Q_2^{(8)}$)	36
4.5.4	Bicubic basis functions in 2D (Q_3)	37
4.5.5	Linear basis functions for triangles in 2D (P_1)	38
4.5.6	Enriched linear basis functions in triangles (P_1^+)	38
4.5.7	Quadratic basis functions for triangles in 2D (P_2)	41
4.5.8	Enriched quadratic basis functions in triangles (P_2^+)	42
4.5.9	Cubic basis functions for triangles (P_3)	44
4.6	Elements and basis functions in 3D	44
4.6.1	Linear basis functions in tetrahedra (P_1)	44
4.6.2	Enriched linear in tetrahedra(P_1^+)	45
4.6.3	Triquadratic basis functions in 3D (Q_2)	47
5	Solving the heat transport equation with linear Finite Elements	48
5.1	The diffusion equation in 1D	48
5.2	The advection-diffusion equation in 1D	55
5.3	The advection-diffusion equation in 2D	57
5.3.1	Dealing with the time discretisation	58
6	Solving the flow equations with the FEM	61
6.1	Strong and weak forms	61
6.2	Which velocity-pressure pair for Stokes?	61
6.2.1	The compatibility condition (or LBB condition)	61
6.2.2	Families	61
6.2.3	The bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$)	62
6.2.4	The bi/tri-quadratic velocity - discontinuous linear pressure element ($Q_2 \times P_{-1}$)	62
6.2.5	The bi/tri-quadratic velocity - bi/tri-linear pressure element ($Q_2 \times Q_1$)	62
6.2.6	The stabilised bi/tri-linear velocity - bi/tri-linear pressure element ($Q_1 \times Q_1\text{-stab}$)	63
6.2.7	The MINI triangular element ($P_1^+ \times P_1$) in 2D	63
6.2.8	The quadratic velocity - linear pressure triangle ($P_2 \times P_1$)	64
6.2.9	The Crouzeix-Raviart triangle ($P_2^+ \times P_{-1}$)	64
6.2.10	The Rannacher-Turek element - rotated $Q_1 \times P_0$	64
6.2.11	Other elements	64
6.3	The penalty approach for viscous flow	65
6.4	The mixed FEM for viscous flow	68
6.4.1	in three dimensions	68
6.4.2	Going from 3D to 2D	73
6.5	Solving the elastic equations	75
6.6	A quick tour of similar literature	75
6.7	The case against the $Q_1 \times P_0$ element	75
7	Additional techniques and features	77
7.1	Dealing with a free surface	77
7.2	Convergence criterion for nonlinear iterations	77
7.3	The SUPG formulation for the energy equation	78
7.3.1	Linear elements	78
7.4	The method of manufactured solutions	80
7.4.1	Analytical benchmark I - "Donea & Huerta"	80
7.4.2	Analytical benchmark II - "Dohrmann & Bochev 2D"	81
7.4.3	Analytical benchmark III - "Dohrmann & Bochev 3D"	82
7.4.4	Analytical benchmark IV - "Bercovier & Engelman"	83
7.4.5	Analytical benchmark V - "VJ"	83
7.4.6	Analytical benchmark VI - "Ilincă & Pelletier"	84

7.4.7	Analytical benchmark VII - "grooves"	85
7.4.8	Analytical benchmark VIII - "Kovasznay"	87
7.5	Geodynamical benchmarks	88
7.6	Assigning values to quadrature points	89
7.7	Matrix (Sparse) storage	92
7.7.1	2D domain - One degree of freedom per node	92
7.7.2	2D domain - Two degrees of freedom per node	93
7.7.3	in fieldstone	94
7.8	Mesh generation	95
7.8.1	Quadrilateral-based meshes	95
7.8.2	Delaunay triangulation and Voronoi cells	96
7.8.3	Tetrahedra	97
7.8.4	Hexahedra	97
7.8.5	Adaptive Mesh Refinement	97
7.9	Visco-Plasticity	101
7.9.1	Tensor invariants	101
7.9.2	Scalar viscoplasticity	102
7.9.3	About the yield stress value Y	102
7.10	Pressure smoothing	103
7.11	Pressure scaling	105
7.12	Pressure normalisation	106
7.12.1	Basic idea and naive implementation	106
7.12.2	Implementation – the real deal	106
7.13	Solving the Stokes system	107
7.13.1	when using the penalty formulation	107
7.13.2	Conjugate gradient and the Schur complement approach	107
7.13.3	Conjugate gradient and the Schur complement approach	109
7.13.4	The Augmented Lagrangian approach	113
7.13.5	The GMRES approach	114
7.14	The consistent boundary flux (CBF)	115
7.14.1	applied to the Stokes equation	115
7.14.2	applied to the heat equation	115
7.14.3	implementation - Stokes equation	116
7.15	The value of the timestep	118
7.16	Mappings	119
7.16.1	Linear mapping on a triangle	119
7.16.2	Linear mapping on a quadrilateral	119
7.17	Exporting data to vtk format	121
7.18	Runge-Kutta methods	123
7.18.1	Using RK methods to advect particles/markers	123
7.19	Am I in or not? - finding reduced coordinates	125
7.19.1	Two-dimensional space	125
7.19.2	Three-dimensional space	125
7.20	Error measurements and convergence rates	128
7.20.1	About extrapolation	129
7.21	The initial temperature field	130
7.21.1	Single layer with imposed temperature b.c.	130
7.21.2	Single layer with imposed heat flux b.c.	131
7.21.3	Single layer with imposed heat flux and temperature b.c.	131
7.21.4	Half cooling space	131
7.21.5	Plate model	131
7.21.6	McKenzie slab	131
7.22	Kinematic boundary conditions	134
7.22.1	In-out flux boundary conditions for lithospheric models	134
7.23	Computing gradients - the recovery process	135
7.23.1	Global recovery	135

7.23.2 Local recovery - centroid average over patch	135
7.23.3 Local recovery - nodal average over patch	135
7.23.4 Local recovery - least squares over patch	135
7.23.5 Link to pressure smoothing	135
7.24 Tracking materials and/or interfaces	136
7.24.1 The Particle-in-cell technique	136
7.24.2 The level set function technique	139
7.24.3 The field/composition technique	139
7.24.4 Hybrid methods	140
7.25 Static condensation	141
7.26 Measuring incompressibility	142
7.27 Periodic boundary conditions	143
7.28 Removing rotational nullspace	144
7.28.1 Three dimensions	145
7.28.2 Two dimensions	145
7.29 Picard and Newton	147
7.29.1 Picard iterations	147
7.30 Defect correction formulation	148
7.31 Parallel or not?	149
7.32 Stream function	150
7.32.1 In Cartesian coordinates	150
7.32.2 In Cylindrical coordinates	150
7.33 Corner flow	151
8 fieldstone_01: simple analytical solution	153
9 fieldstone_02: Stokes sphere	155
10 fieldstone_03: Convection in a 2D box	156
11 fieldstone_04: The lid driven cavity	159
11.1 the lid driven cavity problem (<code>ldc=0</code>)	159
11.2 the lid driven cavity problem - regularisation I (<code>ldc=1</code>)	159
11.3 the lid driven cavity problem - regularisation II (<code>ldc=2</code>)	159
12 fieldstone_05: SolCx benchmark	161
13 fieldstone_06: SolKz benchmark	163
14 fieldstone_07: SolVi benchmark	164
15 fieldstone_08: the indentor benchmark	166
16 fieldstone_09: the annulus benchmark	168
17 fieldstone_10: Stokes sphere (3D) - penalty	170
18 fieldstone_11: stokes sphere (3D) - mixed formulation	171
19 fieldstone_12: consistent pressure recovery	172
20 fieldstone_13: the Particle in Cell technique (1) - the effect of averaging	174
21 fieldstone_f14: solving the full saddle point problem	178
22 fieldstone_f15: saddle point problem with Schur complement approach - benchmark	181

23 fieldstone_f16: saddle point problem with Schur complement approach - Stokes sphere	184
24 fieldstone_17: solving the full saddle point problem in 3D	186
24.0.1 Constant viscosity	188
24.0.2 Variable viscosity	189
25 fieldstone_18: solving the full saddle point problem with $Q_2 \times Q_1$ elements	191
26 fieldstone_19: solving the full saddle point problem with $Q_3 \times Q_2$ elements	193
27 fieldstone_20: the Busse benchmark	195
28 fieldstone_21: The non-conforming $Q_1 \times P_0$ element	197
29 fieldstone_22: The stabilised $Q_1 \times Q_1$ element	198
29.1 The Donea & Huerta benchmark	200
29.2 The Dohrmann & Bochev benchmark	201
29.3 The falling block experiment	201
30 fieldstone_23: compressible flow (1) - analytical benchmark	202
31 fieldstone_24: compressible flow (2) - convection box	205
31.1 The physics	205
31.2 The numerics	205
31.3 The experimental setup	207
31.4 Scaling	207
31.5 Conservation of energy 1	208
31.5.1 under BA and EBA approximations	208
31.5.2 under no approximation at all	209
31.6 Conservation of energy 2	209
31.7 The problem of the onset of convection	210
31.8 results - BA - $Ra = 10^4$	212
31.9 results - BA - $Ra = 10^5$	214
31.10 results - BA - $Ra = 10^6$	215
31.11 results - EBA - $Ra = 10^4$	216
31.12 results - EBA - $Ra = 10^5$	218
31.13 Onset of convection	219
32 fieldstone_25: Rayleigh-Taylor instability (1) - instantaneous	221
33 fieldstone_26: Slab detachment benchmark (1) - instantaneous	223
34 fieldstone_27: Consistent Boundary Flux	225
35 fieldstone_28: convection 2D box - Tosi et al, 2015	229
35.0.1 Case 0: Newtonian case, a la Blankenbach et al., 1989	230
35.0.2 Case 1	231
35.0.3 Case 2	233
35.0.4 Case 3	235
35.0.5 Case 4	237
35.0.6 Case 5	239
36 fieldstone_29: open boundary conditions	241

37 fieldstone_30: conservative velocity interpolation	244
37.1 Couette flow	244
37.2 SolCx	244
37.3 Streamline flow	244
38 fieldstone_31: conservative velocity interpolation 3D	245
39 fieldstone_32: 2D analytical sol. from stream function	246
39.1 Background theory	246
39.2 A simple application	246
40 fieldstone_33: Convection in an annulus	250
41 fieldstone_34: the Cartesian geometry elastic aquarium	252
42 fieldstone_35: 2D analytical sol. in annulus from stream function	254
42.1 Linking with our paper	257
42.2 No slip boundary conditions	257
42.3 Free slip boundary conditions	258
43 fieldstone_36: the annulus geometry elastic aquarium	260
44 fieldstone_37: marker advection and population control	263
45 fieldstone_38: Critical Rayleigh number	264
46 fieldstone_39: chpe15	266
47 fieldstone_40: Rayleigh-Taylor instability	276
48 fieldstone_42: 1D diffusion	278
49 fieldstone_43: the rotating cone	279
50 fieldstone_44: the flat slab	282
51 fieldstone_45: the corner flow	283
52 fieldstone_46: MMS1 with Crouzeix-Raviart elements	286
53 fieldstone: Gravity: buried sphere	287
54 Problems, to do list and projects for students	289
A Three-dimensional applications	291
B Codes in geodynamics	292
C Matrix properties	295
C.1 Symmetric matrices	295
C.2 Schur complement	295
D Dont be a hero - unless you have to	297
E A FANTOM, an ELEFANT and a GHOST	299
F Some useful Python commands	302
F.1 Sparse matrices	302
F.2 condition number	302

G Some useful maths	303
G.1 Inverse of a 3x3 matrix	303
G.2 Inverse of a 3x3 matrix	303

WARNING: this is work in progress

1 Introduction

1.1 Philosophy

This document was writing 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 or statement I make. These are to be found in 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 value code readability over code efficiency. I have also chosen to avoid resorting to multiple code files or even functions 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 [358, 299], CITCOM, LAMEM, PTATIN, PYLITH, ...

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

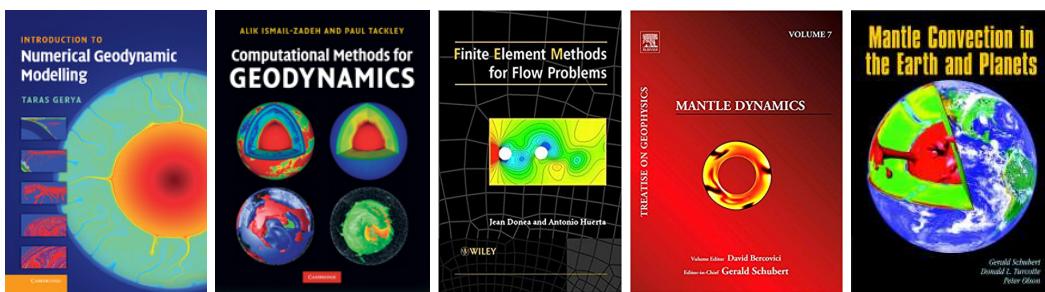
All the python scripts and this document are freely available at

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

1.2 Acknowledgments

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, Rens Elbertsen, Philippe Fullsack, Menno Fraters, Anne Glerum, Timo Heister, Robert Myhill, John Naliboff, E. Gerry Puckett, Melchior Schuh-Senlis, Michael Tetley, Lukas van de Wiel, Arie van den Berg, Tom Weir, and the whole ASPECT family/team.

1.3 Essential literature



<http://www-udc.ig.utexas.edu/external/becker/Geodynamics557.pdf>

1.4 Installation

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

```
python3.6 -m pip install --user numpy scipy matplotlib
```

To install the umfpack solver:

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

1.5 What is a fieldstone?

Simply put, it is 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.

1.6 Why the Finite Element method?

The Finite Element Method (FEM) is by no means the only method to solve PDEs in geodynamics, nor is it necessarily 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) [174, 175, 219], or the Element Free Galerkin Method (EFGM) [292]. I have been using FEM since 2008 and I do not have real experience to speak of in FVM or FDM so I concentrate in this book on what I know best.

1.7 Notations

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 M , σ , 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 $M \cdot \vec{a} = M_{ij} a_j$. 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 product¹. Case in point, $\vec{\nabla} \cdot \vec{v}$ is the velocity divergence while $\vec{\nabla}\vec{v}$ is the velocity gradient tensor.

1.8 Colour maps for visualisation

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 this link and are Perceptually Uniform Colour Maps [357].

Field	colour code
Velocity/displacement	CET-D1A
Pressure	CET-L17
Velocity divergence	CET-L1
Density	CET-D3
Strain rate	CET-R2
Viscosity	CET-R3
Temperature	CET-D9

¹<https://en.wikipedia.org/wiki/Dyadics>

2 List of tutorials

tutorial number	element	outer solver	formulation	physical problem	3D	temperature	time stepping	nonlinear	compressible	analytical benchmark	numerical benchmark	elastomechanics
1	$Q_1 \times P_0$		penalty	analytical benchmark						†		
2	$Q_1 \times P_0$		penalty	Stokes sphere								
3	$Q_1 \times P_0$		penalty	Blankenbach et al., 1989	†	†						
4	$Q_1 \times P_0$		penalty	Lid driven cavity								
5	$Q_1 \times P_0$		penalty	SolCx benchmark								
6	$Q_1 \times P_0$		penalty	SolKz benchmark								
7	$Q_1 \times P_0$		penalty	SolVi benchmark								
8	$Q_1 \times P_0$		penalty	Indentor				†				
9	$Q_1 \times P_0$		penalty	annulus benchmark								
10	$Q_1 \times P_0$		penalty	Stokes sphere	†							
11	$Q_1 \times P_0$	full matrix	mixed	Stokes sphere	†							
12	$Q_1 \times P_0$		penalty	analytical benchmark + consistent press recovery								
13	$Q_1 \times P_0$		penalty	Stokes sphere + markers averaging								
14	$Q_1 \times P_0$	full matrix	mixed	analytical benchmark								
15	$Q_1 \times P_0$	Schur comp. CG	mixed	analytical benchmark								
16	$Q_1 \times P_0$	Schur comp. PCG	mixed	Stokes sphere								
17	$Q_2 \times Q_1$	full matrix	mixed	Burstedde benchmark	†							
18	$Q_2 \times Q_1$	full matrix	mixed	analytical benchmark								
19	$Q_3 \times Q_2$	full matrix	mixed	analytical benchmark								
20	$Q_1 \times P_0$		penalty	Busse et al., 1993	†	†	†					
21	$Q_1 \times P_0$ R-T		penalty	analytical benchmark								
22	$Q_1 \times Q_1$ - stab	full matrix	mixed	analytical benchmark								
23	$Q_1 \times P_0$		mixed	analytical benchmark					†			
24	$Q_1 \times P_0$		mixed	convection box		†	†		†			

tutorial number	element	outer solver	formulation	physical problem	3D	temperature	time stepping	nonlinear	compressible	analytical benchmark	numerical benchmark	elastomechanics
25	$Q_1 \times P_0$	full matrix	mixed	Rayleigh-Taylor instability								
26	$Q_1 \times P_0$	full matrix	mixed	Slab detachment				†				
27	$Q_1 \times P_0$	full matrix	mixed	CBF benchmarks				†		†		
28	$Q_1 \times P_0$	full matrix	mixed	Tosi et al, 2015		†	†	†			†	
29	$Q_1 \times P_0$	full matrix	mixed	Open Boundary conditions						†		
30	Q_1, Q_2	X	X	Cons. Vel. Interp (cvf)		†				†		
31	Q_1, Q_2	X	X	Cons. Vel. Interp (cvf)	†	†				†		
32	$Q_1 \times P_0$	full matrix	mixed	analytical benchmark						†		
33	$Q_1 \times P_0$		penalty	convection in annulus	†	†	†					
34	Q_1			elastic Cartesian aquarium					†	†	†	
35												
36	Q_1			elastic annulus aquarium					†	†	†	
37	Q_1, Q_2	X	X	population control, bmw test		†				†		
38				Critical Rayleigh number								
39	$Q_2 \times Q_1$	full matrix	mixed									
XX												

Analytical benchmark means that an analytical solution exists while numerical benchmark means that a comparison with other code(s) has been carried out.

3 The physical equations

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}	velocity vector	$m \cdot s^{-1}$
ρ	mass density	kg/m^3
η	dynamic viscosity	$Pa \cdot s$
λ	penalty parameter	$Pa \cdot s$
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}
α	thermal expansion coefficient	K^{-1}
k	thermal conductivity	$W/(m \cdot K)$
C_p	Heat capacity	J/K
H	intrinsic specific heat production	W/kg
β_T	isothermal compressibility	Pa^{-1}
τ	deviatoric stress tensor	Pa
σ	full stress tensor	Pa

3.1 Strain rate and spin tensor

The velocity gradient is given in Cartesian coordinates by:

$$\vec{\nabla} \vec{v} = \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial v}{\partial x} & \frac{\partial w}{\partial x} \\ \frac{\partial u}{\partial y} & \frac{\partial v}{\partial y} & \frac{\partial w}{\partial y} \\ \frac{\partial u}{\partial z} & \frac{\partial v}{\partial z} & \frac{\partial w}{\partial z} \end{pmatrix}$$

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

$$\vec{\nabla} \vec{v} = \vec{\nabla}^s \vec{v} + \vec{\nabla}^w \vec{v}$$

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

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

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

$$\vec{R}(\vec{v}) = \frac{1}{2} (\vec{\nabla} \vec{v} - \vec{\nabla} \vec{v}^T)$$

3.2 The heat transport equation - energy conservation equation

Let us start from the heat transport equation as shown in Schubert, Turcotte and Olson [500]:

$$\rho C_p \frac{DT}{Dt} - \alpha T \frac{Dp}{Dt} = \vec{\nabla} \cdot k \vec{\nabla} T + \Phi + \rho H$$

with D/Dt being the total derivatives so that

$$\frac{DT}{Dt} = \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \quad \frac{Dp}{Dt} = \frac{\partial p}{\partial t} + \vec{v} \cdot \vec{\nabla} p$$

Solving for temperature, this equation is often rewritten as follows:

$$\rho C_p \frac{DT}{Dt} - \vec{\nabla} \cdot k \vec{\nabla} T = \alpha T \frac{Dp}{Dt} + \Phi + \rho H$$

A note on the shear heating term Φ : In many publications, Φ is given by $\Phi = \tau_{ij} \partial_j u_i = \boldsymbol{\tau} : \vec{\nabla} \vec{\nabla}$.

$$\begin{aligned}
\Phi &= \tau_{ij} \partial_j u_i \\
&= 2\eta \dot{\varepsilon}_{ij}^d \partial_j u_i \\
&= 2\eta \frac{1}{2} (\dot{\varepsilon}_{ij}^d \partial_j u_i + \dot{\varepsilon}_{ji}^d \partial_i u_j) \\
&= 2\eta \frac{1}{2} (\dot{\varepsilon}_{ij}^d \partial_j u_i + \dot{\varepsilon}_{ij}^d \partial_i u_j) \\
&= 2\eta \dot{\varepsilon}_{ij}^d \frac{1}{2} (\partial_j u_i + \partial_i u_j) \\
&= 2\eta \dot{\varepsilon}_{ij}^d \dot{\varepsilon}_{ij} \\
&= 2\eta \dot{\varepsilon}^d : \dot{\varepsilon} \\
&= 2\eta \dot{\varepsilon}^d : \left(\dot{\varepsilon}^d + \frac{1}{3} (\vec{\nabla} \cdot \vec{v}) \mathbf{1} \right) \\
&= 2\eta \dot{\varepsilon}^d : \dot{\varepsilon}^d + 2\eta \dot{\varepsilon}^d : \mathbf{1} (\vec{\nabla} \cdot \vec{v}) \\
&= 2\eta \dot{\varepsilon}^d : \dot{\varepsilon}^d
\end{aligned} \tag{1}$$

Finally

$$\Phi = \boldsymbol{\tau} : \vec{\nabla} \vec{\nabla} = 2\eta \dot{\varepsilon}^d : \dot{\varepsilon}^d = 2\eta ((\dot{\varepsilon}_{xx}^d)^2 + (\dot{\varepsilon}_{yy}^d)^2 + 2(\dot{\varepsilon}_{xy}^d)^2)$$

3.3 The momentum conservation equations

Because the Prandlt number is virtually zero in Earth science applications the Navier Stokes equations reduce to the Stokes equation:

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

Since

$$\boldsymbol{\sigma} = -p \mathbf{1} + \boldsymbol{\tau}$$

it also writes

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

Using the relationship $\boldsymbol{\tau} = 2\eta \dot{\varepsilon}^d$ we arrive at

$$-\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\varepsilon}^d) + \rho \vec{g} = 0$$

3.4 The mass conservation equations

The mass conservation equation is given by

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

or,

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

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:

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

3.5 The equations in ASPECT manual

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 [500].

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

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

$$\vec{\nabla} \cdot (\rho \vec{\mathbf{v}}) = 0 \quad \text{in } \Omega, \quad (3)$$

$$\begin{aligned} \rho C_p \left(\frac{\partial T}{\partial t} + \vec{\mathbf{v}} \cdot \vec{\nabla} T \right) - \vec{\nabla} \cdot k \vec{\nabla} T &= \rho H \\ &\quad + 2\eta \left(\dot{\boldsymbol{\varepsilon}}(\mathbf{v}) - \frac{1}{3}(\vec{\nabla} \cdot \vec{\mathbf{v}})\mathbf{1} \right) : \left(\dot{\boldsymbol{\varepsilon}}(\mathbf{v}) - \frac{1}{3}(\vec{\nabla} \cdot \vec{\mathbf{v}})\mathbf{1} \right) \\ &\quad + \alpha T (\mathbf{v} \cdot \vec{\nabla} p) \end{aligned} \quad \text{in } \Omega, \quad (4)$$

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

In this set of equations, (449) and (450) represent the compressible Stokes equations in which $\mathbf{v} = \mathbf{v}(\mathbf{x}, t)$ 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 (451) 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 that ASPECT solves, we need to

- neglect the $\partial p / \partial t$. **WHY?**
- neglect the $\partial \rho / \partial t$. **WHY?**

from equations above.

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

$$\Phi = k_B(\nabla \cdot \mathbf{v})^2 + 2\eta \dot{\boldsymbol{\varepsilon}}^d : \dot{\boldsymbol{\varepsilon}}^d$$

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,

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

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

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

3.6 The Boussinesq approximation: an Incompressible flow

[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 (449). The primary result of this assumption is that the continuity equation (450) will now read

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

This implies that the strain rate tensor is deviatoric. Under the Boussinesq approximation, the equations are much simplified:

$$-\nabla \cdot [2\eta\dot{\epsilon}(\mathbf{v})] + \nabla p = \rho\mathbf{g} \quad \text{in } \Omega, \quad (5)$$

$$\nabla \cdot (\rho\mathbf{v}) = 0 \quad \text{in } \Omega, \quad (6)$$

$$\rho_0 C_p \left(\frac{\partial T}{\partial t} + \mathbf{v} \cdot \nabla T \right) - \nabla \cdot k \nabla T = \rho H \quad \text{in } \Omega \quad (7)$$

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

3.7 Stokes equation for elastic medium

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

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

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

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

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

$$\sigma_{ij} = \sum_{kl} C_{ijkl} \epsilon_{kl} \quad (8)$$

where \mathbf{C} is the fourth-order elastic tensor. In the case of an isotropic material, this relationship simplifies to

$$\sigma_{ij} = \lambda \epsilon_{kk} \delta_{ij} + 2\mu \epsilon_{ij} \quad \text{or,} \quad \boldsymbol{\sigma} = \lambda(\nabla \cdot \mathbf{u}) \mathbf{1} + 2\mu \boldsymbol{\epsilon} \quad (9)$$

where λ is the Lamé parameter and μ is the shear modulus². The term $\nabla \cdot \mathbf{u}$ is the isotropic dilation.

The strain tensor is related to the displacement as follows:

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

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

$$\begin{aligned} p &= -\frac{1}{3} \text{Tr}(\boldsymbol{\sigma}) \\ &= -\frac{1}{3} [\lambda(\nabla \cdot \mathbf{u}) \text{Tr}[\mathbf{1}] + 2\mu \text{Tr}[\boldsymbol{\epsilon}]] \\ &= -\frac{1}{3} [\lambda(\nabla \cdot \mathbf{u}) 3 + 2\mu(\nabla \cdot \mathbf{u})] \\ &= -[\lambda + \frac{2}{3}\mu](\nabla \cdot \mathbf{u}) \end{aligned} \quad (10)$$

so that $K = \lambda + \frac{2}{3}\mu$.

Remark. Eq. (8) and (9) 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, \mathbf{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.

²It is also sometimes written G

3.8 The strain rate tensor in all coordinate systems

The strain rate tensor $\dot{\epsilon}$ is given by

$$\dot{\epsilon} = \frac{1}{2}(\vec{\nabla}\vec{v} + \vec{\nabla}\vec{v}^T) \quad (11)$$

3.8.1 Cartesian coordinates

$$\dot{\epsilon}_{xx} = \frac{\partial u}{\partial x} \quad (12)$$

$$\dot{\epsilon}_{yy} = \frac{\partial v}{\partial y} \quad (13)$$

$$\dot{\epsilon}_{zz} = \frac{\partial w}{\partial z} \quad (14)$$

$$\dot{\epsilon}_{yx} = \dot{\epsilon}_{xy} = \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) \quad (15)$$

$$\dot{\epsilon}_{zx} = \dot{\epsilon}_{xz} = \frac{1}{2} \left(\frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \right) \quad (16)$$

$$\dot{\epsilon}_{zy} = \dot{\epsilon}_{yz} = \frac{1}{2} \left(\frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \right) \quad (17)$$

3.8.2 Polar coordinates

$$\dot{\epsilon}_{rr} = \frac{\partial v_r}{\partial r} \quad (18)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{v_r}{r} + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} \quad (19)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(\frac{\partial v_\theta}{\partial r} - \frac{v_\theta}{r} + \frac{1}{r} \frac{\partial v_r}{\partial \theta} \right) \quad (20)$$

3.8.3 Cylindrical coordinates

<http://eml.ou.edu/equation/FLUIDS/STRAIN/STRAIN.HTM>

3.8.4 Spherical coordinates

$$\dot{\epsilon}_{rr} = \frac{\partial v_r}{\partial r} \quad (21)$$

$$\dot{\epsilon}_{\theta\theta} = \frac{v_r}{r} + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} \quad (22)$$

$$\dot{\epsilon}_{\phi\phi} = \frac{1}{r \sin \theta} \frac{\partial v_\phi}{\partial \phi} \quad (23)$$

$$\dot{\epsilon}_{\theta r} = \dot{\epsilon}_{r\theta} = \frac{1}{2} \left(r \frac{\partial}{\partial r} \left(\frac{v_\theta}{r} \right) + \frac{1}{r} \frac{\partial v_r}{\partial \theta} \right) \quad (24)$$

$$\dot{\epsilon}_{\phi r} = \dot{\epsilon}_{r\phi} = \frac{1}{2} \left(\frac{1}{r \sin \theta} \frac{\partial v_r}{\partial \phi} + r \frac{\partial}{\partial r} \left(\frac{v_\phi}{r} \right) \right) \quad (25)$$

$$\dot{\epsilon}_{\phi\theta} = \dot{\epsilon}_{\theta\phi} = \frac{1}{2} \left(\frac{\sin \theta}{r} \frac{\partial}{\partial \theta} \left(\frac{v_\phi}{\sin \theta} \right) + \frac{1}{r \sin \theta} \frac{\partial v_\theta}{\partial \phi} \right) \quad (26)$$

3.9 Boundary conditions

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 on 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.

3.9.1 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. $\sigma \cdot \vec{n} = \vec{0}$. 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 $(\sigma \cdot \vec{n}) \times \vec{n} = \vec{0}$. 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 fieldstone 29.

3.9.2 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).

3.10 Meaningful physical quantities

- 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):

$$v_{rms} = \left(\frac{\int_{\Omega} |\vec{v}|^2 d\Omega}{\int_{\Omega} d\Omega} \right)^{1/2} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} |\vec{v}|^2 d\Omega \right)^{1/2} \quad (27)$$

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

In Cartesian coordinates, for a cuboid domain of size $L_x \times L_y \times L_z$, the v_{rms} is simply given by:

$$v_{rms} = \left(\frac{1}{L_x L_y L_z} \int_0^{L_x} \int_0^{L_y} \int_0^{L_z} (u^2 + v^2 + w^2) dx dy dz \right)^{1/2} \quad (28)$$

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:

$$v_r|_{rms} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} v_r^2 d\Omega \right)^{1/2} \quad (29)$$

$$v_{\theta}|_{rms} = \left(\frac{1}{V_{\Omega}} \int_{\Omega} v_{\theta}^2 d\Omega \right)^{1/2} \quad (30)$$

- Pressure p (Pa):
- Stress tensor σ (Pa):
- Strain tensor ϵ (dimensionless):
- Strain rate tensor $\dot{\epsilon}$ (s⁻¹):
- Rayleigh number Ra (X):
- Prandtl number Pr (X): It is named after the German physicist Ludwig Prandtl and is defined as the ratio of momentum diffusivity to thermal diffusivity. It is given as:

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

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

- Nusselt number N_u (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:

$$Nu = \frac{q}{q_c} \quad (31)$$

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 [64]

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

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

finish, not happy with definition. Look at literature

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 [602] 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 [328], the Nusselt number at the inner and outer boundaries are:

$$\text{Nu}_{inner} = \frac{f \ln f}{1-f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T}{\partial r} \right)_{r=R_1} d\theta \quad (32)$$

$$\text{Nu}_{outer} = \frac{\ln f}{1-f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T}{\partial r} \right)_{r=R_2} d\theta \quad (33)$$

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

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

We then find:

$$\text{Nu}_{inner} = \frac{f \ln f}{1-f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T_c}{\partial r} \right)_{r=R_1} d\theta = \frac{f \ln f}{1-f} \frac{1}{R_1} \frac{1}{\ln f} = 1 \quad (34)$$

$$\text{Nu}_{outer} = \frac{\ln f}{1-f} \frac{1}{2\pi} \int_0^{2\pi} \left(\frac{\partial T_c}{\partial r} \right)_{r=R_2} d\theta = \frac{\ln f}{1-f} \frac{1}{R_2} \frac{1}{\ln f} = 1 \quad (35)$$

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

- Temperature (K):
- Viscosity (Pa.s):
- Density (kg/m^3):
- Heat capacity C_p ($\text{J.K}^{-1}.\text{kg}^{-1}$): It is the measure of the heat energy required to increase the temperature of a unit quantity of a substance by unit degree.
- 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.
- 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.

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

3.11 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.

4 The building blocks of the Finite Element Method

4.1 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

$$\int_a^b f(x)dx$$

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.

4.1.1 in 1D - theory

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.

$$\int_a^b f(x)dx \simeq (b-a)f\left(\frac{a+b}{2}\right)$$

insert here figure

The interpolating function may 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.

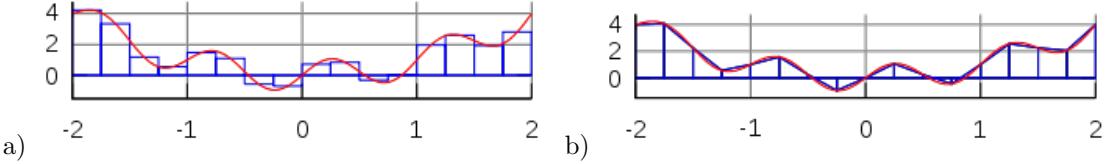
$$\int_a^b f(x)dx \simeq (b-a)\frac{f(a) + f(b)}{2}$$

insert here figure

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. This is called a composite rule, extended rule, or iterated rule. For example, the composite trapezoidal rule can be stated as

$$\int_a^b f(x)dx \simeq \frac{b-a}{n} \left(\frac{f(a)}{2} + \sum_{k=1}^{n-1} f(a + k \frac{b-a}{n}) + \frac{f(b)}{2} \right)$$

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 NewtonCotes formulas, of which the rectangle rule and the trapezoidal rule are examples. 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 NewtonCotes 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 $P_n(x)$.

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.

Number of points, n	Points, x_i	Weights, w_i
1	0	2
2	$\pm\sqrt{\frac{1}{3}}$	1
3	0	$\frac{8}{9}$
	$\pm\sqrt{\frac{3}{5}}$	$\frac{5}{9}$
4	$\pm\sqrt{\frac{3}{7} - \frac{2}{7}\sqrt{\frac{6}{5}}}$	$\frac{18+\sqrt{30}}{36}$
	$\pm\sqrt{\frac{3}{7} + \frac{2}{7}\sqrt{\frac{6}{5}}}$	$\frac{18-\sqrt{30}}{36}$
5	0	$\frac{128}{225}$
	$\pm\frac{1}{3}\sqrt{5 - 2\sqrt{\frac{10}{7}}}$	$\frac{322+13\sqrt{70}}{900}$
	$\pm\frac{1}{3}\sqrt{5 + 2\sqrt{\frac{10}{7}}}$	$\frac{322-13\sqrt{70}}{900}$

Gauss-Legendre points and their weights.

As shown in the above table, it can be shown that the weight values must fulfill the following condition:

$$\sum_{i_q} w_{i_q} = 2 \tag{36}$$

and it is 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

$$r = \frac{2}{b-a}(x-a) - 1$$

This relationship can be reversed such that when r is known, its equivalent coordinate $x \in [a, b]$ can be computed:

$$x = \frac{b-a}{2}(1+r) + a$$

From this it follows that

$$dx = \frac{b-a}{2}dr$$

and then

$$\int_a^b f(x)dx = \frac{b-a}{2} \int_{-1}^{+1} f(r)dr \simeq \frac{b-a}{2} \sum_{i_q=1}^n w_{i_q} f(r_{i_q})$$

4.1.2 in 1D - 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 $f(x) = \pi$. Then we can compute the integral of this function over the interval $[a, b]$ exactly:

$$I = \int_{-1}^{+1} f(x)dx = \pi \int_{-1}^{+1} dx = 2\pi$$

We can now use a Gauss-Legendre formula to compute this same integral:

$$I_{gq} = \int_{-1}^{+1} f(x)dx = \sum_{i_q=1}^{n_q} w_{i_q} f(x_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} \pi = \pi \underbrace{\sum_{i_q=1}^{n_q} w_{i_q}}_{=2} = 2\pi$$

where we have used the property of the weight values of Eq.(36). Since the actual number of points was never specified, this result is valid for all quadrature rules.

example 2 Let us now take $f(x) = mx + p$ and repeat the same exercise:

$$I = \int_{-1}^{+1} f(x)dx = \int_{-1}^{+1} (mx + p)dx = [\frac{1}{2}mx^2 + px]_{-1}^{+1} = 2p$$

$$I_{gq} = \int_{-1}^{+1} f(x)dx = \sum_{i_q=1}^{n_q} w_{i_q} f(x_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} (mx_{i_q} + p) = m \underbrace{\sum_{i_q=1}^{n_q} w_{i_q} x_{i_q}}_{=0} + p \underbrace{\sum_{i_q=1}^{n_q} w_{i_q}}_{=2} = 2p$$

since the quadrature points are symmetric w.r.t. to zero on the x-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 $f(x) = x^2$. We have

$$I = \int_{-1}^{+1} f(x)dx = \int_{-1}^{+1} x^2 dx = [\frac{1}{3}x^3]_{-1}^{+1} = \frac{2}{3}$$

and

$$I_{gq} = \int_{-1}^{+1} f(x)dx = \sum_{i_q=1}^{n_q} w_{i_q} f(x_{i_q}) = \sum_{i_q=1}^{n_q} w_{i_q} x_{i_q}^2$$

- $n_q = 1$: $x_{i_q}^{(1)} = 0$, $w_{i_q} = 2$. $I_{gq} = 0$
- $n_q = 2$: $x_q^{(1)} = -1/\sqrt{3}$, $x_q^{(2)} = 1/\sqrt{3}$, $w_q^{(1)} = w_q^{(2)} = 1$. $I_{gq} = \frac{2}{3}$
- It also works $\forall n_q > 2$!

4.1.3 in 2D/3D - theory

Let us now turn to a two-dimensional integral of the form

$$I = \int_{-1}^{+1} \int_{-1}^{+1} f(x, y) dx dy$$

The equivalent Gaussian quadrature writes:

$$I_{gq} \simeq \sum_{i_q=1}^{n_q} \sum_{j_q}^{n_q} f(x_{i_q}, y_{j_q}) w_{i_q} w_{j_q}$$

4.1.4 quadrature on triangles

Dunavant, 1985 [166].

ip_triangle.m from MILAMIN
BETTER REF???

r_q	s_q	w_q		
1/3	1/3	1/2		
1/6	1/6	1/6		
2/3	1/6	1/6		
1/6	2/3	1/6		
1/3	1/3	-27/96		
0.6	0.2	25/96		
0.2	0.6	25/96		
0.2	0.2	25/96		
1 - 2g ₁	g_1	$w_1/2$	0.108103018168070	0.44594849091596
	g_1	$w_1/2$	0.445948490915965	0.108103018168070
	g_1	$w_1/2$	0.445948490915965	0.445948490915965
1 - 2g ₂	g_2	$w_2/2$	0.816847572980459	0.091576213509771
	g_2	$w_2/2$	0.091576213509771	0.816847572980459
	g_2	$w_2/2$	0.091576213509771	0.091576213509771
			0.091576213509771	0.091576213509771
			0.816847572980459	0.091576213509771
			0.091576213509771	0.816847572980459
			0.445948490915965	0.445948490915965
			0.108103018168070	0.445948490915965
			0.445948490915965	0.108103018168070
			0.091576213509771	0.223381589678011/2.0
			0.816847572980459	0.091576213509771
			0.091576213509771	0.109951743655322/2.0
			0.445948490915965	0.445948490915965
			0.223381589678011/2.0	0.109951743655322/2.0
			0.108103018168070	0.223381589678011/2.0
			0.445948490915965	0.223381589678011/2.0
			0.091576213509771	0.223381589678011/2.0
			0.1012865073235	0.1012865073235
			0.7974269853531	0.1012865073235
			0.1012865073235	0.7974269853531
			0.4701420641051	0.0597158717898
			0.4701420641051	0.4701420641051
			0.0597158717898	0.0597158717898
			0.3333333333333	0.3333333333333
			0.3333333333333	0.1125000000000
			5.01426509658179E - 01	2.49286745170910E - 01
			2.49286745170910E - 01	5.01426509658179E - 01
			2.49286745170910E - 01	5.83931378631895E - 02
			8.73821971016996E - 01	6.30890144915020E - 02
			6.30890144915020E - 02	2.54224531851035E - 02
			6.30890144915020E - 02	6.30890144915020E - 02
			5.31450498448170E - 02	3.10352451033784E - 01
			6.36502499121399E - 01	5.31450498448170E - 02
			3.10352451033784E - 01	4.14255378091870E - 02
			5.31450498448170E - 02	6.36502499121399E - 01
			6.36502499121399E - 01	4.14255378091870E - 02
			3.10352451033784E - 01	4.14255378091870E - 02
			5.31450498448170E - 02	6.36502499121399E - 01
			6.36502499121399E - 01	4.14255378091870E - 02
			3.10352451033784E - 01	5.31450498448170E - 02

where $g_1 = (8 - \sqrt{10}) + \sqrt{38 - 44 * \sqrt{2/5}})/18$ $g_2 = (8 - \sqrt{10}) - \sqrt{38 - 44 * \sqrt{2/5}})/18$
 $w_1 = (620 + \sqrt{213125 - 53320 * \sqrt{10}}))/3720$ $w_2 = (620 - \sqrt{213125 - 53320 * \sqrt{10}}))/3720$

4.1.5 quadrature on tetrahedra

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:

$$\int \int \int_{el} f(x, y, z) dx dy dz = V_{el} \sum_{iq=1}^{nqel} w_{iq} f(\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq})$$

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 $w_{iq} = V_{el} * 0.25 = 1/24 \simeq 04166666666666667$ 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:

$$\begin{aligned} r_{iq} &= \sum_{i=1}^4 \xi_i^{iq} x_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (r_1, r_2, r_3, r_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 1, 0, 0) = \xi_2^{iq} \\ s_{iq} &= \sum_{i=1}^4 \xi_i^{iq} y_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (s_1, s_2, s_3, s_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 0, 1, 0) = \xi_3^{iq} \\ t_{iq} &= \sum_{i=1}^4 \xi_i^{iq} z_i = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (t_1, t_2, t_3, t_4) = (\xi_1^{iq}, \xi_2^{iq}, \xi_3^{iq}, \xi_4^{iq}) \cdot (0, 0, 0, 1) = \xi_4^{iq} \end{aligned}$$

Finally:

	r_q	s_q	t_q	w_q
iq = 1	0.138196601125011	0.138196601125011	0.138196601125011	0.04166666666666667
iq = 2	0.585410196624969	0.138196601125011	0.138196601125011	0.04166666666666667
iq = 3	0.138196601125011	0.585410196624969	0.138196601125011	0.04166666666666667
iq = 4	0.138196601125011	0.138196601125011	0.585410196624969	0.04166666666666667

4.2 The mesh

4.3 A bit of FE terminology

We introduce here some terminology for efficient element descriptions [277]:

- 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

$$u \sim a_1 + a_2x + a_3y + a_4xy + a_5x^2 + a_6y^2$$

with similar approximations for v , and

$$p \sim b_1 + b_2x + b_3y$$

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 $P_2 P_1$ 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

$$u \sim a_1 + a_2x + a_3y + a_4xy + a_5x^2 + a_6y^2 + a_7x^2y + a_8xy^2 + a_9x^2y^2$$

and

$$p \sim b_1 + b_2x + b_3y + b_4xy$$

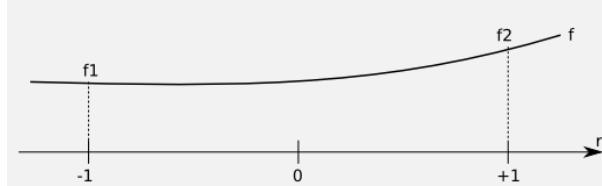
- For these same families, $Q_m \times Q_{-n}$ is as above, except that the pressure approximation is not continuous at element boundaries.
- Again for the same families, $Q_m \times P_{-n}$ indicates the same velocity approximation with a pressure approximation that is a discontinuous complete piecewise 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 [277], conforming velocity elements are those for which the basis functions for 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 ??) 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 [146].

4.4 Elements and basis functions in 1D

4.4.1 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

$$f(r) = a + br \quad (37)$$

Then it must fulfill

$$\begin{aligned} f(r = -1) &= a - b = f_1 \\ f(r = +1) &= a + b = f_2 \end{aligned}$$

This leads to

$$a = \frac{1}{2}(f_1 + f_2) \quad b = \frac{1}{2}(-f_1 + f_2)$$

and then replacing a, b in Eq. (37) by the above values on gets

$$f(r) = \left[\frac{1}{2}(1 - r) \right] f_1 + \left[\frac{1}{2}(1 + r) \right] f_2$$

or

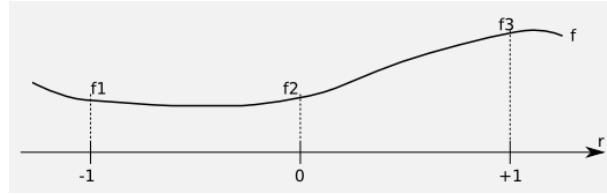
$$f(r) = \sum_{i=1}^2 N_i(r) f_i$$

with

$N_1(r) = \frac{1}{2}(1-r)$ $N_2(r) = \frac{1}{2}(1+r)$	(38)
---	------

4.4.2 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

$$f(r) = a + br + cr^2 \quad (39)$$

Then it must fulfill

$$\begin{aligned} f(r = -1) &= a - b + c = f_1 \\ f(r = 0) &= a = f_2 \\ f(r = +1) &= a + b + c = f_3 \end{aligned}$$

This leads to

$$a = f_2 \quad b = \frac{1}{2}(-f_1 + f_3) \quad c = \frac{1}{2}(f_1 + f_3 - 2f_2)$$

and then replacing a, b, c in Eq. (39) by the above values gets

$$f(r) = \left[\frac{1}{2}r(r-1) \right] f_1 + (1-r^2)f_2 + \left[\frac{1}{2}r(r+1) \right] f_3$$

or,

$$f(r) = \sum_{i=1}^3 N_i(r) f_i$$

with

$N_1(r) = \frac{1}{2}r(r-1)$ $N_2(r) = (1-r^2)$ $N_3(r) = \frac{1}{2}r(r+1)$	(40)
--	------

4.4.3 Cubic basis functions (Q_3)

The 1D basis polynomial is given by

$$f(r) = a + br + cr^2 + dr^3$$

with the nodes at position -1,-1/3, +1/3 and +1.

$$\begin{aligned} f(-1) &= a - b + c - d = f_1 \\ f(-1/3) &= a - \frac{b}{3} + \frac{c}{9} - \frac{d}{27} = f_2 \\ f(+1/3) &= a - \frac{b}{3} + \frac{c}{9} - \frac{d}{27} = f_3 \\ f(+1) &= a + b + c + d = f_4 \end{aligned}$$

Adding the first and fourth equation and the second and third, one arrives at

$$f_1 + f_4 = 2a + 2c \quad f_2 + f_3 = 2a + \frac{2c}{9}$$

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

$$2b + 2d = f_4 - f_1 \quad \frac{2b}{3} + \frac{2d}{27} = f_3 - f_2$$

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,

$$\begin{aligned} f(r) &= a + b + cr^2 + dr^3 \\ &= \frac{1}{16} (-1 + r + 9r^2 - 9r^3) f_1 \\ &\quad + \frac{1}{16} (9 - 27r - 9r^2 + 27r^3) f_2 \\ &\quad + \frac{1}{16} (9 + 27r - 9r^2 - 27r^3) f_3 \\ &\quad + \frac{1}{16} (-1 - r + 9r^2 + 9r^3) f_4 \\ &= \sum_{i=1}^4 N_i(r) f_i \end{aligned}$$

where

$$\begin{aligned}
N_1 &= \frac{1}{16}(-1 + r + 9r^2 - 9r^3) \\
N_2 &= \frac{1}{16}(9 - 27r - 9r^2 + 27r^3) \\
N_3 &= \frac{1}{16}(9 + 27r - 9r^2 - 27r^3) \\
N_4 &= \frac{1}{16}(-1 - r + 9r^2 + 9r^3)
\end{aligned}$$

Verification:

- Let us assume $f(r) = C$, then

$$\hat{f}(r) = \sum_i N_i(r) f_i = \sum_i N_i C = C \sum_i N_i = C$$

so that a constant function is exactly reproduced, as expected.

- 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

$$\begin{aligned}
\hat{f}(r) &= \sum_i N_i(r) f_i \\
&= -N_1(r) - \frac{1}{3}N_2(r) + \frac{1}{3}N_3(r) + N_4(r) \\
&= [-(-1 + r + 9r^2 - 9r^3) \\
&\quad - \frac{1}{3}(9 - 27r - 9r^2 + 27r^3) \\
&\quad + \frac{1}{3}(9 + 27r - 9r^2 - 27r^3) \\
&\quad + (-1 - r + 9r^2 + 9r^3)]/16 \\
&= [-r + 9r + 9r - r]/16 + \dots 0... \\
&= r
\end{aligned} \tag{41}$$

The basis functions derivative are given by

$$\begin{aligned}
\frac{\partial N_1}{\partial r} &= \frac{1}{16}(1 + 18r - 27r^2) \\
\frac{\partial N_2}{\partial r} &= \frac{1}{16}(-27 - 18r + 81r^2) \\
\frac{\partial N_3}{\partial r} &= \frac{1}{16}(+27 - 18r - 81r^2) \\
\frac{\partial N_4}{\partial r} &= \frac{1}{16}(-1 + 18r + 27r^2)
\end{aligned}$$

Verification:

- Let us assume $f(r) = C$, then

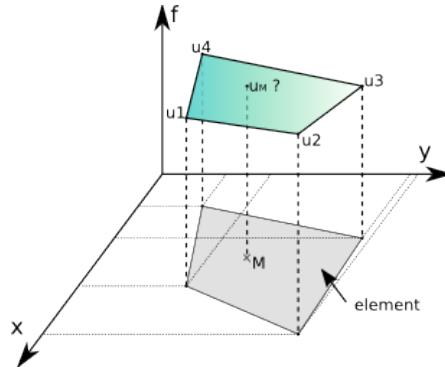
$$\begin{aligned}
\frac{\partial \hat{f}}{\partial r} &= \sum_i \frac{\partial N_i}{\partial r} f_i \\
&= C \sum_i \frac{\partial N_i}{\partial r} \\
&= \frac{C}{16} [(1 + 18r - 27r^2) \\
&\quad + (-27 - 18r + 81r^2) \\
&\quad + (+27 - 18r - 81r^2) \\
&\quad + (-1 + 18r + 27r^2)] \\
&= 0
\end{aligned}$$

- 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

$$\begin{aligned}
\frac{\partial \hat{f}}{\partial r} &= \sum_i \frac{\partial N_i}{\partial r} f_i \\
&= \frac{1}{16} [-(1 + 18r - 27r^2) \\
&\quad - \frac{1}{3}(-27 - 18r + 81r^2) \\
&\quad + \frac{1}{3}(+27 - 18r - 81r^2) \\
&\quad + (-1 + 18r + 27r^2)] \\
&= \frac{1}{16} [-2 + 18 + 54r^2 - 54r^2] \\
&= 1
\end{aligned}$$

4.5 Elements and basis functions in 2D

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 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 $u_M = u(x_M, y_M)$ 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. 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 $u_M \rightarrow u_2$.

In light of this, we could now assume that the weights would depend on the position of M in a continuous fashion:

$$u(x_M, y_M) = \sum_{i=1}^4 N_i(x_M, y_M) u_i$$

where the N_i are continuous ("well behaved") functions which have the property:

$$N_i(x_j, y_j) = \delta_{ij}$$

or, in other words:

$$N_3(x_1, y_1) = 0 \quad (42)$$

$$N_3(x_2, y_2) = 0 \quad (43)$$

$$N_3(x_3, y_3) = 1 \quad (44)$$

$$N_3(x_4, y_4) = 0 \quad (45)$$

The functions N_i are commonly called basis functions.

Omitting the M subscripts for any point inside the element, the velocity components u and v are given by:

$$\hat{u}(x, y) = \sum_{i=1}^4 N_i(x, y) u_i \quad (46)$$

$$\hat{v}(x, y) = \sum_{i=1}^4 N_i(x, y) v_i \quad (47)$$

Rather interestingly, 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 differentiable):

$$\dot{\epsilon}_{xx}(x, y) = \frac{\partial u}{\partial x} = \sum_{i=1}^4 \frac{\partial N_i}{\partial x} u_i \quad (48)$$

$$\dot{\epsilon}_{yy}(x, y) = \frac{\partial v}{\partial y} = \sum_{i=1}^4 \frac{\partial N_i}{\partial y} v_i \quad (49)$$

$$\dot{\epsilon}_{xy}(x, y) = \frac{1}{2} \frac{\partial u}{\partial y} + \frac{1}{2} \frac{\partial v}{\partial x} = \frac{1}{2} \sum_{i=1}^4 \frac{\partial N_i}{\partial y} u_i + \frac{1}{2} \sum_{i=1}^4 \frac{\partial N_i}{\partial x} v_i \quad (50)$$

How we actually obtain the exact form of the basis functions is explained in the coming section.

4.5.1 Bilinear basis functions in 2D (Q_1)

In this section, we place ourselves in the most favorable case, i.e. the element is a square defined by $-1 < r < 1$, $-1 < s < 1$ in the Cartesian coordinates system (r, s) :

```
3=====2
|           |   (r_0,s_0)=(-1,-1)
|           |   (r_1,s_1)=(+1,-1)
|           |   (r_2,s_2)=(+1,+1)
|           |   (r_3,s_3)=(-1,+1)
|
0=====1
```

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 later on. For now, the basis functions in the above reference element and in the reduced coordinates system (r, s) are given by:

$N_1(r, s) = 0.25(1 - r)(1 - s)$
$N_2(r, s) = 0.25(1 + r)(1 - s)$
$N_3(r, s) = 0.25(1 + r)(1 + s)$
$N_4(r, s) = 0.25(1 - r)(1 + s)$

The partial derivatives of these functions with respect to r and s automatically follow:

$\frac{\partial N_1}{\partial r}(r, s) = -0.25(1 - s)$	$\frac{\partial N_1}{\partial s}(r, s) = -0.25(1 - r)$
$\frac{\partial N_2}{\partial r}(r, s) = +0.25(1 - s)$	$\frac{\partial N_2}{\partial s}(r, s) = -0.25(1 + r)$
$\frac{\partial N_3}{\partial r}(r, s) = +0.25(1 + s)$	$\frac{\partial N_3}{\partial s}(r, s) = +0.25(1 + r)$
$\frac{\partial N_4}{\partial r}(r, s) = -0.25(1 + s)$	$\frac{\partial N_4}{\partial s}(r, s) = +0.25(1 - r)$

Let us go back to Eq.(47). 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

$$\hat{v}(r, s) = \sum_{i=1}^4 N_i(r, s) v_i = C \sum_{i=1}^4 N_i(r, s) = C[N_1(r, s) + N_2(r, s) + N_3(r, s) + N_4(r, s)] = C$$

This is a very important property: if the v function used to assign values at the vertices is constant, then the value of \hat{v} anywhere in the element is exactly C . If we now turn to the derivatives of v with respect to r and s :

$$\frac{\partial \hat{v}}{\partial r}(r, s) = \sum_{i=1}^4 \frac{\partial N_i}{\partial r}(r, s) v_i = C \sum_{i=1}^4 \frac{\partial N_i}{\partial r}(r, s) = C[-0.25(1 - s) + 0.25(1 - s) + 0.25(1 + s) - 0.25(1 + s)] = 0$$

$$\frac{\partial \hat{v}}{\partial s}(r, s) = \sum_{i=1}^4 \frac{\partial N_i}{\partial s}(r, s) v_i = C \sum_{i=1}^4 \frac{\partial N_i}{\partial s}(r, s) = C[-0.25(1 - r) - 0.25(1 + r) + 0.25(1 + r) + 0.25(1 - r)] = 0$$

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:

$$\hat{v}(r, s) = \sum_{i=1}^4 N_i(r, s) v_i \tag{51}$$

$$= \sum_{i=1}^4 N_i(r, s)(ar_i + bs_i) \tag{52}$$

$$= \underbrace{a \sum_{i=1}^4 N_i(r, s)r_i}_{r} + \underbrace{b \sum_{i=1}^4 N_i(r, s)s_i}_{s} \tag{53}$$

$$= a[0.25(1 - r)(1 - s)(-1) + 0.25(1 + r)(1 - s)(+1) + 0.25(1 + r)(1 + s)(+1) + 0.25(1 - r)(1 + s)(-1)]$$

$$+ b[0.25(1 - r)(1 - s)(-1) + 0.25(1 + r)(1 - s)(-1) + 0.25(1 + r)(1 + s)(+1) + 0.25(1 - r)(1 + s)(+1)]$$

$$= a[-0.25(1 - r)(1 - s) + 0.25(1 + r)(1 - s) + 0.25(1 + r)(1 + s) - 0.25(1 - r)(1 + s)]$$

$$+ b[-0.25(1 - r)(1 - s) - 0.25(1 + r)(1 - s) + 0.25(1 + r)(1 + s) + 0.25(1 - r)(1 + s)]$$

$$= ar + bs \tag{54}$$

verify above eq. This set of bilinear shape functions is therefore capable of exactly representing a bilinear field. The derivatives are:

$$\frac{\partial \hat{v}}{\partial r}(r, s) = \sum_{i=1}^4 \frac{\partial N_i}{\partial r}(r, s) v_i \quad (55)$$

$$= a \sum_{i=1}^4 \frac{\partial N_i}{\partial r}(r, s) r_i + b \sum_{i=1}^4 \frac{\partial N_i}{\partial r}(r, s) s_i \quad (56)$$

$$= a [-0.25(1-s)(-1) + 0.25(1-s)(+1) + 0.25(1+s)(+1) - 0.25(1+s)(-1)]$$

$$+ b [-0.25(1-s)(-1) + 0.25(1-s)(-1) + 0.25(1+s)(+1) - 0.25(1+s)(+1)]$$

$$= \frac{a}{4} [(1-s) + (1-s) + (1+s) + (1+s)]$$

$$+ \frac{b}{4} [(1-s) - (1-s) + (1+s) - (1+s)]$$

$$= a \quad (57)$$

Here again, we find that the derivative of the bilinear field inside the element is exact: $\frac{\partial \hat{v}}{\partial r} = \frac{\partial v}{\partial r}$.

However, following the same methodology as above, one can easily prove that this is no more true for polynomials of degree strivtly higher than 1. This fact has serious consequences: if the solution to the problem at hand is for instance a parabola, the Q_1 shape 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 themselves quadratic terms.

4.5.2 Biquadratic basis functions in 2D (Q_2)

This element is part of the so-called LAgrange family.

citation needed

Inside an element the local numbering of the nodes is as follows:

```

3=====6=====2
|       |       |   (r_0,s_0)=(-1,-1)   (r_4,s_4)=( 0,-1)
|       |       |   (r_1,s_1)=(+1,-1)   (r_5,s_5)=(+1, 0)
7=====8=====5   (r_2,s_2)=(+1,+1)   (r_6,s_6)=( 0,+1)
|       |       |   (r_3,s_3)=(-1,+1)   (r_7,s_7)=(-1, 0)
|       |       |                           (r_8,s_8)=( 0, 0)
0=====4=====1

```

The basis polynomial is then

$$f(r, s) = a + br + cs + drs + er^2 + fs^2 + gr^2s + hrs^2 + ir^2s^2$$

The velocity shape functions are given by:

$$\begin{aligned}
N_0(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s-1) \\
N_1(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s-1) \\
N_2(r, s) &= \frac{1}{2}r(r+1)\frac{1}{2}s(s+1) \\
N_3(r, s) &= \frac{1}{2}r(r-1)\frac{1}{2}s(s+1) \\
N_4(r, s) &= (1-r^2)\frac{1}{2}s(s-1) \\
N_5(r, s) &= \frac{1}{2}r(r+1)(1-s^2) \\
N_6(r, s) &= (1-r^2)\frac{1}{2}s(s+1) \\
N_7(r, s) &= \frac{1}{2}r(r-1)(1-s^2) \\
N_8(r, s) &= (1-r^2)(1-s^2)
\end{aligned}$$

and their derivatives by:

$$\begin{aligned}
\frac{\partial N_0}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s-1) & \frac{\partial N_0}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s-1) \\
\frac{\partial N_1}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s-1) & \frac{\partial N_1}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s-1) \\
\frac{\partial N_2}{\partial r} &= \frac{1}{2}(2r+1)\frac{1}{2}s(s+1) & \frac{\partial N_2}{\partial s} &= \frac{1}{2}r(r+1)\frac{1}{2}(2s+1) \\
\frac{\partial N_3}{\partial r} &= \frac{1}{2}(2r-1)\frac{1}{2}s(s+1) & \frac{\partial N_3}{\partial s} &= \frac{1}{2}r(r-1)\frac{1}{2}(2s+1) \\
\frac{\partial N_4}{\partial r} &= (-2r)\frac{1}{2}s(s-1) & \frac{\partial N_4}{\partial s} &= (1-r^2)\frac{1}{2}(2s-1) \\
\frac{\partial N_5}{\partial r} &= \frac{1}{2}(2r+1)(1-s^2) & \frac{\partial N_5}{\partial s} &= \frac{1}{2}r(r+1)(-2s) \\
\frac{\partial N_6}{\partial r} &= (-2r)\frac{1}{2}s(s+1) & \frac{\partial N_6}{\partial s} &= (1-r^2)\frac{1}{2}(2s+1) \\
\frac{\partial N_7}{\partial r} &= \frac{1}{2}(2r-1)(1-s^2) & \frac{\partial N_7}{\partial s} &= \frac{1}{2}r(r-1)(-2s) \\
\frac{\partial N_8}{\partial r} &= (-2r)(1-s^2) & \frac{\partial N_8}{\partial s} &= (1-r^2)(-2s)
\end{aligned}$$

4.5.3 Eight node serendipity basis functions in 2D ($Q_2^{(8)}$)

Inside an element the local numbering of the nodes is as follows:

```

3=====6=====2
|       |       |   (r_0,s_0)=(-1,-1)   (r_4,s_4)=( 0,-1)
|       |       |   (r_1,s_1)=(+1,-1)   (r_5,s_5)=(+1, 0)
7=====+=====5   (r_2,s_2)=(+1,+1)   (r_6,s_6)=( 0,+1)
|       |       |   (r_3,s_3)=(-1,+1)   (r_7,s_7)=(-1, 0)
|
0=====4=====1

```

The main difference with the Q_2 element resides in the fact that there is no node in the middle of the element. The basis polynomial is then

$$f(r, s) = a + br + cs + drs + er^2 + fs^2 + gr^2s + hrs^2$$

Note that absence of the r^2s^2 term which was previously associated to the center node. We find that

$$N_0(r, s) = \frac{1}{4}(1-r)(1-s)(-r-s-1) \quad (58)$$

$$N_1(r, s) = \frac{1}{4}(1+r)(1-s)(r-s-1) \quad (59)$$

$$N_2(r, s) = \frac{1}{4}(1+r)(1+s)(r+s-1) \quad (60)$$

$$N_3(r, s) = \frac{1}{4}(1-r)(1+s)(-r+s-1) \quad (61)$$

$$N_4(r, s) = \frac{1}{2}(1-r^2)(1-s) \quad (62)$$

$$N_5(r, s) = \frac{1}{2}(1+r)(1-s^2) \quad (63)$$

$$N_6(r, s) = \frac{1}{2}(1-r^2)(1+s) \quad (64)$$

$$N_7(r, s) = \frac{1}{2}(1-r)(1-s^2) \quad (65)$$

The shape 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 shape functions for corner nodes are modifications of the bilinear quadrilateral element:

[verify those](#)

4.5.4 Bicubic basis functions in 2D (Q_3)

Inside an element the local numbering of the nodes is as follows:

12==13==14==15	$(r, s)_{\{00\}}=(-1, -1)$	$(r, s)_{\{08\}}=(-1, +1/3)$
	$(r, s)_{\{01\}}=(-1/3, -1)$	$(r, s)_{\{09\}}=(-1/3, +1/3)$
08==09==10==11	$(r, s)_{\{02\}}=(+1/3, -1)$	$(r, s)_{\{10\}}=(+1/3, +1/3)$
	$(r, s)_{\{03\}}=(+1, -1)$	$(r, s)_{\{11\}}=(+1, +1/3)$
04==05==06==07	$(r, s)_{\{04\}}=(-1, -1/3)$	$(r, s)_{\{12\}}=(-1, +1)$
	$(r, s)_{\{05\}}=(-1/3, -1/3)$	$(r, s)_{\{13\}}=(-1/3, +1)$
00==01==02==03	$(r, s)_{\{06\}}=(+1/3, -1/3)$	$(r, s)_{\{14\}}=(+1/3, +1)$
	$(r, s)_{\{07\}}=(+1, -1/3)$	$(r, s)_{\{15\}}=(+1, +1)$

The velocity shape functions are given by:

$$N_1(r) = (-1 + r + 9r^2 - 9r^3)/16$$

$$N_1(t) = (-1 + t + 9t^2 - 9t^3)/16$$

$$N_2(r) = (+9 - 27r - 9r^2 + 27r^3)/16$$

$$N_2(t) = (+9 - 27t - 9t^2 + 27t^3)/16$$

$$N_3(r) = (+9 + 27r - 9r^2 - 27r^3)/16$$

$$N_3(t) = (+9 + 27t - 9t^2 - 27t^3)/16$$

$$N_4(r) = (-1 - r + 9r^2 + 9r^3)/16$$

$$N_4(t) = (-1 - t + 9t^2 + 9t^3)/16$$

$N_{01}(r,s) = N_1(r)N_1(s) = (-1 + r + 9r^2 - 9r^3)/16 * (-1 + t + 9s^2 - 9s^3)/16$	(66)
$N_{02}(r,s) = N_2(r)N_1(s) = (+9 - 27r - 9r^2 + 27r^3)/16 * (-1 + t + 9s^2 - 9s^3)/16$	(66)
$N_{03}(r,s) = N_3(r)N_1(s) = (+9 + 27r - 9r^2 - 27r^3)/16 * (-1 + t + 9s^2 - 9s^3)/16$	(66)
$N_{04}(r,s) = N_4(r)N_1(s) = (-1 - r + 9r^2 + 9r^3)/16 * (-1 + t + 9s^2 - 9s^3)/16$	(66)
$N_{05}(r,s) = N_1(r)N_2(s) = (-1 + r + 9r^2 - 9r^3)/16 * (9 - 27s - 9s^2 + 27s^3)/16$	(66)
$N_{06}(r,s) = N_2(r)N_2(s) = (+9 - 27r - 9r^2 + 27r^3)/16 * (9 - 27s - 9s^2 + 27s^3)/16$	(66)
$N_{07}(r,s) = N_3(r)N_2(s) = (+9 + 27r - 9r^2 - 27r^3)/16 * (9 - 27s - 9s^2 + 27s^3)/16$	(66)
$N_{08}(r,s) = N_4(r)N_2(s) = (-1 - r + 9r^2 + 9r^3)/16 * (9 - 27s - 9s^2 + 27s^3)/16$	(66)
$N_{09}(r,s) = N_1(r)N_3(s) =$	(66)
$N_{10}(r,s) = N_2(r)N_3(s) =$	(67)
$N_{11}(r,s) = N_3(r)N_3(s) =$	(68)
$N_{12}(r,s) = N_4(r)N_3(s) =$	(69)
$N_{13}(r,s) = N_1(r)N_4(s) =$	(70)
$N_{14}(r,s) = N_2(r)N_4(s) =$	(71)
$N_{15}(r,s) = N_3(r)N_4(s) =$	(72)
$N_{16}(r,s) = N_4(r)N_4(s) =$	(73)

4.5.5 Linear basis functions for triangles in 2D (P_1)

```

2
\ \
|   \      (r_0,s_0)=(0,0)
|   \      (r_1,s_1)=(1,0)
|   \      (r_2,s_2)=(0,2)
0=====1

```

The basis polynomial is then

$$f(r,s) = a + br + cs$$

and the shape functions:

$$N_0(r,s) = 1 - r - s \quad (74)$$

$$N_1(r,s) = r \quad (75)$$

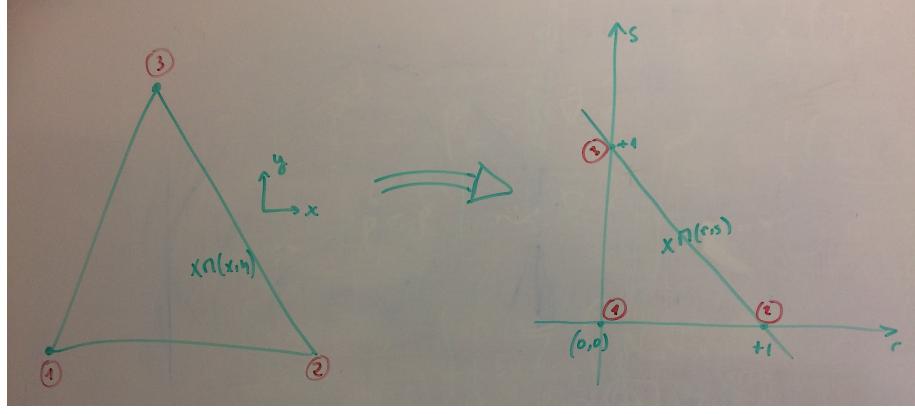
$$N_2(r,s) = s \quad (76)$$

4.5.6 Enriched linear basis functions in triangles (P_1^+)

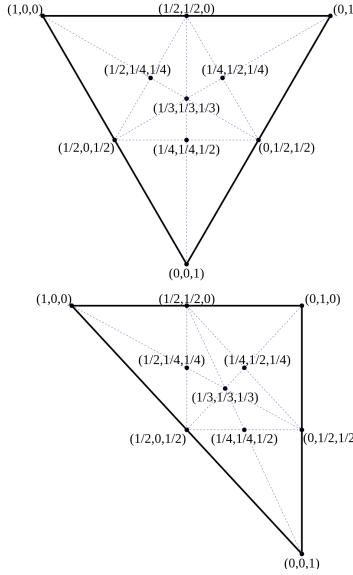
As we will see in Section 6.2.7 the above P_1 can be enriched with a so-called bubble function. The bubble function of the MINI element is described in [22] as being $\lambda_1\lambda_2\lambda_3$ where λ_i are the so-called barycentric coordinates³.

$$\begin{aligned} \lambda_1 &= \frac{(y_2 - y_3)(x - x_3) + (x_3 - x_2)(y - y_3)}{(y_2 - y_3)(x_1 - x_3) + (x_3 - x_2)(y_1 - y_3)} \\ \lambda_2 &= \frac{(y_3 - y_1)(x - x_3) + (x_1 - x_3)(y - y_3)}{(y_2 - y_3)(x_1 - x_3) + (x_3 - x_2)(y_1 - y_3)} \\ \lambda_3 &= 1 - \lambda_1 - \lambda_2 \end{aligned}$$

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



representation of the element in the real coordinate system (x, y) and in the reduced coordinate system (r, s)



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 shape 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

$$d = 27(f_4 - a - \frac{b}{3} - \frac{c}{3}) = 27(f_4 - f_1 - \frac{f_2 - f_1}{3} - \frac{f_3 - f_1}{3}) = 27(f_4 - \frac{f_1}{3} - \frac{f_2}{3} - \frac{f_3}{3})$$

Finally

$$\begin{aligned} f(r, s) &= a + br + cs + d(1 - r - s)rs \\ &= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + 27(f_4 - \frac{f_1}{3} - \frac{f_2}{3} - \frac{f_3}{3})(1 - r - s)rs \\ &= [1 - r - s - 9(1 - r - s)rs]f_1 + [r - 9(1 - r - s)rs]f_2 + [s - 9(1 - r - s)rs]f_3 + [27(1 - r - s)rs]f_4 \end{aligned}$$

so that

$$f(r, s) = \sum_{i=1}^4 N_i(r, s) f_i$$

with

$$\begin{aligned} N_1(r, s) &= 1 - r - s - 9(1 - r - s)rs \\ N_2(r, s) &= r - 9(1 - r - s)rs \\ N_3(r, s) &= s - 9(1 - r - s)rs \\ N_4(r, s) &= 27(1 - r - s)rs \end{aligned}$$

It is trivial to verify that $\sum_i N_i = 1$ for all values of r, s and the gradients of the shape functions are:

$$\frac{\partial N_1}{\partial r}(r, s) = -1 - 9(1 - 2r - s)s \quad (77)$$

$$\frac{\partial N_2}{\partial r}(r, s) = +1 - 9(1 - 2r - s)s \quad (78)$$

$$\frac{\partial N_3}{\partial r}(r, s) = -9(1 - 2r - s)s \quad (79)$$

$$\frac{\partial N_4}{\partial r}(r, s) = 27(1 - 2r - s)s \quad (80)$$

(81)

$$\frac{\partial N_1}{\partial s}(r, s) = -1 - 9(1 - r - 2s)r \quad (82)$$

$$\frac{\partial N_2}{\partial s}(r, s) = -9(1 - r - 2s)r \quad (83)$$

$$\frac{\partial N_3}{\partial s}(r, s) = +1 - 9(1 - r - 2s)r \quad (84)$$

$$\frac{\partial N_4}{\partial s}(r, s) = 27(1 - r - 2s)r \quad (85)$$

We have two coordinate systems for the element: the global coordinates (x, y) and the natural coordinates (r, s) . Inside the element, the relation between the two is given by

$$\begin{aligned} x &= N_1 x_1 + N_2 x_2 + N_3 x_3 + N_4 x_4 = \sum_i N_i(r, s) x_i \\ y &= N_1 y_1 + N_2 y_2 + N_3 y_3 + N_4 y_4 = \sum_i N_i(r, s) y_i \end{aligned} \quad (86)$$

or,

$$\begin{aligned}
x &= [1 - r - s - 9(1 - r - s)rs]x_1 + [r - 9(1 - r - s)rs]x_2 + [s - 9(1 - r - s)rs]x_3 + [27(1 - r - s)rs]x_4 \\
&= x_1 - r(x_1 - x_2) - s(x_1 - x_3) + (1 - r - s)rs(-9x_1 - 9x_2 - 9x_3 + 27x_4) \\
&= x_1 - r(x_1 - x_2) - s(x_1 - x_3) + (1 - r - s)rs(-9x_1 - 9x_2 - 9x_3 + 27(x_1 + x_2 + x_3)/3) \\
&= x_1 - r(x_1 - x_2) - s(x_1 - x_3) \\
&= x_1 - rx_{12} - sx_{13} \\
y &= [1 - r - s - 9(1 - r - s)rs]y_1 + [r - 9(1 - r - s)rs]y_2 + [s - 9(1 - r - s)rs]y_3 + [27(1 - r - s)rs]y_4 \\
&= y_1 - r(y_1 - y_2) - s(y_1 - y_3) + (1 - r - s)rs(-9y_1 - 9y_2 - 9y_3 + 27y_4) \\
&= y_1 - r(y_1 - y_2) - s(y_1 - y_3) + (1 - r - s)rs(-9y_1 - 9y_2 - 9y_3 + 27(y_1 + y_2 + y_3)/3) \\
&= y_1 - r(y_1 - y_2) - s(y_1 - y_3) \\
&= y_1 - ry_{12} - sy_{13}
\end{aligned}$$

4.5.7 Quadratic basis functions for triangles in 2D (P_2)

```

2
| \
|   \      (r_0,s_0)=(0,0)  (r_3,s_3)=(1/2,0)
5   4      (r_1,s_1)=(1,0)  (r_4,s_4)=(1/2,1/2)
|     \    (r_2,s_2)=(0,1)  (r_5,s_5)=(0,1/2)
|       \
0====3====1

```

The basis polynomial is then

$$f(r, s) = c_1 + c_2r + c_3s + c_4r^2 + c_5rs + c_6s^2$$

We have

$$\begin{aligned}
f_1 = f(r_1, s_1) &= c_1 \\
f_2 = f(r_2, s_2) &= c_1 + c_2 + c_4 \\
f_3 = f(r_3, s_3) &= c_1 + c_3 + c_6 \\
f_4 = f(r_4, s_4) &= c_1 + c_2/2 + c_4/4 \\
f_5 = f(r_5, s_5) &= c_1 + c_2/2 + c_3/2 \\
&\quad + c_4/4 + c_5/4 + c_6/4 \\
f_6 = f(r_6, s_6) &= c_1 + c_3/2 + c_6/4
\end{aligned}$$

This can be cast as $\mathbf{f} = \mathbf{A} \cdot \mathbf{c}$ where \mathbf{A} is a 6x6 matrix:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1/2 & 0 & 1/4 & 0 & 0 \\ 1 & 1/2 & 1/2 & 1/4 & 1/4 & 1/4 \\ 1 & 0 & 1/2 & 0 & 0 & 1/4 \end{pmatrix}$$

It is rather trivial to compute the inverse of this matrix:

$$\mathbf{A}^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ -3 & -1 & 0 & 4 & 0 & 0 \\ -3 & 0 & -1 & 0 & 0 & 4 \\ 2 & 2 & 0 & -4 & 0 & 0 \\ 4 & 0 & 0 & -4 & 4 & -4 \\ 2 & 0 & 2 & 0 & 0 & -4 \end{pmatrix}$$

In the end, one obtains:

$$\begin{aligned}
f(r, s) &= f_1 + (-3f_1 - f_2 + 4f_4)r + (-3f_1 - f_3 + 4f_6)s \\
&\quad + (2f_1 + 2f_2 - 4f_4)r^2 + (4f_1 - 4f_4 + 4f_5 - 4f_6)rs \\
&\quad + (2f_1 + 2f_3 - 4f_6)s^2 \\
&= \sum_{i=1}^6 N_i(r, s)f_i
\end{aligned} \tag{87}$$

with

$N_1(r, s)$	$=$	$1 - 3r - 3s + 2r^2 + 4rs + 2s^2$
$N_2(r, s)$	$=$	$-r + 2r^2$
$N_3(r, s)$	$=$	$-s + 2s^2$
$N_4(r, s)$	$=$	$4r - 4r^2 - 4rs$
$N_5(r, s)$	$=$	$4rs$
$N_6(r, s)$	$=$	$4s - 4rs - 4s^2$

4.5.8 Enriched quadratic basis functions in triangles (P_2^+)

This is used by the Crouzeix-Raviart element, see Section 6.2.9.

```

03          (r_1,s_1)=(0,0)
||\ \
||  \ \
||   \ \
06   05      (r_5,s_5)=(1/2,1/2)
|| 07  \ \
||    \ \
01==04==02

```

The shape functions are given by:

[find reference](#)

$$N_1(r, s) = (1 - r - s)(1 - 2r - 2s + 3rs) \tag{88}$$

$$N_2(r, s) = r(2r - 1 + 3s - 3rs - 3s^2) \tag{89}$$

$$N_3(r, s) = s(2s - 1 + 3r - 3r^2 - 3rs) \tag{90}$$

$$N_4(r, s) = 4(1 - r - s)r(1 - 3s) \tag{91}$$

$$N_5(r, s) = 4rs[-2 + 3r + 3s] \tag{92}$$

$$N_6(r, s) = 4(1 - r - s)s(1 - 3r) \tag{93}$$

$$N_7(r, s) = 27(1 - r - s)rs \tag{94}$$

It is then easy to verify that for all shape functions we have $N_i(r_j, s_j) = \delta_{ij}$ where j denotes one of the seven nodes.

The derivatives are as follows:

$$\frac{\partial N_1}{\partial r}(r, s) = r(4 - 6s) - 3s^2 + 7s - 3 \quad (95)$$

$$\frac{\partial N_2}{\partial r}(r, s) = r(4 - 6s) - 3s^2 + 3s - 1 \quad (96)$$

$$\frac{\partial N_3}{\partial r}(r, s) = -3s(2r + s - 1) \quad (97)$$

$$\frac{\partial N_4}{\partial r}(r, s) = 4(3s - 1)(2r + s - 1) \quad (98)$$

$$\frac{\partial N_5}{\partial r}(r, s) = 4s(6r + 3s - 2) \quad (99)$$

$$\frac{\partial N_6}{\partial r}(r, s) = 4s(6r + 3s - 4) \quad (100)$$

$$\frac{\partial N_7}{\partial r}(r, s) = -27s(2r + s - 1) \quad (101)$$

$$\frac{\partial N_1}{\partial s}(r, s) = -3r^2 + r(7 - 6s) + 4s - 3 \quad (102)$$

$$\frac{\partial N_2}{\partial s}(r, s) = -3r(r + 2s - 1) \quad (103)$$

$$\frac{\partial N_3}{\partial s}(r, s) = -3r^2 + r(3 - 6s) + 4s - 1 \quad (104)$$

$$\frac{\partial N_4}{\partial s}(r, s) = 4r(3r + 6s - 4) \quad (105)$$

$$\frac{\partial N_5}{\partial s}(r, s) = 4r(3r + 6s - 2) \quad (106)$$

$$\frac{\partial N_6}{\partial s}(r, s) = 4(3r - 1)(r + 2s - 1) \quad (107)$$

$$\frac{\partial N_7}{\partial s}(r, s) = -27r(r + 2s - 1) \quad (108)$$

Note that the shape functions can also be expressed as a function of the barycentric coordinates, as in the MILAMIN code [151] or in Cuvelier et al, 1986 [150]⁴

```
03
|| \\
||  \\
||  \\
05  04
|| 07 \\
||      \\
01==06==02
```

$$N_1(\lambda_1, \lambda_2, \lambda_3) = \eta_1(2\eta_1 - 1) + 3\eta_1\eta_2\eta_3 \quad (109)$$

$$N_2(\lambda_1, \lambda_2, \lambda_3) = \eta_2(2\eta_2 - 1) + 3\eta_1\eta_2\eta_3 \quad (110)$$

$$N_3(\lambda_1, \lambda_2, \lambda_3) = \eta_3(2\eta_3 - 1) + 3\eta_1\eta_2\eta_3 \quad (111)$$

$$N_4(\lambda_1, \lambda_2, \lambda_3) = 4\eta_2\eta_3 - 12\eta_1\eta_2\eta_3 \quad (112)$$

$$N_5(\lambda_1, \lambda_2, \lambda_3) = 4\eta_1\eta_3 - 12\eta_1\eta_2\eta_3 \quad (113)$$

$$N_6(\lambda_1, \lambda_2, \lambda_3) = 4\eta_1\eta_2 - 12\eta_1\eta_2\eta_3 \quad (114)$$

$$N_7(\lambda_1, \lambda_2, \lambda_3) = 27\eta_1\eta_2\eta_3 \quad (115)$$

VERIFY that when $\eta_1 = 1 - r - s$, $\eta_2 = r$ and $\eta_3 = s$ we find the above r, s shape functions

⁴Note that the numbering of the nodes in the book is different with respect to the one above.

4.5.9 Cubic basis functions for triangles (P_3)

```

2
|\       (r_0,s_0)=(0,0)   (r_5,s_5)=(2/3,1/3)
| \     (r_1,s_1)=(1,0)   (r_6,s_6)=(1/3,2/3)
7   6   (r_2,s_2)=(0,1)   (r_7,s_7)=(0,2/3)
|   \
8   9   5   (r_3,s_3)=(1/3,0) (r_8,s_8)=(0,1/3)
|       \
0==3==4==1

```

The basis polynomial is then

$$f(r,s) = c_1 + c_2r + c_3s + c_4r^2 + c_5rs + c_6s^2 + c_7r^3 + c_8r^2s + c_9rs^2 + c_{10}s^3$$

$$N_0(r,s) = \frac{9}{2}(1-r-s)(1/3-r-s)(2/3-r-s) \quad (116)$$

$$N_1(r,s) = \frac{9}{2}r(r-1/3)(r-2/3) \quad (117)$$

$$N_2(r,s) = \frac{9}{2}s(s-1/3)(s-2/3) \quad (118)$$

$$N_3(r,s) = \frac{27}{2}(1-r-s)r(2/3-r-s) \quad (119)$$

$$N_4(r,s) = \frac{27}{2}(1-r-s)r(r-1/3) \quad (120)$$

$$N_5(r,s) = \frac{27}{2}rs(r-1/3) \quad (121)$$

$$N_6(r,s) = \frac{27}{2}rs(r-2/3) \quad (122)$$

$$N_7(r,s) = \frac{27}{2}(1-r-s)s(s-1/3) \quad (123)$$

$$N_8(r,s) = \frac{27}{2}(1-r-s)s(2/3-r-s) \quad (124)$$

$$N_9(r,s) = 27rs(1-r-s) \quad (125)$$

verify those

4.6 Elements and basis functions in 3D

4.6.1 Linear basis functions in tetrahedra (P_1)

$$(r_0,s_0) = (0,0,0)$$

$$(r_1,s_1) = (1,0,0)$$

$$(r_2,s_2) = (0,2,0)$$

$$(r_3,s_3) = (0,0,1)$$

The basis polynomial is given by

$$f(r,s,t) = c_0 + c_1r + c_2s + c_3t$$

$$f_1 = f(r_1,s_1,t_1) = c_0 \quad (126)$$

$$f_2 = f(r_2,s_2,t_2) = c_0 + c_1 \quad (127)$$

$$f_3 = f(r_3,s_3,t_3) = c_0 + c_2 \quad (128)$$

$$f_4 = f(r_4,s_4,t_4) = c_0 + c_3 \quad (129)$$

which yields:

$$c_0 = f_1 \quad c_1 = f_2 - f_1 \quad c_2 = f_3 - f_1 \quad c_3 = f_4 - f_1$$

$$\begin{aligned} f(r, s, t) &= c_0 + c_1 r + c_2 s + c_3 t \\ &= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + (f_4 - f_1)t \\ &= f_1(1 - r - s - t) + f_2 r + f_3 s + f_4 t \\ &= \sum_i N_i(r, s, t) f_i \end{aligned}$$

Finally,

\$N_1(r, s, t) = 1 - r - s - t\$
\$N_2(r, s, t) = r\$
\$N_3(r, s, t) = s\$
\$N_4(r, s, t) = t\$

4.6.2 Enriched linear in tetrahedra(P_1^+)

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

$$f(r_1, s_1, t_1) = f_1 = a + b r_1 + c s_1 + d t_1 + e r_1 s_1 t_1 (1 - r_1 - s_1 - t_1) \quad (130)$$

$$f(r_2, s_2, t_2) = f_2 = a + b r_2 + c s_2 + d t_2 + e r_2 s_2 t_2 (1 - r_2 - s_2 - t_2) \quad (131)$$

$$f(r_3, s_3, t_3) = f_3 = a + b r_3 + c s_3 + d t_3 + e r_3 s_3 t_3 (1 - r_3 - s_3 - t_3) \quad (132)$$

$$f(r_4, s_4, t_4) = f_4 = a + b r_4 + c s_4 + d t_4 + e r_4 s_4 t_4 (1 - r_4 - s_4 - t_4) \quad (133)$$

$$f(r_5, s_5, t_5) = f_5 = a + b r_5 + c s_5 + d t_5 + e r_5 s_5 t_5 (1 - r_5 - s_5 - t_5) \quad (134)$$

i.e.,

$$f_1 = a \quad (135)$$

$$f_2 = a + b \quad (136)$$

$$f_3 = a + c \quad (137)$$

$$f_4 = a + d \quad (138)$$

$$f_5 = a + b/4 + c/4 + d/4 + e/64(1 - 1/4 - 1/4 - 1/4) \quad (139)$$

$$= a + b/4 + c/4 + d/4 + e/256 \quad (140)$$

Then

$$a = f_1 \quad (141)$$

$$b = f_2 - f_1 \quad (142)$$

$$c = f_3 - f_1 \quad (143)$$

$$d = f_4 - f_1 \quad (144)$$

$$e = 256(f_5 - a - b/4 - c/4 - d/4) \quad (145)$$

$$= 256(f_5 - f_1 - (f_2 - f_1)/4 - (f_3 - f_1)/4 - (f_4 - f_1)/4) \quad (146)$$

$$= 256(-f_1/4 - f_2/4 - f_3/4 - f_4/4 + f_5) \quad (147)$$

$$= 64(-f_1 - f_2 - f_3 - f_4 + 4f_5) \quad (148)$$

Finally:

$$\begin{aligned}
f(r, s, t) &= a + br + cs + dt + erst(1 - r - s - t) \\
&= f_1 + (f_2 - f_1)r + (f_3 - f_1)s + (f_4 - f_1)t + 64(-f_1 - f_2 - f_3 - f_4 + 4f_5)rst(1 - r - s - t) \\
&= f_1[1 - r - s - t - 64rst(1 - r - s - t)] \\
&\quad + f_2[r - 64rst(1 - r - s - t)] \\
&\quad + f_3[s - 64rst(1 - r - s - t)] \\
&\quad + f_4[t - 64rst(1 - r - s - t)] \\
&\quad + f_5[256rst(1 - r - s - t)] \\
&= \sum_{i=1}^5 N_i(r, s, t) f_i
\end{aligned} \tag{149}$$

with

$$N_1(r, s, t) = 1 - r - s - t - 64rst(1 - r - s - t) \tag{150}$$

$$N_2(r, s, t) = r - 64rst(1 - r - s - t) \tag{151}$$

$$N_3(r, s, t) = s - 64rst(1 - r - s - t) \tag{152}$$

$$N_4(r, s, t) = t - 64rst(1 - r - s - t) \tag{153}$$

$$N_5(r, s, t) = +256rst(1 - r - s - t) \tag{154}$$

The derivatives are given by:

$$\frac{\partial N_1}{\partial r}(r, s, t) = -1 - 64st(1 - 2r - s - t)$$

$$\frac{\partial N_2}{\partial r}(r, s, t) = +1 - 64st(1 - 2r - s - t)$$

$$\frac{\partial N_3}{\partial r}(r, s, t) = -64st(1 - 2r - s - t)$$

$$\frac{\partial N_4}{\partial r}(r, s, t) = -64st(1 - 2r - s - t)$$

$$\frac{\partial N_5}{\partial r}(r, s, t) = 256st(1 - 2r - s - t)$$

$$\frac{\partial N_1}{\partial s}(r, s, t) = -1 - 64rt(1 - r - 2s - t)$$

$$\frac{\partial N_2}{\partial s}(r, s, t) = -64rt(1 - r - 2s - t)$$

$$\frac{\partial N_3}{\partial s}(r, s, t) = +1 - 64rt(1 - r - 2s - t)$$

$$\frac{\partial N_4}{\partial s}(r, s, t) = -64rt(1 - r - 2s - t)$$

$$\frac{\partial N_5}{\partial s}(r, s, t) = 256rt(1 - r - 2s - t)$$

$$\frac{\partial N_1}{\partial t}(r, s, t) = -1 - 64rs(1 - r - s - 2t)$$

$$\frac{\partial N_2}{\partial t}(r, s, t) = -64rs(1 - r - s - 2t)$$

$$\frac{\partial N_3}{\partial t}(r, s, t) = -64rs(1 - r - s - 2t)$$

$$\frac{\partial N_4}{\partial t}(r, s, t) = +1 - 64rs(1 - r - s - 2t)$$

$$\frac{\partial N_5}{\partial t}(r, s, t) = 256rs(1 - r - s - 2t)$$

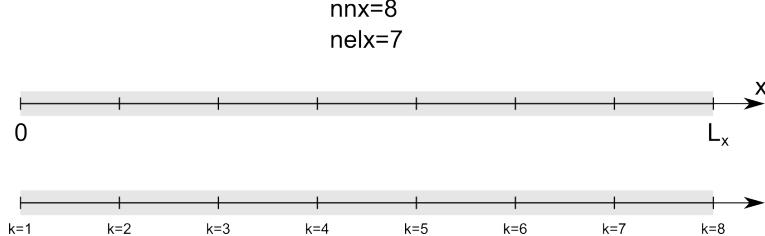
4.6.3 Triquadratic basis functions in 3D (Q_2)

$$\begin{aligned}
N_1 &= 0.5r(r-1) 0.5s(s-1) 0.5t(t-1) \\
N_2 &= 0.5r(r+1) 0.5s(s-1) 0.5t(t-1) \\
N_3 &= 0.5r(r+1) 0.5s(s+1) 0.5t(t-1) \\
N_4 &= 0.5r(r-1) 0.5s(s+1) 0.5t(t-1) \\
N_5 &= 0.5r(r-1) 0.5s(s-1) 0.5t(t+1) \\
N_6 &= 0.5r(r+1) 0.5s(s-1) 0.5t(t+1) \\
N_7 &= 0.5r(r+1) 0.5s(s+1) 0.5t(t+1) \\
N_8 &= 0.5r(r-1) 0.5s(s+1) 0.5t(t+1) \\
N_9 &= (1 - r^2) 0.5s(s-1) 0.5t(t-1) \\
N_{10} &= 0.5r(r+1) (1 - s^2) 0.5t(t-1) \\
N_{11} &= (1 - r^2) 0.5s(s+1) 0.5t(t-1) \\
N_{12} &= 0.5r(r-1) (1 - s^2) 0.5t(t-1) \\
N_{13} &= (1 - r^2) 0.5s(s-1) 0.5t(t+1) \\
N_{14} &= 0.5r(r+1) (1 - s^2) 0.5t(t+1) \\
N_{15} &= (1 - r^2) 0.5s(s+1) 0.5t(t+1) \\
N_{16} &= 0.5r(r-1) (1 - s^2) 0.5t(t+1) \\
N_{17} &= 0.5r(r-1) 0.5s(s-1) (1 - t^2) \\
N_{18} &= 0.5r(r+1) 0.5s(s-1) (1 - t^2) \\
N_{19} &= 0.5r(r+1) 0.5s(s+1) (1 - t^2) \\
N_{20} &= 0.5r(r-1) 0.5s(s+1) (1 - t^2) \\
N_{21} &= (1 - r^2) (1 - s^2) 0.5t(t-1) \\
N_{22} &= (1 - r^2) 0.5s(s-1) (1 - t^2) \\
N_{23} &= 0.5r(r+1) (1 - s^2) (1 - t^2) \\
N_{24} &= (1 - r^2) 0.5s(s+1) (1 - t^2) \\
N_{25} &= 0.5r(r-1) (1 - s^2) (1 - t^2) \\
N_{26} &= (1 - r^2) (1 - s^2) 0.5t(t+1) \\
N_{27} &= (1 - r^2) (1 - s^2) (1 - t^2)
\end{aligned}$$

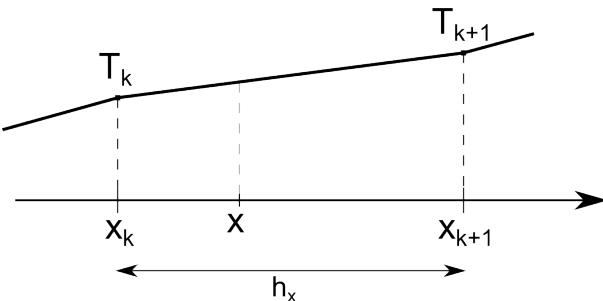
5 Solving the heat transport equation with linear Finite Elements

5.1 The diffusion equation in 1D

Let us consider the following one-dimensional grid:



Its spans the domain Ω of length L_x . It is discretised by means of nnx nodes and $nelx = nnx - 1$ elements. Zooming in on element 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} :



We focus here on the 1D diffusion equation (no advection, no heat sources):

$$\rho C_p \frac{\partial T}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) \quad (155)$$

This is the **strong form** of the ODE to solve. I can multiply this equation by a function⁵ $f(x)$ and integrate it over Ω :

$$\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx = \int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx \quad (156)$$

Looking at the right hand side, it is of the form $\int uv'$ so that I naturally integrate it by parts:

$$\int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx = \left[f(x) k \frac{\partial T}{\partial x} \right]_{\partial\Omega} - \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx \quad (157)$$

Assuming there is no heat flux prescribed on the boundary (i.e. $q_x = -k \partial T / \partial x = 0$),

NOT happy with this statement!!

then:

$$\int_{\Omega} f(x) \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) dx = - \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx \quad (158)$$

We then obtain the **weak form** of the diffusion equation in 1D:

$$\int_{\Omega} f(x) \rho C_p \frac{\partial T}{\partial t} dx + \int_{\Omega} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx = 0$$

(159)

⁵This function should be well-behaved with special properties, but we here assume it is a polynomial function.

We then use the additive property of the integral $\int_{\Omega} \dots = \sum_{elts} \int_{\Omega_e} \dots$ so that

$$\sum_{elts} \left(\underbrace{\int_{\Omega_e} f(x) \rho C_p \frac{\partial T}{\partial t} dx}_{\Lambda_f^e} + \underbrace{\int_{\Omega_e} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx}_{\Upsilon_f^e} \right) = 0 \quad (160)$$

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 have is the temperature at the nodes. For $x \in [x_k, x_{k+1}]$ we need to come up with a way to compute the temperature at this location. It makes sense to think that $T(x)$ will then be a function of the temperature at the nodes, i.e. $T(x) = \alpha T_k + \beta T_{k+1}$ where α and β are coefficients. One over-simplified approach would be to assign $T(x) = (T_k + T_{k+1})/2$ but this would make the temperature discontinuous from element to element. The rather logical solution to this problem is a linear temperature field between T_k and T_{k+1} :

$$T(x) = \underbrace{\frac{x_{k+1} - x}{h_x} T_k}_{N_k^\theta(x)} + \underbrace{\frac{x - x_k}{h_x} T_{k+1}}_{N_{k+1}^\theta(x)}$$

where $N_k^\theta(x)$ is the (temperature) shape function associated to node k and $N_{k+1}^\theta(x)$ is the shape function associated to node $k+1$.

Rather reassuringly, we have:

- $x = x_k$ yields $T(x) = T_k$
- $x = x_{k+1}$ yields $T(x) = T_{k+1}$
- $x = (x_k + x_{k+1})/2$ yields $T(x) = (T_k + T_{k+1})/2$

In what follows we abbreviate $\partial T / \partial x$ by \dot{T} . Let us compute Λ_f^e and Υ_f^e separately.

$$\begin{aligned} \Lambda_f^e &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p \dot{T}(x) dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p [N_k^\theta(x) \dot{T}_k + N_{k+1}^\theta(x) \dot{T}_{k+1}] dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p N_k^\theta(x) \dot{T}_k dx + \int_{x_k}^{x_{k+1}} f(x) \rho C_p N_{k+1}^\theta(x) \dot{T}_{k+1} dx \\ &= \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p N_k^\theta(x) dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p N_{k+1}^\theta(x) dx \right) \dot{T}_{k+1} \end{aligned}$$

Taking $f(x) = N_k^\theta(x)$ and omitting '(x)' in the rhs:

$$\Lambda_{N_k^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p N_k^\theta N_k^\theta dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p N_k^\theta N_{k+1}^\theta dx \right) \dot{T}_{k+1}$$

Taking $f(x) = N_{k+1}^\theta(x)$ and omitting '(x)' in the rhs:

$$\Lambda_{N_{k+1}^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p N_{k+1}^\theta N_k^\theta dx \right) \dot{T}_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p N_{k+1}^\theta N_{k+1}^\theta dx \right) \dot{T}_{k+1}$$

We can rearrange these last two equations as follows:

$$\begin{pmatrix} \Lambda_{N_k^\theta}^e \\ \Lambda_{N_{k+1}^\theta}^e \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} N_k^\theta \rho C_p N_k^\theta dx & \int_{x_k}^{x_{k+1}} N_k^\theta \rho C_p N_{k+1}^\theta dx \\ \int_{x_k}^{x_{k+1}} N_{k+1}^\theta \rho C_p N_k^\theta dx & \int_{x_k}^{x_{k+1}} N_{k+1}^\theta \rho C_p N_{k+1}^\theta dx \end{pmatrix} \cdot \begin{pmatrix} \dot{T}_k \\ \dot{T}_{k+1} \end{pmatrix}$$

and we can take the integrals outside of the matrix:

$$\begin{pmatrix} \Lambda_{N_k^\theta}^e \\ \Lambda_{N_{k+1}^\theta}^e \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} \rho C_p \begin{pmatrix} N_k^\theta N_k^\theta & N_k^\theta N_{k+1}^\theta \\ N_{k+1}^\theta N_k^\theta & N_{k+1}^\theta N_{k+1}^\theta \end{pmatrix} dx \right] \cdot \begin{pmatrix} \dot{T}_k \\ \dot{T}_{k+1} \end{pmatrix}$$

Finally, we can define the vectors

$$\vec{N}^T = \begin{pmatrix} N_k^\theta(x) \\ N_{k+1}^\theta(x) \end{pmatrix}$$

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

$$\begin{pmatrix} \Lambda_{N_k^\theta}^e \\ \Lambda_{N_{k+1}^\theta}^e \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \vec{N} dx \right) \cdot \dot{\vec{T}}^e$$

Back to the diffusion term:

$$\begin{aligned} \Upsilon_f^e &= \int_{x_k}^{x^{k+1}} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx \\ &= \int_{x_k}^{x^{k+1}} \frac{\partial f}{\partial x} k \frac{\partial (N_k^\theta(x) T_k + N_{k+1}^\theta(x) T_{k+1})}{\partial x} dx \\ &= \left(\int_{x_k}^{x^{k+1}} \frac{\partial f}{\partial x} k \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x^{k+1}} \frac{\partial f}{\partial x} k \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1} \end{aligned}$$

Taking $f(x) = N_k^\theta(x)$

$$\Upsilon_{N_k^\theta}^e = \left(\int_{x_k}^{x^{k+1}} k \frac{\partial N_k^\theta}{\partial x} \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x^{k+1}} k \frac{\partial N_k^\theta}{\partial x} \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1}$$

Taking $f(x) = N_{k+1}^\theta(x)$

$$\begin{aligned} \Upsilon_{N_{k+1}^\theta}^e &= \left(\int_{x_k}^{x^{k+1}} k \frac{\partial N_{k+1}^\theta}{\partial x} \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x^{k+1}} k \frac{\partial N_{k+1}^\theta}{\partial x} \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1} \\ \begin{pmatrix} \Upsilon_{N_k^\theta}^e \\ \Upsilon_{N_{k+1}^\theta}^e \end{pmatrix} &= \begin{pmatrix} \int_{x_k}^{x^{k+1}} \frac{\partial N_k^\theta}{\partial x} k \frac{\partial N_k^\theta}{\partial x} dx & \int_{x_k}^{x^{k+1}} \frac{\partial N_k^\theta}{\partial x} k \frac{\partial N_{k+1}^\theta}{\partial x} dx \\ \int_{x_k}^{x^{k+1}} \frac{\partial N_{k+1}^\theta}{\partial x} k \frac{\partial N_k^\theta}{\partial x} dx & \int_{x_k}^{x^{k+1}} \frac{\partial N_{k+1}^\theta}{\partial x} k \frac{\partial N_{k+1}^\theta}{\partial x} dx \end{pmatrix} \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix} \end{aligned}$$

or,

$$\begin{pmatrix} \Upsilon_{N_k^\theta}^e \\ \Upsilon_{N_{k+1}^\theta}^e \end{pmatrix} = \left[\int_{x_k}^{x^{k+1}} k \begin{pmatrix} \frac{\partial N_k^\theta}{\partial x} \frac{\partial N_k^\theta}{\partial x} & \frac{\partial N_k^\theta}{\partial x} \frac{\partial N_{k+1}^\theta}{\partial x} \\ \frac{\partial N_{k+1}^\theta}{\partial x} \frac{\partial N_k^\theta}{\partial x} & \frac{\partial N_{k+1}^\theta}{\partial x} \frac{\partial N_{k+1}^\theta}{\partial x} \end{pmatrix} dx \right] \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

Finally, we can define the vector

$$\vec{B}^T = \begin{pmatrix} \frac{\partial N_k^\theta}{\partial x} \\ \frac{\partial N_{k+1}^\theta}{\partial x} \end{pmatrix}$$

so that

$$\begin{pmatrix} \Upsilon_{N_k^\theta}^e \\ \Upsilon_{N_{k+1}^\theta}^e \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx \right) \cdot \vec{T}^e$$

The weak form discretised over 1 element becomes

$$\underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \vec{N} dx \right)}_{\mathbf{M}^e} \cdot \dot{\vec{T}}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx \right)}_{\mathbf{K}_d^e} \cdot \vec{T}^e = 0$$

or,

$$\boxed{\mathbf{M}^e \cdot \dot{\vec{T}}^e + \mathbf{K}_d^e \cdot \vec{T}^e = 0}$$

or,

$$\boxed{\mathbf{M}^e \cdot \frac{\partial \vec{T}^e}{\partial t} + \mathbf{K}_d^e \cdot \vec{T}^e = 0}$$

Using a backward first order in time discretisation for the time derivative:

$$\dot{\vec{T}} = \frac{\partial \vec{T}}{\partial t} = \frac{\vec{T}^{new} - \vec{T}^{old}}{\delta t}$$

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

$$\mathbf{M}^e = \int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \vec{N} dx \quad \mathbf{K}_d^e = \int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx$$

Let us compute \mathbf{M} for an element:

$$\mathbf{M}^e = \int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \vec{N} dx$$

with

$$\vec{N}^T = \begin{pmatrix} N_k(x) \\ N_{k+1}(x) \end{pmatrix} = \begin{pmatrix} \frac{x_{k+1}-x}{h_x} \\ \frac{x-x_k}{h_x} \end{pmatrix}$$

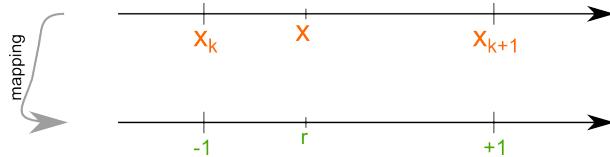
Then

$$\mathbf{M}^e = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \rho C_p N_k^\theta N_k^\theta dx & \int_{x_k}^{x_{k+1}} \rho C_p N_k^\theta N_{k+1}^\theta dx \\ \int_{x_k}^{x_{k+1}} \rho C_p N_{k+1}^\theta N_k^\theta dx & \int_{x_k}^{x_{k+1}} \rho C_p N_{k+1}^\theta N_{k+1}^\theta dx \end{pmatrix}$$

I only need to compute 3 integrals since $M_{12} = M_{21}$. Let us start with M_{11} :

$$M_{11} = \int_{x_k}^{x_{k+1}} \rho C_p N_k^\theta(x) N_k^\theta(x) dx = \int_{x_k}^{x_{k+1}} \rho C_p \frac{x_{k+1}-x}{h_x} \frac{x_{k+1}-x}{h_x} dx$$

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.

$$M_{11} = \rho C_p \int_{x_k}^{x_{k+1}} \frac{x_{k+1} - x}{h_x} \frac{x_{k+1} - x}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 - r)(1 - r) dr = \frac{h_x}{3} \rho C_p$$

Similarly we arrive at

$$M_{12} = \rho C_p \int_{x_k}^{x_{k+1}} \frac{x_{k+1} - x}{h_x} \frac{x - x_k}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 - r)(1 + r) dr = \frac{h_x}{6} \rho C_p$$

and

$$M_{22} = \rho C_p \int_{x_k}^{x_{k+1}} \frac{x - x_k}{h_x} \frac{x - x_k}{h_x} dx = \frac{\rho C_p h_x}{8} \int_{-1}^{+1} (1 + r)(1 + r) dr = \frac{h_x}{3} \rho C_p$$

Finally

$$\boxed{\mathbf{M}^e = \frac{h_x}{3} \rho C_p \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}}$$

In the new coordinate system, the shape functions

$$N_k^\theta(x) = \frac{x_{k+1} - x}{h_x} \quad N_{k+1}^\theta(x) = \frac{x - x_k}{h_x}$$

become

$$N_k^\theta(r) = \frac{1}{2}(1 - r) \quad N_{k+1}^\theta(r) = \frac{1}{2}(1 + r)$$

Also,

$$\frac{\partial N_k^\theta}{\partial x} = -\frac{1}{h_x} \quad \frac{\partial N_{k+1}^\theta}{\partial x} = \frac{1}{h_x}$$

so that

$$\vec{B}^T = \begin{pmatrix} \frac{\partial N_k^\theta}{\partial x} \\ \frac{\partial N_{k+1}^\theta}{\partial x} \end{pmatrix} = \begin{pmatrix} -\frac{1}{h_x} \\ \frac{1}{h_x} \end{pmatrix}$$

We here also assume that k is constant within the element:

$$\mathbf{K}_d = \int_{x_k}^{x_{k+1}} \vec{B}^T k \vec{B} dx = k \int_{x_k}^{x_{k+1}} \vec{B}^T \vec{B} dx$$

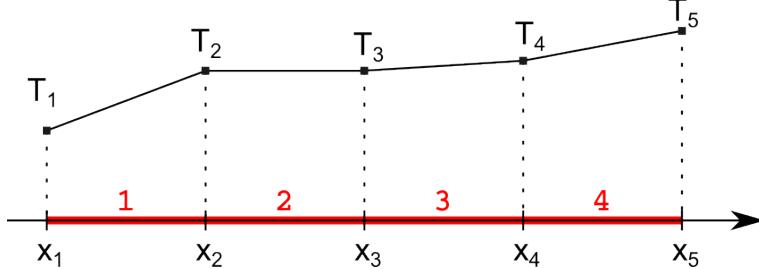
simply becomes

$$\mathbf{K}_d = k \int_{x_k}^{x_{k+1}} \frac{1}{h_x^2} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} dx$$

and then

$$\boxed{\mathbf{K}_d = \frac{k}{h_x} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}}$$

Let us consider this very simple grid consisting of 4 elements/5 nodes:



For each element we have

$$\underbrace{(\mathbf{M}^e + \mathbf{K}_d^e \delta t)}_{\mathbf{A}^e} \cdot \vec{\mathbf{T}}^{new} = \underbrace{\mathbf{M}^e \cdot \vec{\mathbf{T}}^{old}}_{\vec{\mathbf{b}}^e}$$

We can write this equation very explicitly for each element:

- element 1

$$\mathbf{A}^1 \cdot \begin{pmatrix} T_1 \\ T_2 \end{pmatrix} = \vec{\mathbf{b}}^1$$

$$\begin{cases} A_{11}^1 T_1 + A_{12}^1 T_2 = b_x^1 \\ A_{21}^1 T_1 + A_{22}^1 T_2 = b_y^1 \end{cases}$$

- element 2

$$\mathbf{A}^2 \cdot \begin{pmatrix} T_2 \\ T_3 \end{pmatrix} = \vec{\mathbf{b}}^2$$

$$\begin{cases} A_{11}^2 T_2 + A_{12}^2 T_3 = b_1^2 \\ A_{21}^2 T_2 + A_{22}^2 T_3 = b_2^2 \end{cases}$$

- element 3

$$\mathbf{A}^3 \cdot \begin{pmatrix} T_3 \\ T_4 \end{pmatrix} = \vec{\mathbf{b}}^3$$

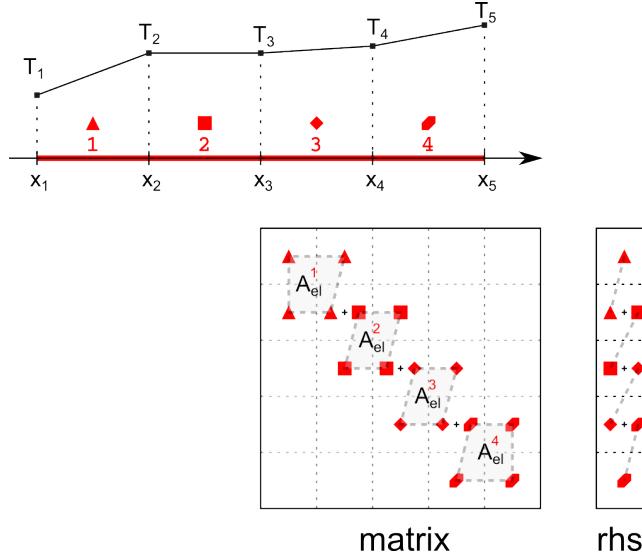
$$\begin{cases} A_{11}^3 T_3 + A_{12}^3 T_4 = b_1^3 \\ A_{21}^3 T_3 + A_{22}^3 T_4 = b_2^3 \end{cases}$$

- element 4

$$\mathbf{A}^4 \cdot \begin{pmatrix} T_4 \\ T_5 \end{pmatrix} = \vec{\mathbf{b}}^4$$

$$\begin{cases} A_{11}^4 T_4 + A_{12}^4 T_5 = b_1^4 \\ A_{21}^4 T_4 + A_{22}^4 T_5 = b_2^4 \end{cases}$$

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:

$$\begin{pmatrix} A_{11}^1 & A_{12}^1 & 0 & 0 & 0 \\ A_{21}^1 & A_{22}^1 + A_{11}^2 & A_{12}^2 & 0 & 0 \\ 0 & A_{21}^2 & A_{22}^2 + A_{11}^3 & A_{12}^3 & 0 \\ 0 & 0 & A_{21}^3 & A_{22}^3 + A_{11}^4 & A_{12}^4 \\ 0 & 0 & 0 & A_{21}^4 & A_{22}^4 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} b_1^1 \\ b_2^1 + b_1^2 \\ b_2^2 + b_1^3 \\ b_2^3 + b_1^4 \\ b_2^4 \end{pmatrix}$$

Ultimately the assembled matrix system also takes the form

$$\begin{pmatrix} A_{11} & A_{12} & 0 & 0 & 0 \\ A_{21} & A_{22} & A_{23} & 0 & 0 \\ 0 & A_{32} & A_{33} & A_{34} & 0 \\ 0 & 0 & A_{43} & A_{44} & A_{45} \\ 0 & 0 & 0 & A_{54} & A_{55} \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{pmatrix}$$

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

$$T_2 = T^{bc}$$

This equation can be cast as

$$(0 \ 1 \ 0 \ 0 \ 0) \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} 0 \\ T^{bc} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

This replaces the second line in the previous matrix equation:

$$\begin{pmatrix} A_{11} & A_{12} & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & A_{32} & A_{33} & A_{34} & 0 \\ 0 & 0 & A_{43} & A_{44} & A_{45} \\ 0 & 0 & 0 & A_{54} & A_{55} \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \end{pmatrix} = \begin{pmatrix} b_1 \\ T^{bc} \\ b_3 \\ b_4 \\ b_5 \end{pmatrix}$$

That's it, we have a linear system of equations which can be solved!

5.2 The advection-diffusion equation in 1D

We start with the 1D advection-diffusion equation

$$\rho C_p \left(\frac{\partial T}{\partial t} + u \frac{\partial T}{\partial x} \right) = \frac{\partial}{\partial x} \left(k \frac{\partial T}{\partial x} \right) + H \quad (161)$$

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 $\int_{\Omega} \cdots = \sum_{elts} \int_{\Omega_e} \cdots$

$$\sum_{elts} \left(\underbrace{\int_{\Omega_e} f(x) \rho C_p \frac{\partial T}{\partial t} dx}_{\Lambda_f^e} + \underbrace{\int_{\Omega_e} f(x) \rho C_p u \frac{\partial T}{\partial x} dx}_{\Sigma_f^e} + \underbrace{\int_{\Omega_e} \frac{\partial f}{\partial x} k \frac{\partial T}{\partial x} dx}_{\Upsilon_f^e} - \underbrace{\int_{\Omega_e} f(x) H dx}_{\Omega_f^e} \right) = 0$$

In the element, we have seen that the temperature can be written:

$$T(x) = N_k^\theta(x) T_k + N_{k+1}^\theta(x) T_{k+1}$$

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 .

$$\begin{aligned} \Sigma_f^e &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial T}{\partial x} dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial [N_k^\theta(x) T_k + N_{k+1}^\theta(x) T_{k+1}]}{\partial x} dx \\ &= \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial N_k^\theta}{\partial x} T_k dx + \int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial N_{k+1}^\theta}{\partial x} T_{k+1} dx \\ &= \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} f(x) \rho C_p u \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1} \end{aligned}$$

Taking $f(x) = N_k^\theta(x)$ and omitting '(x)' in the rhs:

$$\Sigma_{N_k^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p u N_k^\theta \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p u N_{k+1}^\theta \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1}$$

Taking $f(x) = N_{k+1}^\theta(x)$ and omitting '(x)' in the rhs:

$$\Sigma_{N_{k+1}^\theta}^e = \left(\int_{x_k}^{x_{k+1}} \rho C_p u N_{k+1}^\theta \frac{\partial N_k^\theta}{\partial x} dx \right) T_k + \left(\int_{x_k}^{x_{k+1}} \rho C_p u N_{k+1}^\theta \frac{\partial N_{k+1}^\theta}{\partial x} dx \right) T_{k+1}$$

$$\begin{pmatrix} \Sigma_{N_k^\theta} \\ \Sigma_{N_{k+1}^\theta} \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} \rho C_p u N_k^\theta \frac{\partial N_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \rho C_p u N_k^\theta \frac{\partial N_{k+1}^\theta}{\partial x} dx \\ \int_{x_k}^{x_{k+1}} \rho C_p u N_{k+1}^\theta \frac{\partial N_k^\theta}{\partial x} dx & \int_{x_k}^{x_{k+1}} \rho C_p u N_{k+1}^\theta \frac{\partial N_{k+1}^\theta}{\partial x} dx \end{pmatrix} \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

or,

$$\begin{pmatrix} \Sigma_{N_k^\theta} \\ \Sigma_{N_{k+1}^\theta} \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} \rho C_p u \begin{pmatrix} N_k^\theta \frac{\partial N_k^\theta}{\partial x} & N_k^\theta \frac{\partial N_{k+1}^\theta}{\partial x} \\ N_{k+1}^\theta \frac{\partial N_k^\theta}{\partial x} & N_{k+1}^\theta \frac{\partial N_{k+1}^\theta}{\partial x} \end{pmatrix} dx \right] \cdot \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

Finally, we have already defined the vectors

$$\vec{N}^T = \begin{pmatrix} N_k^\theta(x) \\ N_{k+1}^\theta(x) \end{pmatrix} \quad \vec{B}^T = \begin{pmatrix} \frac{\partial N_k^\theta}{\partial x} \\ \frac{\partial N_{k+1}^\theta}{\partial x} \end{pmatrix} \quad \vec{T}^e = \begin{pmatrix} T_k \\ T_{k+1} \end{pmatrix}$$

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 $f(x) = N_k^\theta(x)$:

$$\Omega_{N_k^\theta} = \int_{x_k}^{x_{k+1}} N_k^\theta(x) H(x) dx$$

Taking $f(x) = N_{k+1}^\theta(x)$:

$$\Omega_{N_{k+1}^\theta} = \int_{x_k}^{x_{k+1}} N_{k+1}^\theta(x) H(x) dx$$

We can rearrange both equations as follows:

$$\begin{pmatrix} \Omega_{N_k^\theta} \\ \Omega_{N_{k+1}^\theta} \end{pmatrix} = \begin{pmatrix} \int_{x_k}^{x_{k+1}} N_k^\theta(x) H(x) dx \\ \int_{x_k}^{x_{k+1}} N_{k+1}^\theta(x) H(x) dx \end{pmatrix}$$

or,

$$\begin{pmatrix} \Omega_{N_k^\theta} \\ \Omega_{N_{k+1}^\theta} \end{pmatrix} = \left[\int_{x_k}^{x_{k+1}} \begin{pmatrix} N_k^\theta(x) H(x) \\ N_{k+1}^\theta(x) H(x) \end{pmatrix} dx \right]$$

so that

$$\begin{pmatrix} \Omega_{N_k^\theta} \\ \Omega_{N_{k+1}^\theta} \end{pmatrix} = \left(\int_{x_k}^{x_{k+1}} \vec{N}^T H(x) dx \right)$$

The weak form discretised over 1 element becomes

$$\underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p \mathbf{N} dx \right)}_{\mathbf{M}^e} \cdot \dot{\vec{T}}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T \rho C_p u \mathbf{B} dx \right)}_{\mathbf{K}_a^e} \cdot \vec{T}^e + \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{B}^T k \mathbf{B} dx \right)}_{\mathbf{K}_d^e} \cdot \vec{T}^e = \underbrace{\left(\int_{x_k}^{x_{k+1}} \vec{N}^T H(x) dx \right)}_{\vec{F}^e}$$

or,

$$\mathbf{M}^e \cdot \dot{\vec{T}}^e + (\mathbf{K}_a^e + \mathbf{K}_d^e) \cdot \vec{T}^e = \vec{F}^e$$

or,

$$\mathbf{M}^e \cdot \frac{\partial \vec{T}^e}{\partial t} + (\mathbf{K}_a^e + \mathbf{K}_d^e) \cdot \vec{T}^e = \vec{F}^e$$

5.3 The advection-diffusion equation in 2D

We start from the 'bare-bones' heat transport equation (source terms are omitted):

$$\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot k \vec{\nabla} T \quad (162)$$

In what follows we assume that the velocity field \vec{v} is known so that temperature is the only unknown. Let N^θ be the temperature basis functions so that the temperature inside an element is given by⁶:

$$T^h(\vec{r}) = \sum_{i=1}^{m_T} N_i^\theta(\vec{r}) T_i = \vec{N}^\theta \cdot \vec{T} \quad (163)$$

where \vec{T} is a vector of length m_T . The weak form is then

$$\begin{aligned} \int_{\Omega} N_i^\theta \left[\rho C_p \left(\frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) \right] d\Omega &= \int_{\Omega} N_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T d\Omega \\ \underbrace{\int_{\Omega} N_i^\theta \rho C_p \frac{\partial T}{\partial t} d\Omega}_{I} + \underbrace{\int_{\Omega} N_i^\theta \rho C_p \vec{v} \cdot \vec{\nabla} T d\Omega}_{II} &= \underbrace{\int_{\Omega} N_i^\theta \vec{\nabla} \cdot k \vec{\nabla} T d\Omega}_{III} \quad i = 1, m_T \end{aligned} \quad (164)$$

Looking at the first term:

$$\int_{\Omega} N_i^\theta \rho C_p \frac{\partial T}{\partial t} d\Omega = \int_{\Omega} N_i^\theta \rho C_p \vec{N}^\theta \cdot \dot{\vec{T}} d\Omega \quad (165)$$

$$(166)$$

so that when we assemble all contributions for $i = 1, m_T$ we get:

$$I = \int_{\Omega} \vec{N}^\theta \rho C_p \vec{N}^\theta \cdot \dot{\vec{T}} d\Omega = \left(\int_{\Omega} \rho C_p \vec{N}^\theta \vec{N}^\theta d\Omega \right) \cdot \dot{\vec{T}} = \mathbf{M}^T \cdot \dot{\vec{T}}$$

where \mathbf{M}^T is the mass matrix of the system of size $(m_T \times m_T)$ with

$$M_{ij}^T = \int_{\Omega} \rho C_p N_i^\theta N_j^\theta d\Omega$$

Turning now to the second term:

$$\int_{\Omega} N_i^\theta \rho C_p \vec{v} \cdot \vec{\nabla} T d\Omega = \int_{\Omega} N_i^\theta \rho C_p (u \frac{\partial T}{\partial x} + v \frac{\partial T}{\partial y}) d\Omega \quad (167)$$

$$= \int_{\Omega} N_i^\theta \rho C_p (u \frac{\partial \vec{N}^\theta}{\partial x} + v \frac{\partial \vec{N}^\theta}{\partial y}) \cdot \vec{T} d\Omega \quad (168)$$

$$(169)$$

⁶the θ superscript has been chosen to denote temperature so as to avoid confusion with the transpose operator

so that when we assemble all contributions for $i = 1, m_T$ we get:

$$II = \left(\int_{\Omega} \rho C_p \vec{N}^{\theta} (u \frac{\partial \vec{N}^{\theta}}{\partial x} + v \frac{\partial \vec{N}^{\theta}}{\partial y}) d\Omega \right) \cdot \vec{T} = \mathbf{K}_a \cdot \vec{T}$$

where \mathbf{K}_a is the advection term matrix of size $(m_T \times m_T)$ with

$$(K_a)_{ij} = \int_{\Omega} \rho C_p N_i^{\theta} \left(u \frac{\partial N_j^{\theta}}{\partial x} + v \frac{\partial N_j^{\theta}}{\partial y} \right) d\Omega$$

Now looking at the third term, we carry out an integration by part and neglect the surface term for now, so that

$$\int_{\Omega} N_i^{\theta} \vec{\nabla} \cdot k \vec{\nabla} T d\Omega = - \int_{\Omega} k \vec{\nabla} N_i^{\theta} \cdot \vec{\nabla} T d\Omega \quad (170)$$

$$= - \int_{\Omega} k \vec{\nabla} N_i^{\theta} \cdot \vec{\nabla} (\vec{N}^{\theta} \cdot \vec{T}) d\Omega \quad (171)$$

(172)

with

$$\vec{\nabla} \vec{N}^{\theta} = \begin{pmatrix} \partial_x N_1^{\theta} & \partial_x N_2^{\theta} & \dots & \partial_x N_{m_T}^{\theta} \\ \partial_y N_1^{\theta} & \partial_y N_2^{\theta} & \dots & \partial_y N_{m_T}^{\theta} \end{pmatrix}$$

so that finally:

$$III = - \left(\int_{\Omega} k (\vec{\nabla} \vec{N}^{\theta})^T \cdot \vec{\nabla} \vec{N}^{\theta} d\Omega \right) \cdot \vec{T} = - \mathbf{K}_d \cdot \vec{T}$$

where \mathbf{K}_d is the diffusion term matrix:

$$\mathbf{K}_d = \int_{\Omega} k (\vec{\nabla} \vec{N}^{\theta})^T \cdot \vec{\nabla} \vec{N}^{\theta} d\Omega$$

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!!

5.3.1 Dealing with the time discretisation

Essentially we have to solve a PDE of the type:

$$\frac{\partial T}{\partial t} = \mathcal{F}(\vec{v}, T, \vec{\nabla} T, \Delta T)$$

with $\mathcal{F} = \frac{1}{\rho C_p} (-\vec{v} \cdot \vec{\nabla} T + \vec{\nabla} \cdot k \vec{\nabla} T)$.

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:

$$\frac{T^{n+1} - T^n}{\delta t} = \frac{1}{2} [\mathcal{F}^n(T, \vec{\nabla} T, \Delta T) + \mathcal{F}^{n+1}(T, \vec{\nabla} T, \Delta T)]$$

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,

$$\boxed{(\mathbf{M}^{n+1} + (\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1})\delta t) \cdot \vec{T}^{n+1} = \mathbf{M}^n \cdot \vec{T}^n}$$

- Crank-Nicolson

$$\frac{1}{\delta t}(\mathbf{M}^{n+1} \cdot \vec{T}^{n+1} - \mathbf{M}^n \cdot \vec{T}^n) = \frac{1}{2} [-(\mathbf{K}_a^{n+1} + \mathbf{K}_d^{n+1}) \cdot \vec{T}^{n+1} - (\mathbf{K}_a^n + \mathbf{K}_d^n) \cdot \vec{T}^n]$$

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 [288] or Wikipedia⁷. The second-order BDF (or BDF-2) as shown in [?] 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} :

$$\frac{\partial T}{\partial t}(t^n) = \frac{1}{\tau_n} \left(\frac{2\tau_n + \tau_{n-1}}{\tau_n + \tau_{n-1}} T(t^n) - \frac{\tau_n + \tau_{n-1}}{\tau_{n-1}} T(t^{n-1}) + \frac{\tau_n^2}{\tau_{n-1}(\tau_n + \tau_{n-1})} T(t^{n-2}) \right) \quad (173)$$

where $\tau_n = t^n - t^{n-1}$. Starting again from $\mathbf{M}^\theta \cdot \dot{\vec{T}} + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T} = \vec{0}$, we write

$$\mathbf{M}^\theta \cdot \frac{1}{\tau_n} \left(\frac{2\tau_n + \tau_{n-1}}{\tau_n + \tau_{n-1}} \vec{T}^n - \frac{\tau_n + \tau_{n-1}}{\tau_{n-1}} \vec{T}^{n-1} + \frac{\tau_n^2}{\tau_{n-1}(\tau_n + \tau_{n-1})} \vec{T}^{n-2} \right) + (\mathbf{K}_a + \mathbf{K}_d) \cdot \vec{T}^n = \vec{0}$$

and finally:

$$\left[\frac{2\tau_n + \tau_{n-1}}{\tau_n + \tau_{n-1}} \mathbf{M}^\theta + \tau_n (\mathbf{K}_a + \mathbf{K}_d) \right] \cdot \vec{T}^n = \frac{\tau_n + \tau_{n-1}}{\tau_{n-1}} \mathbf{M}^\theta \cdot \vec{T}^{n-1} - \frac{\tau_n^2}{\tau_{n-1}(\tau_n + \tau_{n-1})} \mathbf{M}^\theta \cdot \vec{T}^{n-2}$$

Note that if all timesteps are equal, i.e. $\tau_n = \tau_{n-1} = \delta t$, this equation 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,

$$\left[\mathbf{M}^\theta + \frac{2}{3} \delta t (\mathbf{K}_a + \mathbf{K}_d) \right] \cdot \vec{T}^n = \mathbf{M}^\theta \cdot \left(\frac{4}{3} \vec{T}^{n-1} - \frac{1}{3} \vec{T}^{n-2} \right)$$

As mentioned before the backward differentiation formula (BDF) is a family of implicit methods for the integration of ODEs. Each BDF- s method achieves order s . The BDF-1 is simply the backward Euler method as seen above:

$$T^{n+1} - T^n = \delta t \mathcal{F}^{n+1}$$

⁷https://en.wikipedia.org/wiki/Backward_differentiation_formula

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

$$T^{n+4} - \frac{48}{25}T^{n+1} + \frac{36}{25}T^{n+1} - \frac{16}{25}T^{n+1} + \frac{3}{25}T^n = \frac{12}{25}\delta t\mathcal{F}^{n+4}$$

6 Solving the flow equations with the FEM

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 word 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) [165]. 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 6.3) and the so-called mixed finite element method 6.4.

6.1 Strong and weak forms

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.

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.

6.2 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.

6.2.1 The compatibility condition (or LBB condition)

'LBB stable' elements assure the existence of a unique solution and assure convergence at the optimal rate.

6.2.2 Families

The family of Taylor-Hood finite element spaces on triangular/tetrahedral grids is given by $P_k \times P_{k-1}$ with $k \geq 2$, and on quadrilateral/hexahedral grids by $Q_k \times Q_{k-1}$ with $k \geq 2$. This means that the pressure is then approximated by continuous functions.

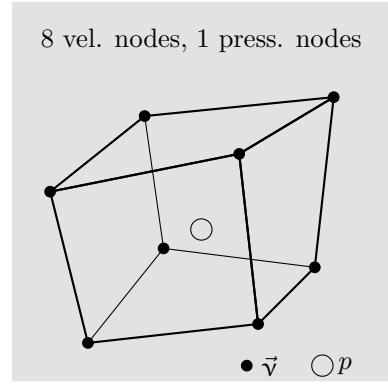
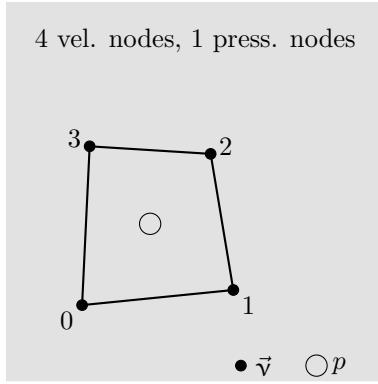
These finite elements are very popular, in particular the pairs for $k = 2$, i.e. $Q_2 \times Q_1$ and $P_2 \times P_1$. The reason why $k \geq 2$ comes from the fact that the $Q_1 \times Q_0$ (i.e. $Q_1 \times P_0$) and $P_2 \times P_1$ are not stable elements (they are not inf-sup stable).

Remark. Note that a similar element to $Q_2 \times Q_1$ has been proposed and used successfully used [536, 306]: it is denoted by $Q_2^{(8)} \times Q_1$ 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.

The Raviart-Thomas family on triangles and quadrilaterals.

find literature

6.2.3 The bi/tri-linear velocity - constant pressure element ($Q_1 \times P_0$)



discussed in example 3.71 of [330]

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 [308]p211. In [277] it is qualified as follows: slightly unstable but highly usable.

The $Q_1 \times P_0$ 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. [182, p235].

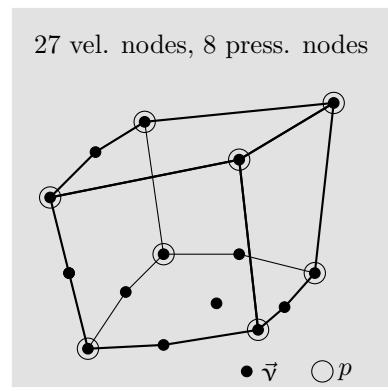
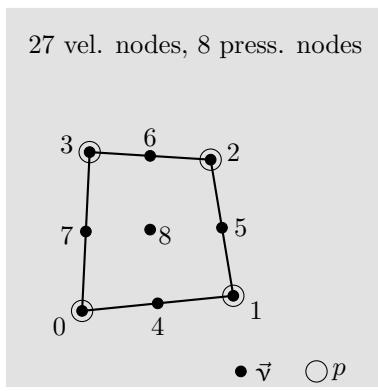
This element is discussed in [204], [205] and in [450] in the context of multigrid use.

This element is plagued by so-called pressure checkerboard modes which have been thoroughly analysed [279], [125], [491, 492]. These can be filtered out [125]. Smoothing techniques are also discussed in [365].

6.2.4 The bi/tri-quadratic velocity - discontinuous linear pressure element ($Q_2 \times P_{-1}$)

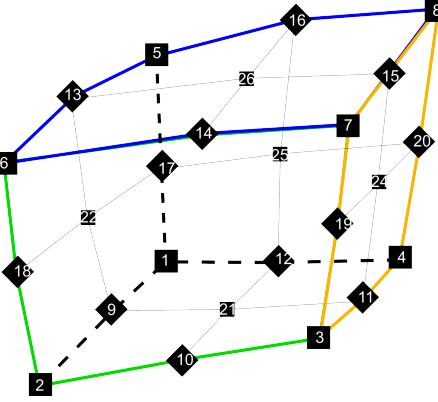
This element is crowned "probably the most accurate 2D element" in [277].

6.2.5 The bi/tri-quadratic velocity - bi/tri-linear pressure element ($Q_2 \times Q_1$)

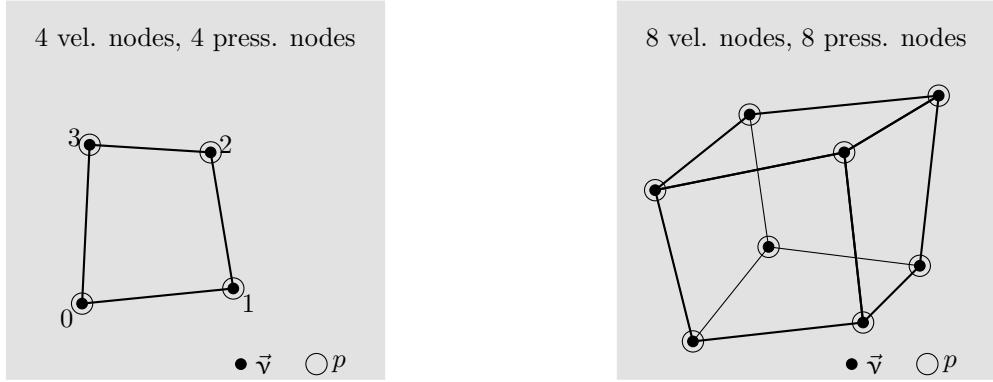


In [277] Gresho & Sani write that in their opinion $\text{div}(\vec{v}) = 0$ is not strong enough.

This element, implemented in penalised form, is discussed in [53] and the follow-up paper [54]. CHECK Biquadratic velocities, bilinear pressure. See Hood and Taylor. The element satisfies the inf-sup condition [308]p215.



6.2.6 The stabilised bi/tri-linear velocity - bi/tri-linear pressure element ($Q_1 \times Q_1$ -stab)



6.2.7 The MINI triangular element ($P_1^+ \times P_1$) in 2D

The MINI element was first introduced in Arnold et al, 1984 [22]. It is also discussed in section 3.6.1 of [330]. It is schematically represented hereunder:

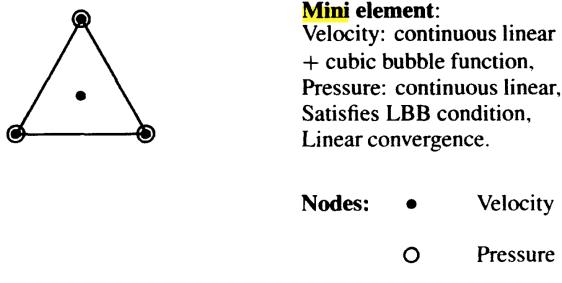
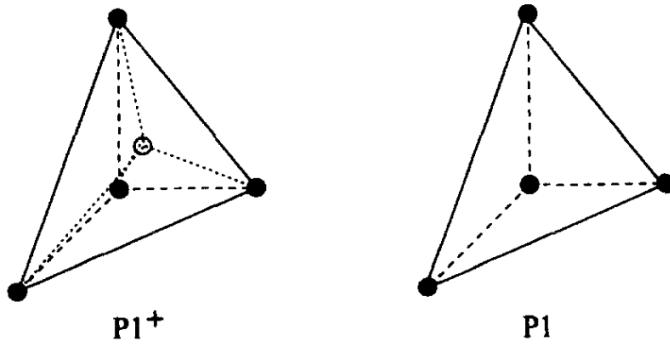


Figure taken from Donea and Huerta [165]

Remark. Note that [206] 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 [?].

Remark. According to Braess[76], since the support of the 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.

The 3D MINI element is not very common but it is used for instance in [447].



Velocity and pressure nodes for the 3D MINI element, taken from [447]

6.2.8 The quadratic velocity - linear pressure triangle ($P_2 \times P_1$)

From [501]. Taylor-Hood elements (Taylor and Hood 1973) are characterized by the fact that the pressure is continuous in the region Ω . A typical example is the quadratic triangle (P_2P_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.

6.2.9 The Crouzeix-Raviart triangle ($P_2^+ \times P_{-1}$)

This element was first introduced in [146]. It is the element used in the MILAMIN code [151]. 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 [150]. This element is LBB stable and no additional stabilization techniques are required[182]. 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 [150] recommend a 6-point or 7-point quadrature rule for this element.

Remark. Segal [501] 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 [150].

Remark. The simplest Crouzeix-Raviart element is the non-conforming linear triangle with constant pressure ($P_1 \times P_0$) [150].

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

6.2.10 The Rannacher-Turek element - rotated $Q_1 \times P_0$

p. 722 of [330]

6.2.11 Other elements

- P1P0 example 3.70 in [330]
- P1P1
- Q2P0: Quadratic velocities, constant pressure. The element satisfies the inf-sup condition, but the constant pressure assumption may require fine discretisation.
- Q2Q2: This element is never used, probably because a) it is unstable, b) it is very costly. There is one reference to it in [310].

- P2P2
- the MINI quadrilateral element $Q_1^+ \times Q_1$.

6.3 The penalty approach for viscous flow

In order to impose the incompressibility constraint, two widely used procedures are available, namely the Lagrange multiplier method and the penalty method [37, 308]. The latter is implemented in ELEFANT, which 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 [150], [476] or [284].

The penalty formulation of the mass conservation equation is based on a relaxation of the incompressibility constraint and writes

$$\vec{\nabla} \cdot \vec{v} + \frac{p}{\lambda} = 0 \quad (174)$$

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 [165, 311].

Equation (174) can be used to eliminate the pressure in the momentum equation so that the mass and momentum conservation equations fuse to become :

$$\vec{\nabla} \cdot (2\eta \dot{\epsilon}(\vec{v})) + \lambda \vec{\nabla}(\vec{\nabla} \cdot \vec{v}) = \rho \mathbf{g} = 0 \quad (175)$$

[393] 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 [51] 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 [262] .

list codes which use this approach

Since the penalty formulation is only valid for incompressible flows, then $\dot{\epsilon} = \dot{\epsilon}^d$ so that the d superscript is omitted in what follows. Because the stress tensor is symmetric one can also rewrite it the

following vector format:

$$\begin{aligned}
\begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} &= \begin{pmatrix} -p \\ -p \\ -p \\ 0 \\ 0 \\ 0 \end{pmatrix} + 2\eta \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yz} \end{pmatrix} \\
&= \lambda \begin{pmatrix} \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} \\ 0 \\ 0 \\ 0 \end{pmatrix} + 2\eta \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ \dot{\epsilon}_{xy} \\ \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yz} \end{pmatrix} \\
&= \left[\underbrace{\lambda \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}}_{K} + \underbrace{\eta \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_C \right] \cdot \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix}.
\end{aligned}$$

Remember that

$$\frac{\partial u}{\partial x} = \sum_{i=1}^4 \frac{\partial N_i}{\partial x} u_i \quad \frac{\partial v}{\partial y} = \sum_{i=1}^4 \frac{\partial N_i}{\partial y} v_i \quad \frac{\partial w}{\partial z} = \sum_{i=1}^4 \frac{\partial N_i}{\partial z} w_i$$

and

$$\begin{aligned}
\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} &= \sum_{i=1}^4 \frac{\partial N_i}{\partial y} u_i + \sum_{i=1}^4 \frac{\partial N_i}{\partial x} v_i \\
\frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} &= \sum_{i=1}^4 \frac{\partial N_i}{\partial z} u_i + \sum_{i=1}^4 \frac{\partial N_i}{\partial x} w_i \\
\frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} &= \sum_{i=1}^4 \frac{\partial N_i}{\partial z} v_i + \sum_{i=1}^4 \frac{\partial N_i}{\partial y} w_i
\end{aligned}$$

so that

$$\begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial N_1}{\partial x} & 0 & 0 & \frac{\partial N_2}{\partial x} & 0 & 0 & \frac{\partial N_3}{\partial x} & 0 & 0 & \dots & \frac{\partial N_4}{\partial x} & 0 & 0 \\ 0 & \frac{\partial N_1}{\partial y} & 0 & 0 & \frac{\partial N_2}{\partial y} & 0 & 0 & \frac{\partial N_3}{\partial y} & 0 & \dots & 0 & \frac{\partial N_4}{\partial y} & 0 \\ 0 & 0 & \frac{\partial N_1}{\partial z} & 0 & 0 & \frac{\partial N_2}{\partial z} & 0 & 0 & \frac{\partial N_3}{\partial z} & \dots & 0 & 0 & \frac{\partial N_4}{\partial z} \\ \frac{\partial N_1}{\partial y} & \frac{\partial N_1}{\partial x} & 0 & \frac{\partial N_2}{\partial y} & \frac{\partial N_2}{\partial x} & 0 & \frac{\partial N_3}{\partial y} & \frac{\partial N_3}{\partial x} & 0 & \dots & \frac{\partial N_4}{\partial y} & \frac{\partial N_4}{\partial x} & 0 \\ \frac{\partial N_1}{\partial z} & 0 & \frac{\partial N_1}{\partial x} & \frac{\partial N_2}{\partial z} & 0 & \frac{\partial N_2}{\partial x} & \frac{\partial N_3}{\partial z} & 0 & \frac{\partial N_3}{\partial x} & \dots & \frac{\partial N_4}{\partial z} & 0 & \frac{\partial N_4}{\partial x} \\ 0 & \frac{\partial N_1}{\partial z} & \frac{\partial N_1}{\partial y} & 0 & \frac{\partial N_2}{\partial z} & \frac{\partial N_2}{\partial y} & 0 & \frac{\partial N_3}{\partial z} & \frac{\partial N_3}{\partial y} & \dots & 0 & \frac{\partial N_4}{\partial z} & \frac{\partial N_4}{\partial y} \end{pmatrix}}_{B(6 \times 24)} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ w_1 \\ u_2 \\ v_2 \\ w_2 \\ u_3 \\ v_3 \\ w_3 \\ \dots \\ u_8 \\ v_8 \\ w_8 \end{pmatrix}}_{\vec{V}(24 \times 1)}$$

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 N_i 's 'regular enough', we can write:

$$\int_{\Omega_e} N_i \vec{\nabla} \cdot \boldsymbol{\sigma} d\Omega + \int_{\Omega_e} N_i \vec{b} d\Omega = 0$$

We can integrate by parts and drop the surface term⁸:

$$\int_{\Omega_e} \vec{\nabla} N_i \cdot \boldsymbol{\sigma} d\Omega = \int_{\Omega_e} N_i \vec{b} d\Omega$$

or,

$$\int_{\Omega_e} \begin{pmatrix} \frac{\partial N_i}{\partial x} & 0 & 0 & \frac{\partial N_i}{\partial y} & \frac{\partial N_i}{\partial z} & 0 \\ 0 & \frac{\partial N_i}{\partial y} & 0 & \frac{\partial N_i}{\partial x} & 0 & \frac{\partial N_i}{\partial z} \\ 0 & 0 & \frac{\partial N_i}{\partial z} & 0 & \frac{\partial N_i}{\partial x} & \frac{\partial N_i}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega = \int_{\Omega_e} N_i \vec{b} d\Omega$$

Let $i = 1, 2, 3, 4, \dots, 8$ and stack the resulting eight equations on top of one another.

$$\begin{aligned} \int_{\Omega_e} \begin{pmatrix} \frac{\partial N_1}{\partial x} & 0 & 0 & \frac{\partial N_1}{\partial y} & \frac{\partial N_1}{\partial z} & 0 \\ 0 & \frac{\partial N_1}{\partial y} & 0 & \frac{\partial N_1}{\partial x} & 0 & \frac{\partial N_1}{\partial z} \\ 0 & 0 & \frac{\partial N_1}{\partial z} & 0 & \frac{\partial N_1}{\partial x} & \frac{\partial N_1}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega &= \int_{\Omega_e} N_1 \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} d\Omega \\ \int_{\Omega_e} \begin{pmatrix} \frac{\partial N_2}{\partial x} & 0 & 0 & \frac{\partial N_2}{\partial y} & \frac{\partial N_2}{\partial z} & 0 \\ 0 & \frac{\partial N_2}{\partial y} & 0 & \frac{\partial N_2}{\partial x} & 0 & \frac{\partial N_2}{\partial z} \\ 0 & 0 & \frac{\partial N_2}{\partial z} & 0 & \frac{\partial N_2}{\partial x} & \frac{\partial N_2}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega &= \int_{\Omega_e} N_2 \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} d\Omega \\ &\dots \\ \int_{\Omega_e} \begin{pmatrix} \frac{\partial N_8}{\partial x} & 0 & 0 & \frac{\partial N_8}{\partial y} & \frac{\partial N_8}{\partial z} & 0 \\ 0 & \frac{\partial N_8}{\partial y} & 0 & \frac{\partial N_8}{\partial x} & 0 & \frac{\partial N_8}{\partial z} \\ 0 & 0 & \frac{\partial N_8}{\partial z} & 0 & \frac{\partial N_8}{\partial x} & \frac{\partial N_8}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega &= \int_{\Omega_e} N_8 \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} d\Omega \end{aligned} \quad (176)$$

We easily recognize \mathbf{B}^T inside the integrals! Let us define

$$\vec{N}_b^T = (N_1 b_x, N_1 b_y, N_1 b_z, \dots, N_8 b_x, N_8 b_y, N_8 b_z)$$

⁸We will come back to this at a later stage

then we can write

$$\int_{\Omega_e} \mathbf{B}^T \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega = \int_{\Omega_e} \vec{N}_b d\Omega$$

and finally:

$$\int_{\Omega_e} \mathbf{B}^T \cdot [\lambda \mathbf{K} + \eta \mathbf{C}] \cdot \mathbf{B} \cdot \vec{V} d\Omega = \int_{\Omega_e} \vec{N}_b d\Omega$$

Since \vec{V} contains is the vector of unknowns (i.e. the velocities at the corners), it does not depend on the x or y coordinates so it can be taking outside of the integral:

$$\underbrace{\left(\int_{\Omega_e} \mathbf{B}^T \cdot [\lambda \mathbf{K} + \eta \mathbf{C}] \cdot \mathbf{B} d\Omega \right)}_{\mathbf{A}_{el}(24 \times 24)} \cdot \underbrace{\vec{V}}_{(24 \times 1)} = \underbrace{\int_{\Omega_e} \vec{N}_b d\Omega}_{\vec{B}_{el}(24 \times 1)}$$

or,

$$\left[\underbrace{\left(\int_{\Omega_e} \lambda \mathbf{B}^T \cdot \mathbf{K} \cdot \mathbf{B} d\Omega \right)}_{\mathbf{A}_{el}^\lambda(24 \times 24)} + \underbrace{\left(\int_{\Omega_e} \eta \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} d\Omega \right)}_{\mathbf{A}_{el}^\eta(24 \times 24)} \right] \cdot \underbrace{\vec{V}}_{(24 \times 1)} = \underbrace{\int_{\Omega_e} \vec{N}_b d\Omega}_{\vec{B}_{el}(24 \times 1)}$$

reduced integration

reduced integration [311]

write about 3D to 2D

6.4 The mixed FEM for viscous flow

6.4.1 in three dimensions

In what follows the flow is assumed 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

$$\vec{v}^h(\vec{r}) = \sum_{i=1}^{m_v} N_i^v(\vec{r}) \vec{v}_i \quad (177)$$

where N_i^v are the polynomial basis functions for the velocity, and the summation runs over the m_v nodes composing the element. A similar expression is used for pressure:

$$p^h(\vec{r}) = \sum_{i=1}^{m_p} N_i^p(\vec{r}) p_i \quad (178)$$

Note that the velocity is a vector of size while pressure (and temperature) is a scalar. There are then $ndof_v$ velocity degrees of freedom per node and $ndof_p$ 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 co-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).

insert here link(s) to manual and literature

Other notations are sometimes used for Eqs.(177) and (178):

$$u^h(\vec{r}) = \vec{N}^v \cdot \vec{u} \quad v^h(\vec{r}) = \vec{N}^v \cdot \vec{v} \quad w^h(\vec{r}) = \vec{N}^v \cdot \vec{w} \quad p^h(\vec{r}) = \vec{N}^p \cdot \vec{p} \quad (179)$$

where $\vec{v} = (u, v, w)$ and \vec{N}^v is the vector containing all basis functions evaluated at location \vec{r} :

$$\vec{N}^v = (N_1^v(\vec{r}), N_2^v(\vec{r}), N_3^v(\vec{r}), \dots N_{m_v}^v(\vec{r})) \quad (180)$$

$$\vec{N}^p = (N_1^p(\vec{r}), N_2^p(\vec{r}), N_3^p(\vec{r}), \dots N_{m_p}^p(\vec{r})) \quad (181)$$

and with

$$\vec{u} = (u_1, u_2, u_3, \dots u_{m_v}) \quad (182)$$

$$\vec{v} = (v_1, v_2, v_3, \dots v_{m_v}) \quad (183)$$

$$\vec{w} = (w_1, w_2, w_3, \dots w_{m_v}) \quad (184)$$

$$\vec{p} = (p_1, p_2, p_3, \dots p_{m_p}) \quad (185)$$

We will now establish the weak form of the momentum conservation equation. We start again from

$$\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{b} = \vec{0} \quad (186)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (187)$$

For the N_i^v 's and N_i^p 'regular enough', we can write:

$$\int_{\Omega_e} N_i^v \vec{\nabla} \cdot \boldsymbol{\sigma} d\Omega + \int_{\Omega_e} N_i^v \vec{b} d\Omega = \vec{0} \quad (188)$$

$$\int_{\Omega_e} N_i^p \vec{\nabla} \cdot \vec{v} d\Omega = 0 \quad (189)$$

We can integrate by parts and drop the surface term⁹:

$$\int_{\Omega_e} \vec{\nabla} N_i^v \cdot \boldsymbol{\sigma} d\Omega = \int_{\Omega_e} N_i^v \vec{b} d\Omega \quad (190)$$

$$\int_{\Omega_e} N_i^p \vec{\nabla} \cdot \vec{v} d\Omega = 0 \quad (191)$$

or,

$$\int_{\Omega_e} \begin{pmatrix} \frac{\partial N_i^v}{\partial x} & 0 & 0 & \frac{\partial N_i^v}{\partial y} & \frac{\partial N_i^v}{\partial z} & 0 \\ 0 & \frac{\partial N_i^v}{\partial y} & 0 & \frac{\partial N_i^v}{\partial x} & 0 & \frac{\partial N_i^v}{\partial z} \\ 0 & 0 & \frac{\partial N_i^v}{\partial z} & 0 & \frac{\partial N_i^v}{\partial x} & \frac{\partial N_i^v}{\partial y} \end{pmatrix} \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega = \int_{\Omega_e} N_i^v \vec{b} d\Omega \quad (192)$$

As before (see section XXX) the above equation can ultimately be written:

$$\int_{\Omega_e} \boldsymbol{B}^T \cdot \begin{pmatrix} \sigma_{xx} \\ \sigma_{yy} \\ \sigma_{zz} \\ \sigma_{xy} \\ \sigma_{xz} \\ \sigma_{yz} \end{pmatrix} d\Omega = \int_{\Omega_e} \vec{N}_b d\Omega \quad (193)$$

⁹We will come back to this at a later stage

We have previously established that the strain rate vector $\vec{\dot{\varepsilon}}$ is:

$$\vec{\dot{\varepsilon}} = \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial w}{\partial z} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \\ \frac{\partial u}{\partial z} + \frac{\partial w}{\partial x} \\ \frac{\partial v}{\partial z} + \frac{\partial w}{\partial y} \end{pmatrix} = \begin{pmatrix} \sum_i \frac{\partial N_i^Y}{\partial x} u_i \\ \sum_i \frac{\partial N_i^Y}{\partial y} v_i \\ \sum_i \frac{\partial N_i^Y}{\partial z} w_i \\ \sum_i (\frac{\partial N_i^Y}{\partial y} u_i + \frac{\partial N_i^Y}{\partial x} v_i) \\ \sum_i (\frac{\partial N_i^Y}{\partial z} u_i + \frac{\partial N_i^Y}{\partial x} w_i) \\ \sum_i (\frac{\partial N_i^Y}{\partial z} v_i + \frac{\partial N_i^Y}{\partial y} w_i) \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial N_1^Y}{\partial x} & 0 & 0 & \dots & \frac{\partial N_{m_v}^Y}{\partial x} & 0 & 0 \\ 0 & \frac{\partial N_1^Y}{\partial y} & 0 & \dots & 0 & \frac{\partial N_{m_v}^Y}{\partial y} & 0 \\ 0 & 0 & \frac{\partial N_1^Y}{\partial z} & \dots & 0 & 0 & \frac{\partial N_{m_v}^Y}{\partial z} \\ \frac{\partial N_1^Y}{\partial y} & \frac{\partial N_1^Y}{\partial x} & 0 & \dots & \frac{\partial N_{m_v}^Y}{\partial x} & \frac{\partial N_{m_v}^Y}{\partial z} & 0 \\ \frac{\partial N_1^Y}{\partial z} & 0 & \frac{\partial N_1^Y}{\partial x} & \dots & \frac{\partial N_{m_v}^Y}{\partial z} & 0 & \frac{\partial N_{m_v}^Y}{\partial x} \\ 0 & \frac{\partial N_1^Y}{\partial z} & \frac{\partial N_1^Y}{\partial y} & \dots & 0 & \frac{\partial N_{m_v}^Y}{\partial z} & \frac{\partial N_{m_v}^Y}{\partial y} \end{pmatrix}}_{\mathbf{B}} \cdot \underbrace{\begin{pmatrix} u_1 \\ v_1 \\ w_1 \\ u_2 \\ v_2 \\ w_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \\ w_{m_v} \end{pmatrix}}_{\vec{V}} \quad (194)$$

or, $\vec{\dot{\varepsilon}} = \mathbf{B} \cdot \vec{V}$ where \mathbf{B} is the gradient matrix and \vec{V} is the vector of all vector degrees of freedom for the element. The matrix \mathbf{B} is then of size $3 \times m_v \times ndim$ and the vector \vec{V} is $m_v * ndof$ long. we have

$$\sigma_{xx} = -p + 2\eta\dot{\varepsilon}_{xx}^d \quad (195)$$

$$\sigma_{yy} = -p + 2\eta\dot{\varepsilon}_{yy}^d \quad (196)$$

$$\sigma_{zz} = -p + 2\eta\dot{\varepsilon}_{zz}^d \quad (197)$$

$$\sigma_{xy} = 2\eta\dot{\varepsilon}_{xy}^d \quad (198)$$

$$\sigma_{xz} = 2\eta\dot{\varepsilon}_{xz}^d \quad (199)$$

$$\sigma_{yz} = 2\eta\dot{\varepsilon}_{yz}^d \quad (200)$$

Since we here only consider incompressible flow, we have $\dot{\varepsilon}^d = \dot{\varepsilon}$ so

$$\vec{\sigma} = - \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} p + \mathbf{C} \cdot \vec{\dot{\varepsilon}} = - \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \vec{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (201)$$

with

$$\mathbf{C} = \eta \begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \vec{\dot{\varepsilon}} = \begin{pmatrix} \dot{\varepsilon}_{xx} \\ \dot{\varepsilon}_{yy} \\ \dot{\varepsilon}_{zz} \\ 2\dot{\varepsilon}_{xy} \\ 2\dot{\varepsilon}_{xz} \\ 2\dot{\varepsilon}_{yz} \end{pmatrix} \quad (202)$$

Let us define matrix \mathbf{N}^p of size $6 \times m_p$:

$$\mathbf{N}^p = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \vec{N}^p = \begin{pmatrix} \vec{N}^p \\ \vec{N}^p \\ \vec{N}^p \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (203)$$

so that

$$\vec{\sigma} = -\mathbf{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (204)$$

finally

$$\int_{\Omega_e} \mathbf{B}^T \cdot [-\mathbf{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V}] d\Omega = \int_{\Omega_e} \mathbf{N}_b d\Omega \quad (205)$$

or,

$$\underbrace{\left(- \int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{N}^p d\Omega \right)}_{\mathbb{G}} \cdot \vec{P} + \underbrace{\left(\int_{\Omega_e} \mathbf{B}^T \cdot \mathbf{C} \cdot \mathbf{B} d\Omega \right)}_{\mathbb{K}} \cdot \vec{V} = \underbrace{\int_{\Omega_e} \mathbf{N}_b d\Omega}_{\vec{f}} \quad (206)$$

where the matrix \mathbb{K} is of size $(m_v * ndof_v \times m_v * ndof_v)$, and matrix \mathbb{G} is of size $(m_v * ndof_v \times m_p * ndof_p)$. Turning now to the mass conservation equation:

$$\begin{aligned} \vec{0} &= \int_{\Omega_e} \vec{N}^p \vec{\nabla} \cdot \vec{v} d\Omega \\ &= \int_{\Omega_e} \vec{N}^p \sum_{i=1}^{m_v} \left(\frac{\partial N_i^v}{\partial x} u_i + \frac{\partial N_i^v}{\partial y} v_i + \frac{\partial N_i^v}{\partial z} w_i \right) d\Omega \\ &= \int_{\Omega_e} \begin{pmatrix} N_1^p \left(\sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial z} w_i \right) \\ N_2^p \left(\sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial z} w_i \right) \\ N_3^p \left(\sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial z} w_i \right) \\ \vdots \\ N_{m_p}^p \left(\sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial x} u_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial y} v_i + \sum_{i=1}^{m_v} \frac{\partial N_i^v}{\partial z} w_i \right) \end{pmatrix} d\Omega \\ &= \int_{\Omega_e} \begin{pmatrix} N_1^p & N_1^p & N_1^p & 0 & 0 & 0 \\ N_2^p & N_2^p & N_2^p & 0 & 0 & 0 \\ N_3^p & N_3^p & N_3^p & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ N_{m_p}^p & N_{m_p}^p & N_{m_p}^p & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} \sum_i \frac{\partial N_i^v}{\partial x} u_i \\ \sum_i \frac{\partial N_i^v}{\partial y} v_i \\ \sum_i \frac{\partial N_i^v}{\partial z} w_i \\ \sum_i (\frac{\partial N_i^v}{\partial y} u_i + \frac{\partial N_i^v}{\partial x} v_i) \\ \sum_i (\frac{\partial N_i^v}{\partial z} u_i + \frac{\partial N_i^v}{\partial x} w_i) \\ \sum_i (\frac{\partial N_i^v}{\partial z} v_i + \frac{\partial N_i^v}{\partial y} w_i) \end{pmatrix} d\Omega \\ &= \int_{\Omega_e} \underbrace{\begin{pmatrix} N_1^p & N_1^p & N_1^p & 0 & 0 & 0 \\ N_2^p & N_2^p & N_2^p & 0 & 0 & 0 \\ N_3^p & N_3^p & N_3^p & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ N_{m_p}^p & N_{m_p}^p & N_{m_p}^p & 0 & 0 & 0 \end{pmatrix}}_{\mathbf{N}^p} \cdot \vec{\varepsilon} d\Omega \\ &= \left(\int \mathbf{N}^p \cdot \mathbf{B} d\Omega \right) \cdot \vec{V} \\ &= -\mathbb{G}_e^T \cdot \vec{V} \end{aligned} \quad (207)$$

Note that it is common to actually start from $-\vec{\nabla} \cdot \vec{v} = 0$ (see Eq.(3) in [3]) 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:

$$\begin{pmatrix} \mathbb{K}_e & \mathbb{G}_e \\ \mathbb{G}_e^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{V} \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f}_e \\ 0 \end{pmatrix}$$

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:

$$-\vec{\nabla}p + \vec{\nabla} \cdot (2\eta\dot{\varepsilon}) + \rho\mathbf{g} = 0 \quad (208)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad (209)$$

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 obtained by multiplying the strong form equations by the (dimensionless) basis functions and integrating over the domain, so that it follows that

$$[\mathbb{K} \cdot \vec{V}] = [\mathbb{G} \cdot \vec{P}] = [\vec{f}] = ML^{-2}T^{-2}L^3 = MLT^{-2}$$

We can then easily deduce:

$$[\mathbb{K}] = MT^{-1} \quad [\mathbb{G}] = L^2$$

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 [277], 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.

It is then worth noting that *only* discontinuous pressure elements assure an element-level mass balance [277]: 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 C matrix The relationship between deviatoric stress and deviatoric strain rate tensor is

$$\tau = 2\eta \dot{\epsilon}^d \quad (210)$$

$$= 2\eta \left(\dot{\epsilon} - \frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} \right) \quad (211)$$

$$= 2\eta \left[\begin{pmatrix} \dot{\epsilon}_{xx} & \dot{\epsilon}_{xy} & \dot{\epsilon}_{xz} \\ \dot{\epsilon}_{yx} & \dot{\epsilon}_{yy} & \dot{\epsilon}_{yz} \\ \dot{\epsilon}_{zx} & \dot{\epsilon}_{zy} & \dot{\epsilon}_{zz} \end{pmatrix} - \frac{1}{3}(\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz}) \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right] \quad (212)$$

$$= \frac{2}{3}\eta \begin{pmatrix} 2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} - \dot{\epsilon}_{zz} & 3\dot{\epsilon}_{xy} & 3\dot{\epsilon}_{xz} \\ 3\dot{\epsilon}_{yx} & -\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy} - \dot{\epsilon}_{yy} & 3\dot{\epsilon}_{yz} \\ 3\dot{\epsilon}_{zx} & 3\dot{\epsilon}_{zy} & -\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} - 2\dot{\epsilon}_{zz} \end{pmatrix} \quad (213)$$

so that

$$\vec{\tau} = \frac{2}{3}\eta \begin{pmatrix} 2\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} - \dot{\epsilon}_{zz} \\ -\dot{\epsilon}_{yy} + 2\dot{\epsilon}_{yy} - \dot{\epsilon}_{yy} \\ -\dot{\epsilon}_{xx} - \dot{\epsilon}_{yy} + 2\dot{\epsilon}_{zz} \\ 3\dot{\epsilon}_{xy} \\ 3\dot{\epsilon}_{xz} \\ 3\dot{\epsilon}_{yz} \end{pmatrix} = \underbrace{\frac{\eta}{3} \begin{pmatrix} 4 & -2 & -2 & 0 & 0 & 0 \\ -2 & 4 & -2 & 0 & 0 & 0 \\ -2 & -2 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{pmatrix}}_{C^d} \cdot \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} = \mathbf{C}^d \cdot \vec{\epsilon} \quad (214)$$

In two dimensions, we have

$$\vec{\tau} = \frac{1}{3}\eta \underbrace{\begin{pmatrix} 4 & -2 & 0 \\ -2 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix}}_{C^d} \cdot$$

In the case where we assume incompressible flow from the beginning, i.e. $\dot{\epsilon} = \dot{\epsilon}^d$, then

$$\vec{\tau} = \eta \underbrace{\begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_C \cdot \begin{pmatrix} \dot{\epsilon}_{xx} \\ \dot{\epsilon}_{yy} \\ \dot{\epsilon}_{zz} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} = \mathbf{C} \cdot \vec{\epsilon} \quad (215)$$

Two slightly different formulations The momentum conservation equation can be written as follows:

$$\vec{\nabla} \cdot (2\eta \vec{\epsilon}) - \vec{\nabla} p + \vec{b} = \vec{0}$$

When the viscosity η is constant 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 [165, Eq. 6.24] and this can have consequences for the Neumann boundary conditions.

On the 'forgotten' surface terms

6.4.2 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:

$$\dot{\boldsymbol{\varepsilon}}^d = \dot{\boldsymbol{\varepsilon}} - \frac{1}{3}(\vec{\nabla} \cdot \vec{v})\mathbf{1} \quad (216)$$

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:

$$\boldsymbol{\varepsilon} = \begin{pmatrix} \dot{\varepsilon}_{xx} & \dot{\varepsilon}_{xy} & \dot{\varepsilon}_{xz} \\ \dot{\varepsilon}_{yx} & \dot{\varepsilon}_{yy} & \dot{\varepsilon}_{yz} \\ \dot{\varepsilon}_{zx} & \dot{\varepsilon}_{zy} & \dot{\varepsilon}_{zz} \end{pmatrix} = \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{1}{2}\left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}\right) & 0 \\ \frac{1}{2}\left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}\right) & \frac{\partial v}{\partial y} & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (217)$$

$$\dot{\boldsymbol{\varepsilon}}^d = \frac{1}{3} \begin{pmatrix} 2\frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} & \frac{1}{2}\left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}\right) & 0 \\ \frac{1}{2}\left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x}\right) & -\frac{\partial u}{\partial x} + 2\frac{\partial v}{\partial y} & 0 \\ 0 & 0 & -\frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} \end{pmatrix} \quad (218)$$

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:

$$\vec{\nabla} \cdot 2\eta\dot{\boldsymbol{\varepsilon}}^d =$$

FINISH!

In two dimensions the velocity is then $\vec{v} = (u, v)$ and the FEM building blocks and matrices are simply:

$$\vec{\dot{\boldsymbol{\varepsilon}}} = \begin{pmatrix} \dot{\varepsilon}_{xx} \\ \dot{\varepsilon}_{yy} \\ 2\dot{\varepsilon}_{xy} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \end{pmatrix}}_B = \underbrace{\begin{pmatrix} \frac{\partial N_1^\gamma}{\partial x} & 0 & \frac{\partial N_2^\gamma}{\partial x} & 0 & \frac{\partial N_3^\gamma}{\partial x} & 0 & \dots & \frac{\partial N_{m_v}^\gamma}{\partial x} & 0 \\ 0 & \frac{\partial N_1^\gamma}{\partial y} & 0 & \frac{\partial N_2^\gamma}{\partial y} & 0 & \frac{\partial N_3^\gamma}{\partial y} & \dots & 0 & \frac{\partial N_{m_v}^\gamma}{\partial y} \\ \frac{\partial N_1^\gamma}{\partial y} & \frac{\partial N_1^\gamma}{\partial x} & \frac{\partial N_2^\gamma}{\partial y} & \frac{\partial N_2^\gamma}{\partial x} & \frac{\partial N_3^\gamma}{\partial y} & \frac{\partial N_3^\gamma}{\partial x} & \dots & \frac{\partial N_{m_v}^\gamma}{\partial y} & \frac{\partial N_{m_v}^\gamma}{\partial x} \end{pmatrix}}_B \cdot \begin{pmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ u_3 \\ v_3 \\ \dots \\ u_{m_v} \\ v_{m_v} \end{pmatrix} \underbrace{\vec{v}}$$
(219)

we have

$$\sigma_{xx} = -p + 2\eta\dot{\varepsilon}_{xx} \quad (220)$$

$$\sigma_{yy} = -p + 2\eta\dot{\varepsilon}_{yy} \quad (221)$$

$$\sigma_{xy} = +2\eta\dot{\varepsilon}_{xy} \quad (222)$$

so

$$\vec{\sigma} = - \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} p + \mathbf{C} \cdot \vec{\dot{\boldsymbol{\varepsilon}}} = - \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \vec{N}^p \cdot \vec{P} + \mathbf{C} \cdot \mathbf{B} \cdot \vec{V} \quad (223)$$

with

$$\mathbf{C} = \eta \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{or} \quad \mathbf{C} = \frac{\eta}{3} \begin{pmatrix} 4 & -2 & 0 \\ -2 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad (224)$$

check the right \mathbf{C}

Finally the matrix \mathbf{N}^p is of size $3 \times m_p$:

$$\mathbf{N}^p = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \vec{\mathbf{N}}^p = \begin{pmatrix} \vec{\mathbf{N}}^p \\ \vec{\mathbf{N}}^p \\ 0 \end{pmatrix} \quad (225)$$

6.5 Solving the elastic equations

6.6 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

6.7 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

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{P} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ 0 \end{pmatrix}$$

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 $\mathbb{S}^* = \int \eta^{-1} \vec{\mathbf{N}}^T \vec{\mathbf{N}} dv$ (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 Additional techniques and features

7.1 Dealing with a free surface

7.2 Convergence criterion for nonlinear iterations

7.3 The SUPG formulation for the energy equation

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.

[84]

7.3.1 Linear elements

When using linear elements, its implementation is rather trivial, as shown in the DOUAR paper [80] or the FANTOM paper [542]. The advection matrix is simply modified and computed as follows:

$$(\mathbf{K}_a^e)_{SUPG} = \int_{x_k}^{x_{k+1}} (\mathbf{N}^*)^T \rho C_p \vec{v} \cdot \mathbf{B} dx \quad \text{with} \quad \mathbf{N}^* = \mathbf{N} + \tau \vec{v} \cdot \mathbf{B}$$

Note that we can also write

$$(\mathbf{K}_a^e)_{SUPG} = \int_{x_k}^{x_{k+1}} \mathbf{N}^T \rho C_p \vec{v} \cdot \mathbf{B} dx + \int_{x_k}^{x_{k+1}} \tau (\vec{v} \cdot \mathbf{B})^T \rho C_p (\vec{v} \cdot \mathbf{B}) dx$$

and we see that the SUPG method introduces and 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 :

$$(\mathbf{K}_a^e)_{SUPG} = \mathbf{K}_a^e + \rho C_p \frac{\tau u^2}{h} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

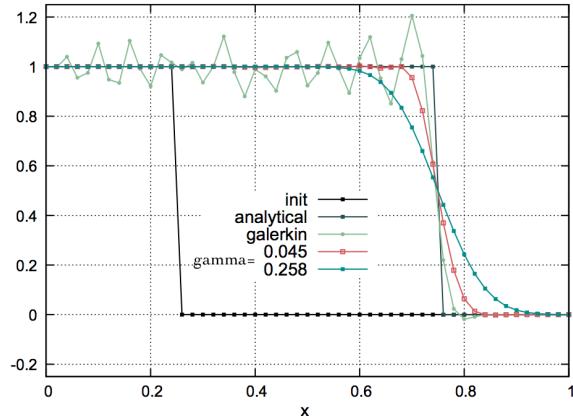
The additional matrix has the same structure as the 1D diffusion matrix matrix in 5.1.

The parameter τ is chosen as follows:

$$\tau = \gamma \frac{h}{\nu} \quad (226)$$

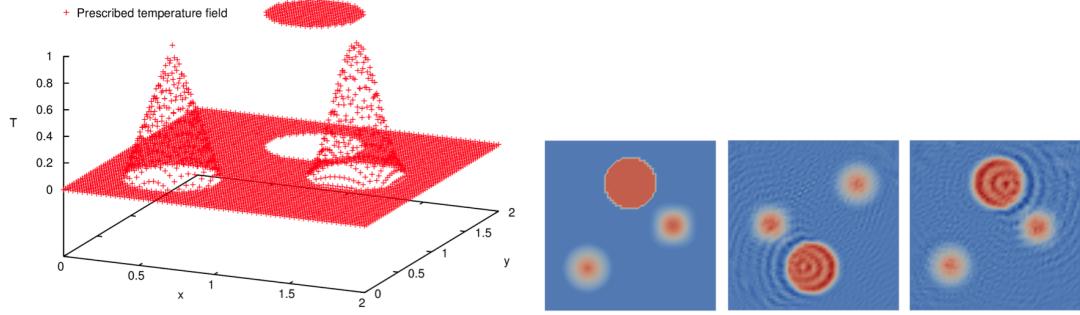
where γ is a user chosen parameter (see Appendix A of [542]).

A typical test case for testing a advection scheme is the step advection benchmark (see for instance [165]). 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:



The prescribed velocity is $\nu = 1$, 50 elements are used and 250 time steps are carried out with $\delta t = 0.1h/\nu = 0.002$. As discussed in [542], using Equation 226, one arrives to $\gamma = 0.045$, which leads to a desired removal of the oscillations through a small amount of numerical diffusion. Braun [79] argues for a constant $\gamma = 1/\sqrt{15} = 0.258$ (after [309]), which effect is also shown in the figure above. This value is arguably too large and introduces undesirable diffusion.

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 [165]).



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[358]. It instead relied on the Entropy Viscosity formulation [282]. 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 63¹⁰.

¹⁰https://www.dealii.org/developer/doxygen/deal.II/step_63.html

7.4 The method of manufactured solutions

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 [165, 113, 65].

7.4.1 Analytical benchmark I - "Donea & Huerta"

Taken from [165]. 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 velocity field $\vec{v} = (u, v)$ and the pressure p such that

$$\eta \Delta \vec{v} - \vec{\nabla} p + \vec{b} = \vec{0} \quad \text{in } \Omega \quad (227)$$

$$\vec{\nabla} \cdot \vec{v} = 0 \quad \text{in } \Omega \quad (228)$$

$$\vec{v} = \vec{0} \quad \text{on } \Gamma_D \quad (229)$$

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

$$\begin{aligned} u(x, y) &= x^2(1-x)^2(2y-6y^2+4y^3) \\ v(x, y) &= -y^2(1-y)^2(2x-6x^2+4x^3) \\ p(x, y) &= x(1-x) - 1/6 \end{aligned}$$

Note that the pressure obeys $\int_{\Omega} p d\Omega = 0$. One can turn to the spatial derivatives of the fields:

$$\frac{\partial u}{\partial x} = (2x - 6x^2 + 4x^3)(2y - 6y^2 + 4y^3) \quad (230)$$

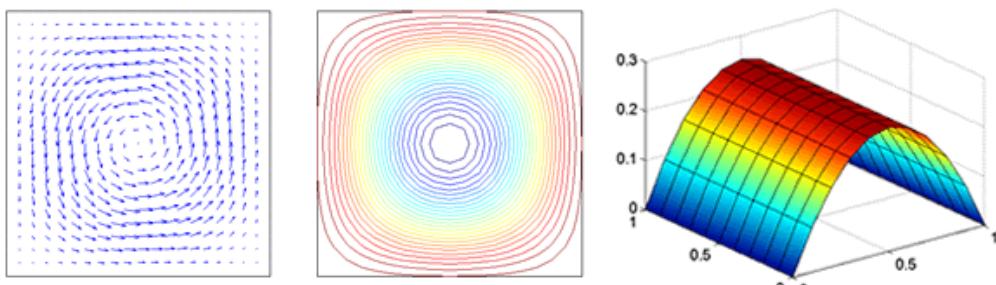
$$\frac{\partial v}{\partial y} = -(2x - 6x^2 + 4x^3)(2y - 6y^2 + 4y^3) \quad (231)$$

with of course $\vec{\nabla} \cdot \vec{v} = 0$ and

$$\frac{\partial p}{\partial x} = 1 - 2x \quad (232)$$

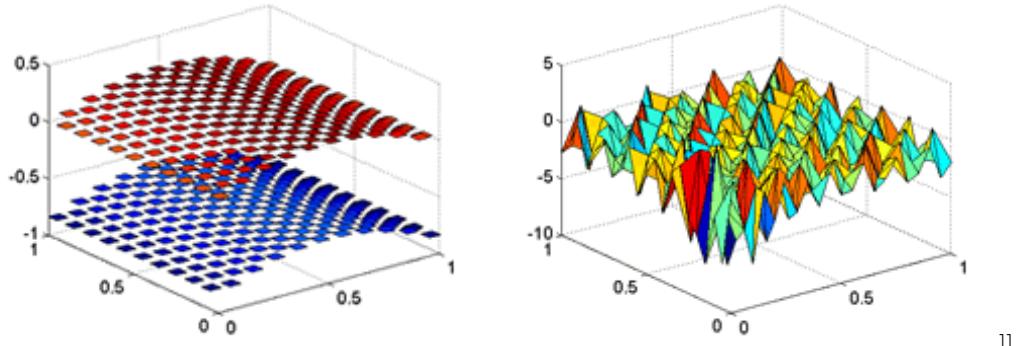
$$\frac{\partial p}{\partial y} = 0 \quad (233)$$

The velocity and pressure fields look like:



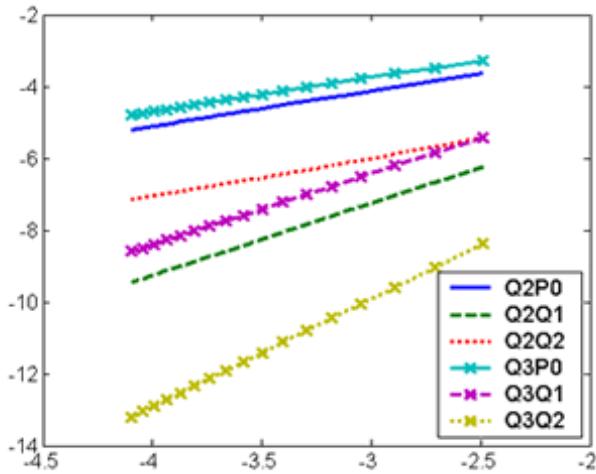
http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible_Ex1.htm

As shown in [165], If the LBB condition is not satisfied, spurious oscillations spoil the pressure approximation. Figures below show results obtained with a mesh of 20x20 Q1P0 (left) and P1P1 (right) elements:



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

7.4.2 Analytical benchmark II - "Dohrmann & Bochev 2D"

Taken from [164]. It is for a unit square with $\nu = \mu/\rho = 1$ and the smooth exact solution is

$$u(x, y) = x + x^2 - 2xy + x^3 - 3xy^2 + x^2y \quad (234)$$

$$v(x, y) = -y - 2xy + y^2 - 3x^2y + y^3 - xy^2 \quad (235)$$

$$p(x, y) = xy + x + y + x^3y^2 - 4/3 \quad (236)$$

Note that the pressure obeys $\int_{\Omega} p d\Omega = 0$

$$b_x = -(1 + y - 3x^2y^2) \quad (237)$$

$$b_y = -(1 - 3x - 2x^3y) \quad (238)$$

7.4.3 Analytical benchmark III - "Dohrmann & Bochev 3D"

This benchmark begins by postulating a polynomial solution to the 3D Stokes equation [164]:

$$\mathbf{v} = \begin{pmatrix} x + x^2 + xy + x^3y \\ y + xy + y^2 + x^2y^2 \\ -2z - 3xz - 3yz - 5x^2yz \end{pmatrix} \quad (239)$$

and

$$p = xyz + x^3y^3z - 5/32 \quad (240)$$

While it is then trivial to verify that this velocity field is divergence-free, the corresponding body force of the Stokes equation can be computed by inserting this solution into the momentum equation with a given viscosity μ (constant or position/velocity/strain rate dependent). The domain is a unit cube and velocity boundary conditions simply use Eq. (422). Note that the pressure fulfills

$$\int_{\Omega} p(\vec{r}) d\Omega = 0.$$

Constant viscosity In this case, the right hand side writes:

$$\begin{aligned} \mathbf{f} &= -\nabla p + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix} \\ &= - \begin{pmatrix} yz + 3x^2y^3z \\ xz + 3x^3y^2z \\ xy + x^3y^3 \end{pmatrix} + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix} \end{aligned}$$

We can compute the components of the strainrate tensor:

$$\dot{\varepsilon}_{xx} = 1 + 2x + y + 3x^2y \quad (241)$$

$$\dot{\varepsilon}_{yy} = 1 + x + 2y + 2x^2y \quad (242)$$

$$\dot{\varepsilon}_{zz} = -2 - 3x - 3y - 5x^2y \quad (243)$$

$$\dot{\varepsilon}_{xy} = \frac{1}{2}(x + y + 2xy^2 + x^3) \quad (244)$$

$$\dot{\varepsilon}_{xz} = \frac{1}{2}(-3z - 10xyz) \quad (245)$$

$$\dot{\varepsilon}_{yz} = \frac{1}{2}(-3z - 5x^2z) \quad (246)$$

Note that we of course have $\dot{\varepsilon}_{xx} + \dot{\varepsilon}_{yy} + \dot{\varepsilon}_{zz} = 0$.

Variable viscosity In this case, the right hand side is obtained through

$$\begin{aligned} \mathbf{f} &= -\nabla p + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix} \\ &+ \begin{pmatrix} 2\dot{\varepsilon}_{xx} \\ 2\dot{\varepsilon}_{xy} \\ 2\dot{\varepsilon}_{xz} \end{pmatrix} \frac{\partial \mu}{\partial x} + \begin{pmatrix} 2\dot{\varepsilon}_{xy} \\ 2\dot{\varepsilon}_{yy} \\ 2\dot{\varepsilon}_{yz} \end{pmatrix} \frac{\partial \mu}{\partial y} + \begin{pmatrix} 2\dot{\varepsilon}_{xz} \\ 2\dot{\varepsilon}_{yz} \\ 2\dot{\varepsilon}_{zz} \end{pmatrix} \frac{\partial \mu}{\partial z} \end{aligned} \quad (247)$$

The viscosity can be chosen to be a smooth varying function:

$$\mu = \exp(1 - \beta(x(1-x) + y(1-y) + z(1-z))) \quad (248)$$

Choosing $\beta = 0$ yields a constant velocity $\mu = e^1$ (and greatly simplifies the right-hand side). One can easily show that the ratio of viscosities μ^* in the system follows $\mu^* = \exp(-3\beta/4)$ so that choosing $\beta = 10$

yields $\mu^* \simeq 1808$ and $\beta = 20$ yields $\mu^* \simeq 3.269 \times 10^6$. In this case

$$\frac{\partial \mu}{\partial x} = -4\beta(1-2x)\mu(x,y,z) \quad (249)$$

$$\frac{\partial \mu}{\partial y} = -4\beta(1-2y)\mu(x,y,z) \quad (250)$$

$$\frac{\partial \mu}{\partial z} = -4\beta(1-2z)\mu(x,y,z) \quad (251)$$

[113] has carried out this benchmark for $\beta = 4$, i.e.:

$$\mu(x,y,z) = \exp(1 - 4(x(1-x) + y(1-y) + z(1-z)))$$

In a unit cube, this yields a variable viscosity such that $0.1353 < \mu < 2.7182$, i.e. a ratio of approx. 20 within the domain. We then have:

$$\frac{\partial \mu}{\partial x} = -4(1-2x)\mu(x,y,z) \quad (252)$$

$$\frac{\partial \mu}{\partial y} = -4(1-2y)\mu(x,y,z) \quad (253)$$

$$\frac{\partial \mu}{\partial z} = -4(1-2z)\mu(x,y,z) \quad (254)$$

sort out mess wrt Eq 26 of busa13

7.4.4 Analytical benchmark IV - "Bercovier & Engelman"

From [53]. The two-dimensional domain is a unit square. The body forces are:

$$\begin{aligned} f_x &= 128[x^2(x-1)^212(2y-1) + 2(y-1)(2y-1)y(12x^2-12x+2)] \\ f_y &= 128[y^2(y-1)^212(2x-1) + 2(x-1)(2x-1)y(12y^2-12y+2)] \end{aligned} \quad (255)$$

The solution is

$$\begin{aligned} u &= -256x^2(x-1)^2y(y-1)(2y-1) \\ v &= 256x^2(y-1)^2x(x-1)(2x-1) \\ p &= 0 \end{aligned} \quad (256)$$

Another choice:

$$\begin{aligned} f_x &= 128[x^2(x-1)^212(2y-1) + 2(y-1)(2y-1)y(12x^2-12x+2)] + y - 1/2 \\ f_y &= 128[y^2(y-1)^212(2x-1) + 2(x-1)(2x-1)y(12y^2-12y+2)] + x - 1/2 \end{aligned} \quad (257)$$

The solution is

$$\begin{aligned} u &= -256x^2(x-1)^2y(y-1)(2y-1) \\ v &= 256x^2(y-1)^2x(x-1)(2x-1) \\ p &= (x-1/2)(y-1/2) \end{aligned} \quad (258)$$

7.4.5 Analytical benchmark V - "VJ"

This is taken from Appendix D1 of [330].

The domain Ω is a unit square. We consider the stream function

$$\phi(x,y) = 1000x^2(1-x)^4y^3(1-y)^2$$

The velocity field is defined by

$$u(x, y) = \partial_y \phi = 1000(x^2(1-x)^4y^2(1-y)(3-5y)) \quad (259)$$

$$v(x, y) = -\partial_x \phi = 1000(-2x(1-x)^3(1-3x)y^3(1-y)^2) \quad (260)$$

and it is easy to verify that $\vec{\nabla} \cdot \vec{v} = 0$.

The pressure is given by:

$$p(x, y) = \pi^2(xy^3 \cos(2\pi x^2 y) - x^2 y \sin(2\pi x y)) + \frac{1}{8}$$

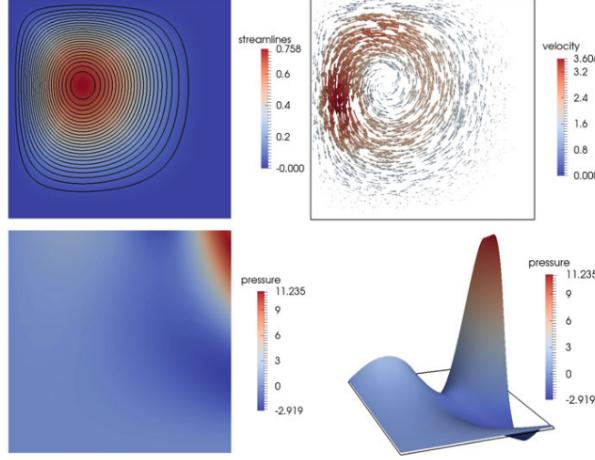


Fig. D.1 Example D.3. Stream function (top left) velocity (top right) and pressure (bottom). These plots are based on results obtained with numerical simulations

Taken from [330].

7.4.6 Analytical benchmark VI - "Ilinca & Pelletier"

This is taken from [319].

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 Φ 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

$$T_M = T(y = 0) = T_0 + \frac{1}{3} \frac{\eta u_0^2}{k}$$

7.4.7 Analytical benchmark VII - "grooves"

This benchmark was designed by Dave May. The velocity and pressure fields are given by

$$\begin{aligned} u(x, y) &= x^3y + x^2 + xy + x \\ v(x, y) &= -\frac{3}{2}x^2y^2 - 2xy - \frac{1}{2}y^2 - y \\ p(x, y) &= x^2y^2 + xy + 5 + p_0 \end{aligned} \quad (261)$$

where p_0 is a constant to be determined based on the type of pressure normalisation. The viscosity is chosen to be

$$\eta(x, y) = -\sin(p) + 1 + \epsilon = -\sin(x^2y^2 + xy + 5) + 1 + \epsilon \quad (262)$$

where ϵ actually controls the viscosity contrast. Note that inserting the polynomial expression of the pressure inside the viscosity expression makes the problem linear. We have

$$\begin{aligned} \dot{\varepsilon}_{xx} &= \frac{\partial u}{\partial x} = 3x^2y + 2x + y + 1 \\ \dot{\varepsilon}_{yy} &= \frac{\partial v}{\partial y} = -3x^2y - 2x - y - 1 \\ \dot{\varepsilon}_{xy} &= \frac{1}{2} \left(\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \right) = \frac{1}{2} (x^3 + x - 3xy^2 - 2y) \end{aligned} \quad (263)$$

and we can verify that the velocity field is incompressible since $\vec{\nabla} \cdot \vec{v} = \dot{\varepsilon}_{xx} + \dot{\varepsilon}_{yy} = 0$. 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

$$\begin{aligned} -\frac{\partial p}{\partial x} + \frac{\partial s_{xx}}{\partial x} + \frac{\partial s_{yx}}{\partial y} + f_x &= 0 \\ -\frac{\partial p}{\partial y} + \frac{\partial s_{xy}}{\partial x} + \frac{\partial s_{yy}}{\partial y} + f_y &= 0 \end{aligned} \quad (264)$$

with

$$\begin{aligned} \frac{\partial s_{xx}}{\partial x} &= \frac{\partial(2\eta\dot{\varepsilon}_{xx})}{\partial x} = 2\eta \frac{\partial\dot{\varepsilon}_{xx}}{\partial x} + 2\frac{\partial\eta}{\partial x}\dot{\varepsilon}_{xx} \\ \frac{\partial s_{zx}}{\partial z} &= \frac{\partial(2\eta\dot{\varepsilon}_{zx})}{\partial z} = 2\eta \frac{\partial\dot{\varepsilon}_{zx}}{\partial z} + 2\frac{\partial\eta}{\partial z}\dot{\varepsilon}_{zx} \\ \frac{\partial s_{xz}}{\partial x} &= \frac{\partial(2\eta\dot{\varepsilon}_{xz})}{\partial x} = 2\eta \frac{\partial\dot{\varepsilon}_{xz}}{\partial x} + 2\frac{\partial\eta}{\partial x}\dot{\varepsilon}_{xz} \\ \frac{\partial s_{zz}}{\partial z} &= \frac{\partial(2\eta\dot{\varepsilon}_{zz})}{\partial z} = 2\eta \frac{\partial\dot{\varepsilon}_{zz}}{\partial z} + 2\frac{\partial\eta}{\partial z}\dot{\varepsilon}_{zz} \\ \frac{\partial\eta}{\partial x} &= -z(2xz + 1)\cos(x^2z^2 + xz + 5) \\ \frac{\partial\eta}{\partial z} &= -x(2xz + 1)\cos(x^2z^2 + xz + 5) \\ \frac{\partial\dot{\varepsilon}_{xx}}{\partial x} &= 6xz + 2 \\ \frac{\partial\dot{\varepsilon}_{zx}}{\partial z} &= -3xz - 1 \\ \frac{\partial\dot{\varepsilon}_{xz}}{\partial x} &= \frac{1}{2}(3x^2 + 1 - 3z^2) \\ \frac{\partial\dot{\varepsilon}_{zz}}{\partial z} &= -3x^2 - 1 \end{aligned}$$

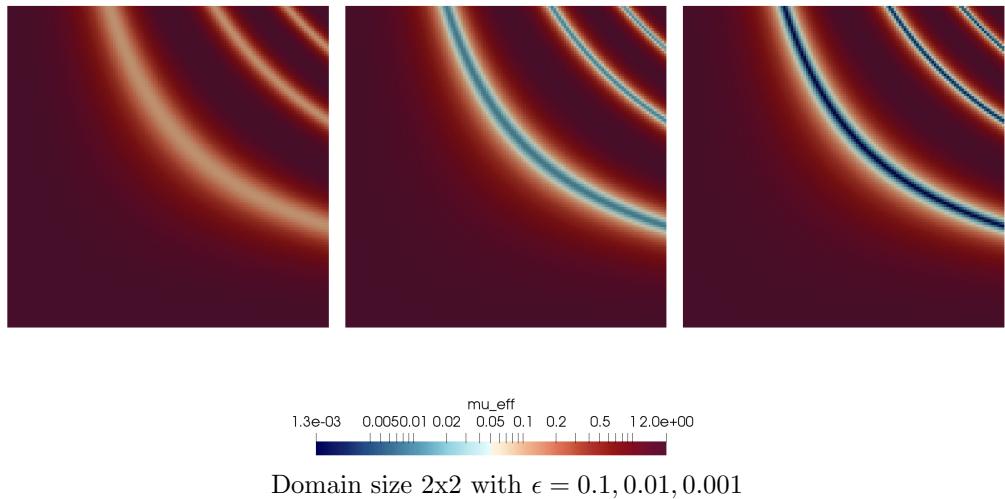
Velocity boundary conditions are prescribed on all four boundaries so that the pressure is known up to a constant (the pressure solution has a nullspace), and the p_0 constant can be determined by requiring that

$$\int_0^L \int_0^L p(x, y) dx dy = \int_0^L \int_0^L (x^2 y^2 + xy + 5) dx dy + \int_0^L \int_0^L p_0 dx dy = \int_0^L \int_0^L (x^2 y^2 + xy + 5) dx dy + p_0 L^2 = 0$$

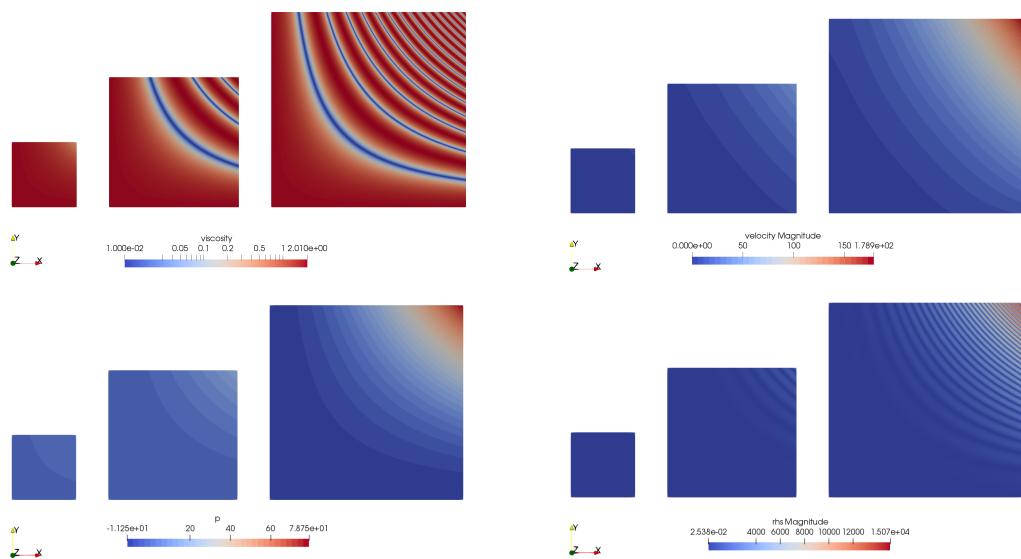
where L is the size of the square domain. Then

$$p_0 = -\frac{1}{L^2} \int_0^L \int_0^L (x^2 y^2 + xy + 5) dx dy = -\frac{L^4}{9} - \frac{L^2}{4} - 5$$

As seen in the following figure, the value of ϵ 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 ϵ .



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:

The screenshot shows the WolframAlpha interface. The input query is: $\int_{-1}^1 \int_{-1}^1 (x^3 y + x^2 + x y + x)^2 + (-1.5 x^2 y^2 - 2 x y - 0.5 y^2 - y)^2 dx dy$. The result is displayed as $\frac{861752}{1575} \approx 547.144$. Below the result, there are buttons for "More digits", "Open code", and "Surprise Me".

and we end up with (for $L = 2$)

$$v_{rms} = \sqrt{\frac{1}{L^2} \frac{861752}{1575}} = \sqrt{\frac{215438}{1575}} \simeq 11.6955560683$$

7.4.8 Analytical benchmark VIII - "Kovasznay"

This flow was published by L.I.G. Kovasznay in 1948 [356]. 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.

$$u(x, y) = 1 - \exp(\lambda x) \cos(2\pi y) \quad v(x, y) = \frac{\lambda}{2\pi} \exp(\lambda x) \sin(2\pi y) \quad \lambda = \frac{Re}{2} - \sqrt{\frac{Re^2}{4} + 4\pi^2}$$

Following step-55 of deal.II¹¹ 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.

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)))
+ 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));
```

¹¹https://www.dealii.org/current/doxygen/deal.II/step_55.html

7.5 Geodynamical benchmarks

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 [366, 337, 473]
- 2D Rayleigh-Benard convection (Blankenbach) [64, 553, 129, 344, 370, 586, 552]
- 2D Rayleigh-Taylor convection/instability [456, 552, 513, 529, 71, 38, 473, 489, 370, 586, 572, 1]
- subduction problems [497, 570]
- numerical sandbox [91, 96]
- the Stokes sphere [360]
- 2D compressible Stokes flow problem [369]
- 3D convection at infinite Prandtl number (Busse) [114, 553]
- Free surface evolution [144]

go through my papers and add relevant ones here

7.6 Assigning values to quadrature points

As we have seen in Section 6, 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 [299] 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 shape 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 à 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 [162] and [171] 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:

$$M_p(a_1, \dots, a_n) = \left(\frac{1}{n} \sum_{i=1}^n a_i^p \right)^{1/p} \quad (265)$$

and it is trivial to verify that we then have the special cases:

$$M_{-\infty} = \lim_{p \rightarrow -\infty} M_p = \min(a_1, \dots, a_n) \quad (\text{minimum}) \quad (266)$$

$$M_{-1} = \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \dots + \frac{1}{a_n}} \quad (\text{harm. avrg.}) \quad (267)$$

$$M_0 = \lim_{p \rightarrow 0} M_p = \left(\prod_{i=1}^n a_i \right)^{1/n} \quad (\text{geom. avrg.}) \quad (268)$$

$$M_{+1} = \frac{1}{n} \sum_{i=1}^n a_i \quad (\text{arithm. avrg.}) \quad (269)$$

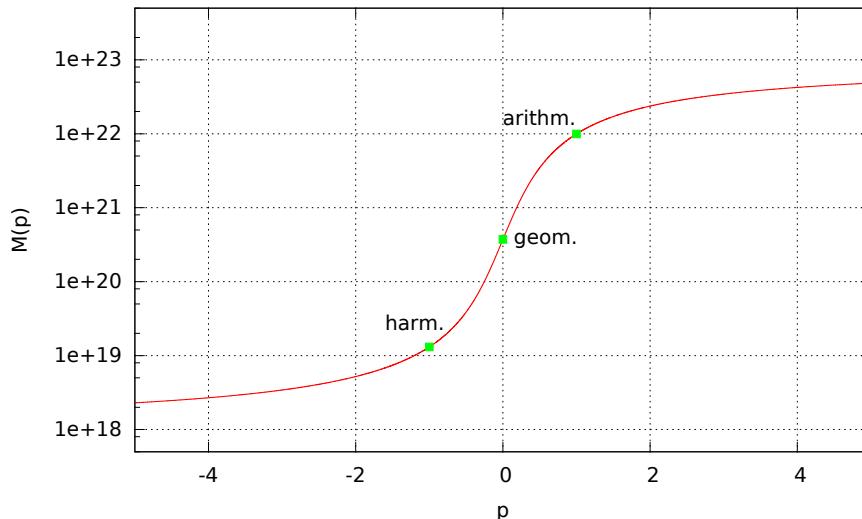
$$M_{+2} = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2} \quad (\text{root mean square}) \quad (270)$$

$$M_{+\infty} = \lim_{p \rightarrow +\infty} M_p = \max(a_1, \dots, a_n) \quad (\text{maximum}) \quad (271)$$

Note that the proofs of the limit convergence are given in [99].

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 [497].

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
```

7.7 Matrix (Sparse) storage

The FE matrix 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 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. [487]

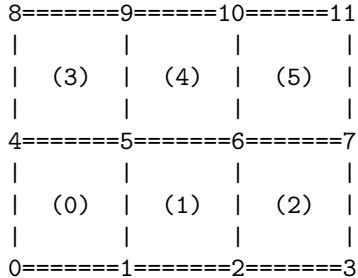
There are several standard formats:

- compressed sparse row (CSR) format
- compressed sparse column format (CSC)
- the Coordinate Format (COO)
- Skyline Storage Format
- ...

I focus on the CSR format in what follows.

7.7.1 2D domain - 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 will look like this:

$$\left(\begin{array}{ccccccccc} X & X & & X & X & & & & \\ X & X & X & X & X & X & & & \\ & X & X & X & X & X & X & & \\ & & X & X & & X & X & & \\ X & X & & X & X & & X & X & \\ X & X & X & X & X & X & X & X & \\ & X & X & X & X & X & X & X & \\ & & X & X & & X & X & & \\ & & & X & X & X & X & & \\ & & & & X & X & X & X & \\ & & & & & X & X & X & \\ & & & & & & X & X & \\ \end{array} \right)$$

where the X stand for non-zero terms. This matrix structure stems from the fact that

- node 0 sees nodes 0,1,4,5
- node 1 sees nodes 0,1,2,4,5,6
- node 2 sees nodes 1,2,3,5,6,7
- ...

- node 5 sees nodes 0,1,2,4,5,6,8,9,10
- ...
- node 10 sees nodes 5,6,7,9,10,11
- node 11 sees nodes 6,7,10,11

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 is then:

$$NZ = 4 \times 4 + 4 \times 6 + 2 \times 6 + 2 \times 9 = 70$$

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 $nnx = nny = n$. Then the matrix size (total number of unknowns) is $N = n^2$ 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:

$$COLIND = (0, 1, 4, 5, 0, 1, 2, 4, 5, 6, 1, 2, 3, 5, 6, 7, \dots, 6, 7, 10, 11)$$

$$RWPTR = (0, 4, 10, 16, \dots)$$

7.7.2 2D domain - 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

<code>NfemV=nnp*ndofV</code>

In the case of the small grid above, we have `NfemV=24`. 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

- $(nnx-2) \times (nny-2)$ 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:

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:

$$COLIND = (0, 1, 2, 3, 8, 9, 10, 11, 0, 1, 2, 3, 8, 9, 10, 11, \dots)$$

$$RW PTR = (0, 8, 16, 28, \dots)$$

7.7.3 in fieldstone

The majority of the 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)
```

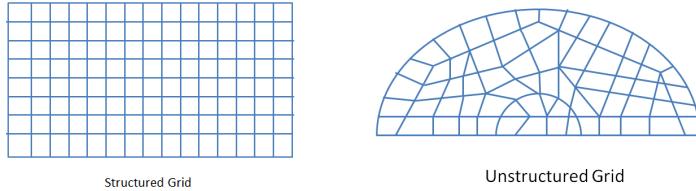
Note that linked list storages can be used (`lil_matrix`). Substantial memory savings but much longer compute times.

7.8 Mesh generation

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 [544]. 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/softwares must be used (see [548, 211, 604]).

We usually distinguish between two broad classes of grids: structured grids (with a regular connectivity) and unstructured grids (with an irregular connectivity).



7.8.1 Quadrilateral-based meshes

Let us now focus on the case of a rectangular computational domain of size $Lx \times Ly$ with a regular mesh composed of $\text{nelx} \times \text{nely} = \text{nel}$ quadrilaterals. There are then $\text{nnx} \times \text{nny} = \text{nnp}$ grid points. The elements are of size $hx \times hy$ with $hx = Lx / \text{nelx}$.

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:

$\begin{array}{ccccccc} 8 & ===== & 9 & ===== & 10 & ===== & 11 \\ & & & & & & \\ & (3) & & (4) & & (5) & \\ & & & & & & \\ 4 & ===== & 5 & ===== & 6 & ===== & 7 \\ & & & & & & \\ & (0) & & (1) & & (2) & \\ & & & & & & \\ 0 & ===== & 1 & ===== & 2 & ===== & 3 \end{array}$	$\begin{array}{ccccccc} 2 & ===== & 5 & ===== & 8 & ===== & 11 \\ & & & & & & \\ & (1) & & (3) & & (5) & \\ & & & & & & \\ 1 & ===== & 4 & ===== & 7 & ===== & 10 \\ & & & & & & \\ & (0) & & (2) & & (4) & \\ & & & & & & \\ 0 & ===== & 3 & ===== & 6 & ===== & 9 \end{array}$
"row by row"	"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 $\text{nnx}-1$ first for $j=0, 1, \dots$ up to $\text{nny}-1$ 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 $m \times \text{nel}$ 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×nely×nelz=nel` elements which represent a cuboid of size `Lx×Ly×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×nel` with `m=8`:

```

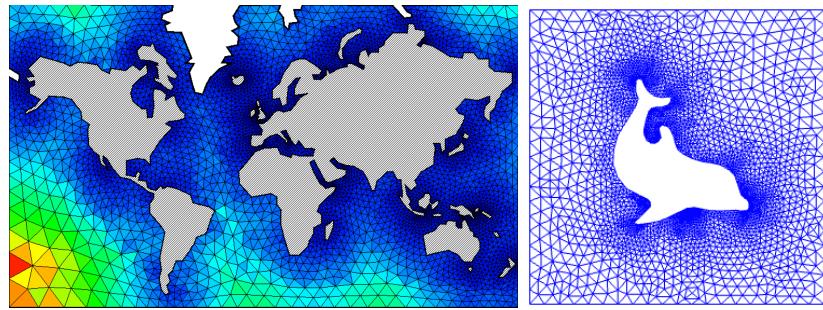
icon =np.zeros((m, nel), dtype=np.int16)
counter = 0
for i in range(0,nelx):
    for j in range(0,nely):
        for k in range(0,nelz):
            icon[0,counter]=nny*nnz*(i    )+nnz*(j    )+k
            icon[1,counter]=nny*nnz*(i+1)+nnz*(j    )+k
            icon[2,counter]=nny*nnz*(i+1)+nnz*(j+1)+k
            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

```

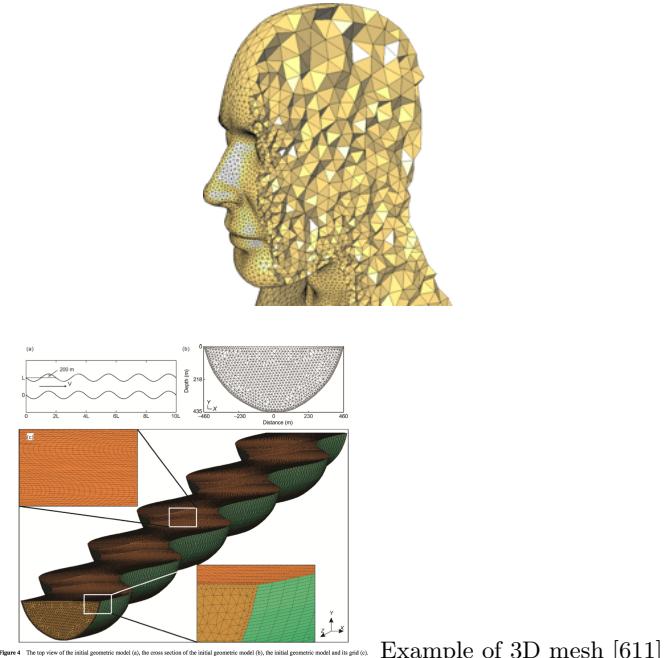
produce drawing of node numbering

7.8.2 Delaunay triangulation and Voronoi cells

Triangle-based meshes are obviously better suited for simulations of complex geometries:



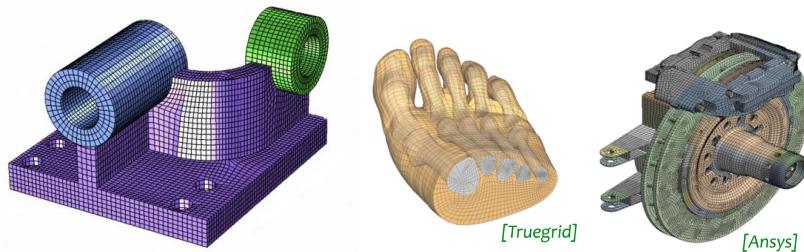
7.8.3 Tetrahedra



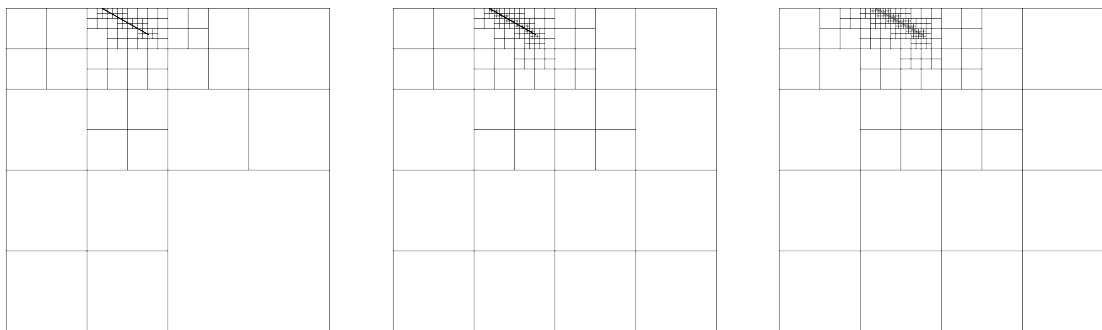
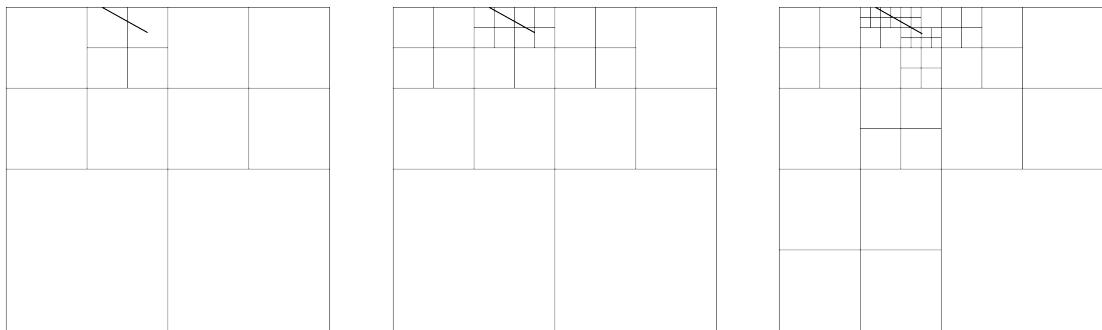
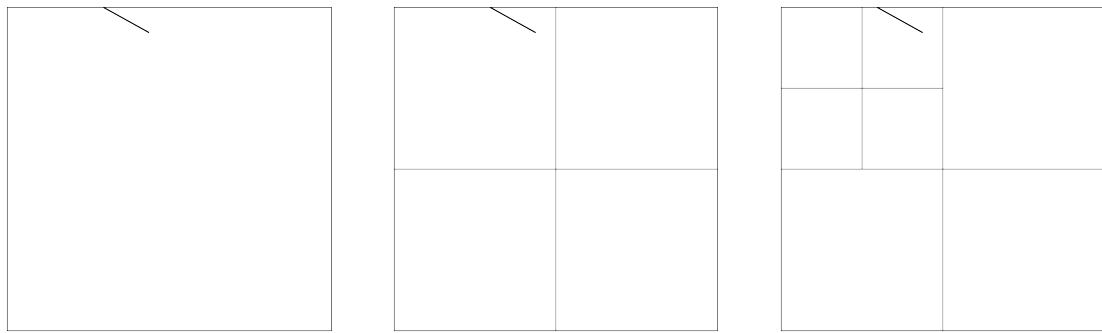
Example of 3D mesh [611]

7.8.4 Hexahedra

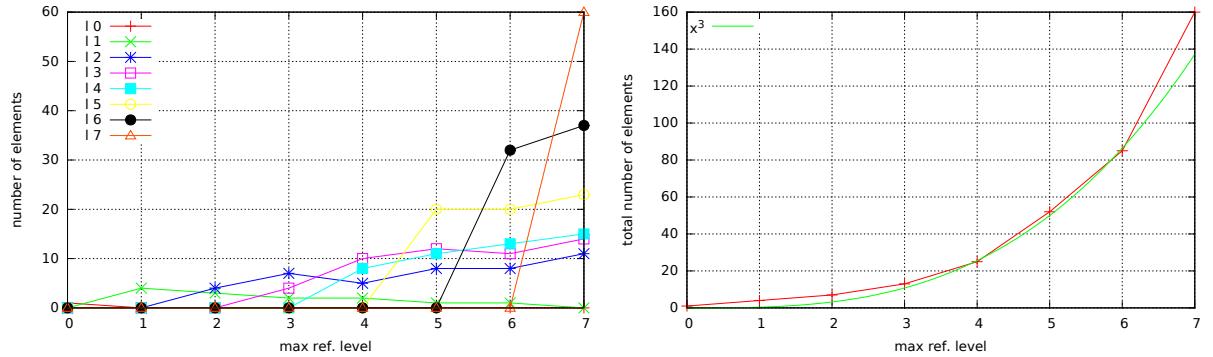
A hexahedron is a convex polytope isomorphic to the cube $[0, 1]^3$. Edges are line segments, facets are strictly **planar** convex polygons.



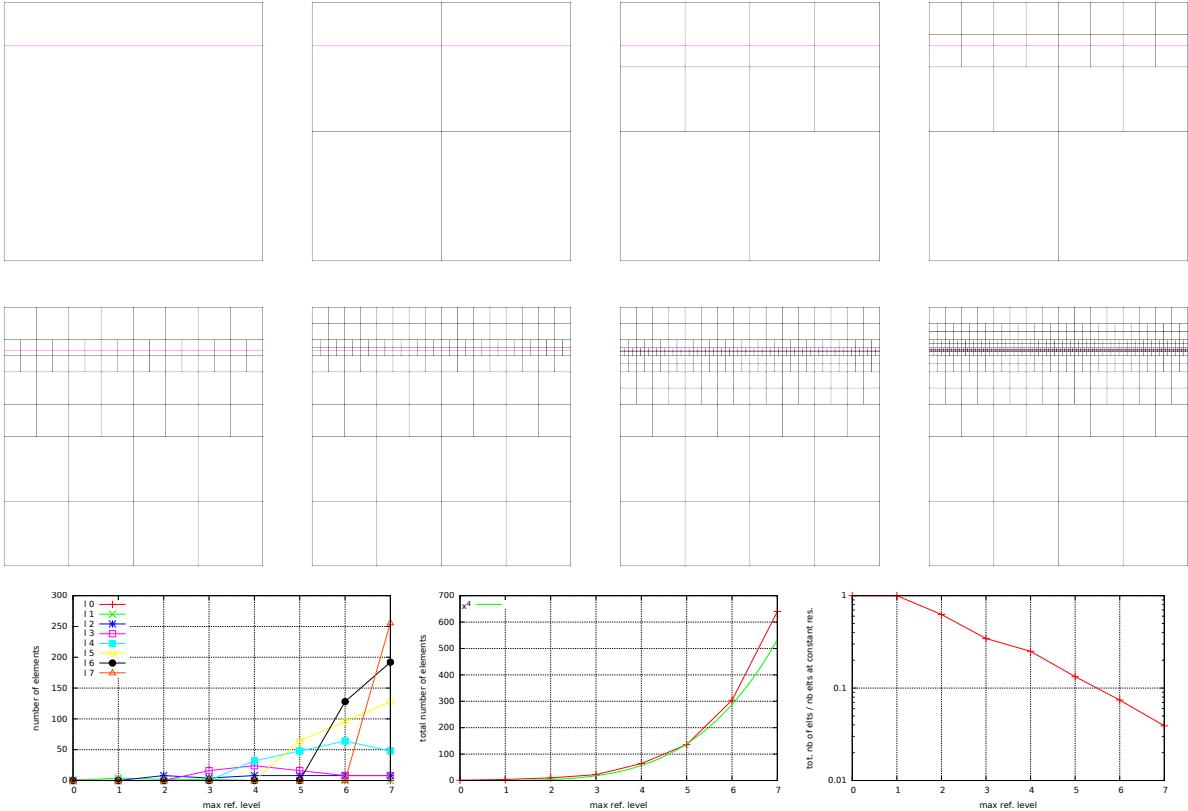
7.8.5 Adaptive Mesh Refinement



	# 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.



7.9 Visco-Plasticity

7.9.1 Tensor invariants

Before we dive into the world of nonlinear rheologies it is necessary to introduce the concept of tensor invariants since they are needed further on. Unfortunately there are many different notations used in the literature and these can prove to be confusing. Note that we only consider symmetric tensors in what follows.

Given a tensor \mathbf{T} , one can compute its (moment) invariants as follows:

- first invariant :

$$\begin{aligned} T_I|^{2D} &= \text{Tr}[\mathbf{T}] = T_{xx} + T_{yy} \\ T_I|^{3D} &= \text{Tr}[\mathbf{T}] = T_{xx} + T_{yy} + T_{zz} \end{aligned}$$

- second invariant :

$$\begin{aligned} T_{II}|^{2D} &= \frac{1}{2} \text{Tr}[\mathbf{T}^2] = \frac{1}{2} \sum_{ij} T_{ij} T_{ji} = \frac{1}{2} (T_{xx}^2 + T_{yy}^2) + T_{xy}^2 \\ T_{II}|^{3D} &= \frac{1}{2} \text{Tr}[\mathbf{T}^2] = \frac{1}{2} \sum_{ij} T_{ij} T_{ji} = \frac{1}{2} (T_{xx}^2 + T_{yy}^2 + T_{zz}^2) + T_{xy}^2 + T_{xz}^2 + T_{yz}^2 \end{aligned}$$

- third invariant :

$$T_{III} = \frac{1}{3} \text{Tr}[\mathbf{T}^3] = \frac{1}{3} \sum_{ijk} T_{ij} T_{jk} T_{ki}$$

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

$$\begin{aligned} \tau_{II}|^{2D} &= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2) + \tau_{xy}^2 \\ &= \frac{1}{4} (\sigma_{xx} - \sigma_{yy})^2 + \sigma_{xy}^2 \\ &= \frac{1}{4} (\sigma_1 - \sigma_2)^2 \\ \tau_{II}|^{3D} &= \frac{1}{2} (\tau_{xx}^2 + \tau_{yy}^2 + \tau_{zz}^2) + \tau_{xy}^2 + \tau_{xz}^2 + \tau_{yz}^2 \\ &= \frac{1}{6} [(\sigma_{xx} - \sigma_{yy})^2 + (\sigma_{yy} - \sigma_{zz})^2 + (\sigma_{xx} - \sigma_{zz})^2] + \sigma_{xy}^2 + \sigma_{xz}^2 + \sigma_{yz}^2 \\ &= \frac{1}{6} [(\sigma_1 - \sigma_2)^2 + (\sigma_2 - \sigma_3)^2 + (\sigma_1 - \sigma_3)^2] \\ \varepsilon_{II}|^{2D} &= \frac{1}{2} [(\dot{\varepsilon}_{xx}^d)^2 + (\dot{\varepsilon}_{yy}^d)^2] + (\dot{\varepsilon}_{xy}^d)^2 \\ &= \frac{1}{2} \left[\frac{1}{4} (\dot{\varepsilon}_{xx} - \dot{\varepsilon}_{yy})^2 + \frac{1}{4} (\dot{\varepsilon}_{yy} - \dot{\varepsilon}_{xx})^2 \right] + \dot{\varepsilon}_{xy}^2 \\ &= \frac{1}{4} (\dot{\varepsilon}_{xx} - \dot{\varepsilon}_{yy})^2 + \dot{\varepsilon}_{xy}^2 \\ \varepsilon_{II}|^{3D} &= \frac{1}{2} [(\dot{\varepsilon}_{xx}^d)^2 + (\dot{\varepsilon}_{yy}^d)^2 + (\dot{\varepsilon}_{zz}^d)^2] + (\dot{\varepsilon}_{xy}^d)^2 + (\dot{\varepsilon}_{xz}^d)^2 + (\dot{\varepsilon}_{yz}^d)^2 \\ &= \frac{1}{6} [(\dot{\varepsilon}_{xx} - \dot{\varepsilon}_{yy})^2 + (\dot{\varepsilon}_{yy} - \dot{\varepsilon}_{zz})^2 + (\dot{\varepsilon}_{xx} - \dot{\varepsilon}_{zz})^2] + \dot{\varepsilon}_{xy}^2 + \dot{\varepsilon}_{xz}^2 + \dot{\varepsilon}_{yz}^2 \end{aligned}$$

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

$$\underline{\tau}_{II} = \sqrt{\tau_{II}} \quad \dot{\underline{\varepsilon}}_{II} = \sqrt{\dot{\varepsilon}_{II}}$$

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

7.9.2 Scalar viscoplasticity

This formulation is quite easy to implement. It is widely used, e.g. [597, 546, 517], and relies on the assumption that a scalar quantity η_p (the 'effective plastic viscosity') exists such that the deviatoric stress tensor

$$\boldsymbol{\tau} = 2\eta_p \dot{\boldsymbol{\varepsilon}} \quad (272)$$

is bounded by some yield stress value Y . From Eq. (272) it follows that $\underline{\tau}_{II} = 2\eta_p \dot{\underline{\varepsilon}}_{II} = Y$ which yields

$$\eta_p = \frac{Y}{2\dot{\underline{\varepsilon}}_{II}}$$

This approach has also been coined the Viscosity Rescaling Method (VRM) [334].

insert here the rederivation 2.1.1 of spmw16

It is at this stage important to realise that (i) in areas where the strainrate is low, the resulting effective viscosity will be large, and (ii) in areas where the strainrate is high, the resulting effective viscosity will be low. This is not without consequences since (effective) viscosity contrasts up to 8-10 orders of magnitude have been observed/obtained with this formulation and it makes the FE matrix very stiff, leading to (iterative) solver convergence issues. In order to contain these viscosity contrasts one usually resorts to viscosity limiters η_{min} and η_{max} such that

$$\eta_{min} \leq \eta_p \leq \eta_{max}$$

Caution must be taken when choosing both values as they may influence the final results.

▷ `python_codes/fieldstone_indentor`

7.9.3 About the yield stress value Y

In geodynamics the yield stress value is often given as a simple function. It can be constant (in space and time) and in this case we are dealing with a von Mises plasticity yield criterion. . We simply assume $Y_{vM} = C$ where C is a constant cohesion independent of pressure, strainrate, deformation history, etc ...

Another model is often used: the Drucker-Prager plasticity model. A friction angle ϕ is then introduced and the yield value Y takes the form

$$Y_{DP} = p \sin \phi + C \cos \phi$$

and therefore depends on the pressure p . Because ϕ is with the range $[0^\circ, 45^\circ]$, Y is found to increase with depth (since the lithostatic pressure often dominates the overpressure).

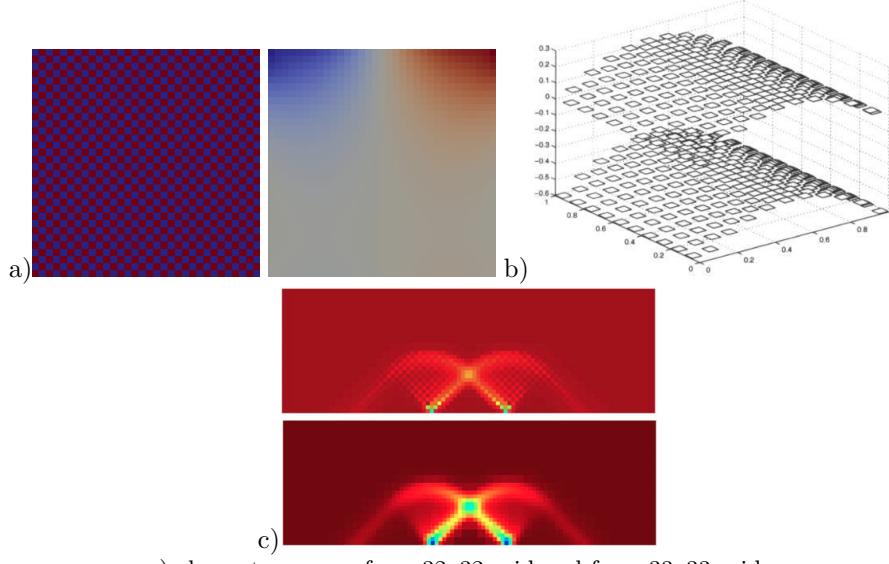
Note that a slightly modified version of this plasticity model has been used: the total pressure p is then replaced by the lithostatic pressure p_{lith} .

7.10 Pressure smoothing

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 [279, 125, 491, 492]. They can be filtered out [125] or simply smoothed [365].

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 is not visible 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 [165, p307] for a manufactured solution;
 c) elemental pressure and smoothed pressure for the punch experiment [546]

The easiest post-processing step that can be used (especially when a regular grid is used) is explained in [546]: "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."

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 as

```
q=np.zeros(nnp, dtype=np.float64)
```

The element-to-node algorithm is then simply (in 2D):

```
count=np.zeros(nnp, dtype=np.int16)
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
```

Pressure smoothing is further discussed in [311].

[produce figure to explain this](#)

[link to proto paper](#)

[link to least square and nodal derivatives](#)

7.11 Pressure scaling

As perfectly explained in the step 32 of deal.ii¹², we often need to scale the \mathbb{G} term since it is many orders of magnitude smaller than \mathbb{K} , 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

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix}$$

and introduce the scaling coefficient as follows (which in fact does not alter the solution at all):

$$\begin{pmatrix} \mathbb{K} & \frac{\eta}{L}\mathbb{G} \\ \frac{\eta}{L}\mathbb{G}^T & -\frac{\eta^2}{L^2}\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \underline{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \underline{h} \end{pmatrix}$$

We then end up with the modified Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \underline{\mathbb{G}}^T & \underline{\mathbb{C}} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \underline{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \underline{h} \end{pmatrix}$$

where

$$\underline{\mathbb{G}} = \frac{\eta}{L}\mathbb{G} \quad \underline{\mathcal{P}} = \frac{L}{\eta}\vec{\mathcal{P}} \quad \underline{\mathbb{C}} = \frac{\eta^2}{L^2}\mathbb{C} \quad \underline{h} = \frac{\eta}{L}\vec{h}$$

After the solve phase, we recover the real pressure with $\vec{\mathcal{P}} = \frac{\eta}{L}\underline{\mathcal{P}}$.

¹²https://www.dealii.org/9.0.0/doxygen/deal.II/step_32.html

7.12 Pressure normalisation

7.12.1 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.

$$\int_{\Omega} p dV = 0 \quad (273)$$

Another possibility is to impose the pressure value at a single node.

Let us assume for example that we are using $Q_1 \times P_0$ elements. Then the pressure is constant inside each element. The integral above becomes:

$$\int_{\Omega} p dV = \sum_e \int_{\Omega_e} p dV = \sum_e p_e \int_{\Omega_e} dV = \sum_e p_e A_e = 0 \quad (274)$$

where the sum runs over all elements e of area A_e . This can be rewritten

$$\mathbb{L}^T \cdot \vec{\mathcal{P}} = 0$$

and it is a constraint on the pressure solution which couples *all* pressure dofs. As we have seen before ??, we can associate to it a Lagrange multiplier λ so that we must solve the modified Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} & 0 \\ \mathbb{G}^T & 0 & \mathbb{L} \\ 0 & \mathbb{L}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \\ \lambda \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \\ 0 \end{pmatrix}$$

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 7.13.

7.12.2 Implementation – the real deal

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 fulfills 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.

7.13 Solving the Stokes system

Let us start again from the (full) Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (275)$$

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.

In what follow I cover:

- solving when the penalty approach is used
- the Schur complement approach
- the FGMRES approach
- the Augmented Lagrangian approach

7.13.1 when using the penalty formulation

In this case we are only solving for velocity since pressure is 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 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 [218], DOUAR [80], or FANTOM [542] relying on the penalty formulation all use direct solvers. The most popular are BLKFCT, MUMPS, WSMP, UMFPACK, SuperLU, PARDISO, CholMod.

Braun et al [80] list the following features of such 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 100^3 element mesh. This explains why direct solvers are often used for 2D problems and rarely in 3D with noticeable exceptions [546, 608, 81, 386, 14, 15, 16, 596, 427].

7.13.2 Conjugate gradient and the Schur complement approach

Let us write the above system as two equations:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (276)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} = \vec{h} \quad (277)$$

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:

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} = \mathbb{G}^T \cdot [\mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}})] = \vec{h} \quad (278)$$

or,

$$(\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}) \cdot \vec{\mathcal{P}} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h} \quad (279)$$

The matrix $\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}$ is called the Schur complement. It is Symmetric (since \mathbb{K} is symmetric) and Positive-Definite¹³ (SPD) if $\text{Ker}(\mathbb{G}) = 0$. [look in donea-huerta book for details](#) Having solved this equation (we have obtained $\vec{\mathcal{P}}$), the velocity can be recovered by solving $\mathbb{K} \cdot \vec{\mathcal{V}} = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}$.

For now, let us assume that we have built the \mathbb{S} matrix and the right hand side $\underline{\vec{f}} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h}$. We must solve $\mathbb{S} \cdot \vec{\mathcal{P}} = \underline{\vec{f}}$.

One can resort to so-called Richardson iterations, defined as follows (e.g., see [580], p141): in solving the matrix equation $\mathbf{A} \cdot \vec{X} = \vec{b}$, the Richardson iterative method is defined by:

$$\vec{X}_{k+1} = \vec{X}_k + \alpha_k (-\mathbf{A} \cdot \vec{X}_k + \vec{b}) \quad m \geq 0 \quad (280)$$

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 $\mathbf{A} \cdot \vec{X} = \vec{b}$ is satisfied. In our case, it writes:

$$\begin{aligned} \vec{\mathcal{P}}_{k+1} &= \vec{\mathcal{P}}_k + \alpha_k (-\mathbb{S} \cdot \vec{\mathcal{P}}_k + \underline{\vec{f}}) \\ &= \vec{\mathcal{P}}_k + \alpha_k (-\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{\mathcal{P}}_k + \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h}) \\ &= \vec{\mathcal{P}}_k + \alpha_k [\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (-\mathbb{G} \cdot \vec{\mathcal{P}}_k + \vec{f}) - \vec{h}] \\ &= \vec{\mathcal{P}}_k + \alpha_k [\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (\mathbb{K} \cdot \vec{\mathcal{V}}_k) - \vec{h}] \\ &= \vec{\mathcal{P}}_k + \alpha_k (\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \vec{h}) \end{aligned} \quad (281)$$

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 [76, p221].

Uzawa algorithm (1):

$$\text{solve} \quad \mathbb{K} \cdot \vec{\mathcal{V}}_k = \underline{\vec{f}} - \mathbb{G} \cdot \vec{\mathcal{P}}_{k-1} \quad (282)$$

$$\vec{\mathcal{P}}_k = \vec{\mathcal{P}}_{k-1} + \alpha_k (\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \vec{h}) \quad k = 1, 2, \dots \quad (283)$$

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 subsection.

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.

Richardson iterations are part of the family of stationary iterative methods, since it can be rewritten

$$\vec{X}_{k+1} = (\mathbf{I} - \alpha_k \mathbf{A}) \cdot \vec{X}_k + \alpha_k \vec{b} \quad (284)$$

which is the definition of a stationary method.

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 [76, p221]. For the analysis of the Uzawa algorithm, we define the residue

$$\vec{\mathcal{R}}_k = \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}_k$$

In addition, suppose the solution of the saddle point problem is denoted by $(\mathcal{V}^*, \mathcal{P}^*)$. Now substituting the iteration formula for \mathcal{V}_k , we get

$$\mathcal{R}_k = \mathbb{G}^T \cdot \vec{\mathcal{V}}^* - \mathbb{G}^T \cdot \mathbb{K}^{-1} (\vec{f} - \mathbb{G} \cdot \mathcal{P}_{k-1}) \quad (285)$$

$$= \mathbb{G}^T \cdot \vec{\mathcal{V}}^* - \mathbb{G}^T \cdot \mathbb{K}^{-1} (\mathbb{K} \cdot \vec{\mathcal{V}}^* + \mathbb{G} \cdot \vec{\mathcal{P}}^* - \mathbb{G} \cdot \mathcal{P}_{k-1}) \quad (286)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot (\vec{\mathcal{P}}_{k-1} - \vec{\mathcal{P}}^*) \quad (287)$$

¹³ M positive definite $\iff x^T M x > 0 \forall x \in \mathbb{R}^n \setminus \mathbf{0}$

From Eq. 283 it follows that:

$$\mathcal{P}_k - \mathcal{P}_{k-1} = \alpha(\mathbb{G}^T \cdot \vec{\mathcal{V}}_k - \vec{h}) \quad (288)$$

$$= -\alpha \vec{\mathcal{R}}_k \quad (289)$$

$$= -\alpha \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot (\vec{\mathcal{P}}_{k-1} - \vec{\mathcal{P}}^*) \quad (290)$$

$$= \alpha \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot (\vec{\mathcal{P}}^* - \vec{\mathcal{P}}_{k-1}) \quad (291)$$

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 $\alpha < 2\|\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}\|^{-1}$ and one can show that the good step size α_k is given by

$$\alpha_k = \frac{\mathcal{R}_k \cdot \mathcal{R}_k}{(\mathbb{G}q_k) \cdot (\mathbb{K}^{-1}\mathbb{G}q_k)} \quad (292)$$

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 in [260] it is stated: the convergence of this algorithm is proved for $\alpha \in (0, 2\mu/d)$ (where d is the number of dimensions).

check this, and report page number

Note that this algorithm is presented in [636] in the context of viscoplastic flow.

As mentioned above, there is a way to rework the original Uzawa algorithm to include Eq. (292). It yields a modified Uzawa algorithm [76, p221]:

Uzawa algorithm (2): Solve $\mathbb{K} \cdot \vec{\mathcal{V}}_1 = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0$. For $k = 1, 2, \dots$, compute

$$\vec{q}_k = \vec{h} - \mathbb{G}^T \cdot \vec{\mathcal{V}}_k \quad (293)$$

$$\vec{p}_k = \mathbb{G} \cdot q_k \quad (294)$$

$$\vec{H}_k = \mathbb{K}^{-1} \cdot \vec{p}_k \quad (295)$$

$$\alpha_k = \frac{\vec{q}_k \cdot \vec{q}_k}{\vec{p}_k \cdot \vec{H}_k} \quad (296)$$

$$\vec{\mathcal{P}}_k = \vec{\mathcal{P}}_{k-1} - \alpha_k \vec{q}_k \quad (297)$$

$$\vec{\mathcal{V}}_{k+1} = \vec{\mathcal{V}}_k + \alpha_k \vec{H}_k \quad (298)$$

7.13.3 Conjugate gradient and the Schur complement approach

Since \mathbb{S} is SPD, the Conjugate Gradient (CG) method 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:

```

r0 := b - Ax0
if r0 is sufficiently small, then return x0 as the result
p0 := r0
k := 0
repeat
     $\alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$ 
     $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$ 
     $\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$ 
    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

```

Algorithm as obtained from Wikipedia¹⁴

This algorithm is of course explained in detail in many textbooks such as [487]

[add biblio](#)

Let us look at this algorithm up close. The parts which may prove to be somewhat tricky are those involving the matrix inverse (in our case the Schur complement). We start the iterations with a guess pressure \vec{P}_0 (and an initial guess velocity which could be obtained by solving $\mathbb{K} \cdot \vec{V}_0 = \vec{f} - \mathbb{G} \cdot \vec{P}_0$).

$$\vec{r}_0 = \vec{f} - \mathbb{S} \cdot \vec{P}_0 \quad (299)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} - \vec{h} - (\mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G}) \cdot \vec{P}_0 \quad (300)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{P}_0) - \vec{h} \quad (301)$$

$$= \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{K} \cdot \vec{V}_0 - \vec{h} \quad (302)$$

$$= \mathbb{G}^T \cdot \vec{V}_0 - \vec{h} \quad (303)$$

$$= \mathbb{G}^T \cdot \vec{V}_0 - \vec{h} \quad (304)$$

We now turn to the α_k coefficient:

$$\alpha_k = \frac{\vec{r}_k^T \cdot \vec{r}_k}{\vec{p}_k \cdot \mathbb{S} \cdot \vec{p}_k} = \frac{\vec{r}_k^T \cdot \vec{r}_k}{\vec{p}_k \cdot \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{p}_k} = \frac{\vec{r}_k^T \cdot \vec{r}_k}{(\mathbb{G} \cdot \vec{p}_k)^T \cdot \mathbb{K}^{-1} \cdot (\mathbb{G} \cdot \vec{p}_k)}$$

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 $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k)$

Then we need to look at the term $\mathbb{S} \cdot \vec{p}_k$:

$$\mathbb{S} \cdot \vec{p}_k = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{p}_k = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \tilde{\vec{p}}_k = \mathbb{G}^T \cdot \vec{d}_k$$

We can then rewrite the CG algorithm as follows [629]:

- $\vec{r}_0 = \mathbb{G}^T \cdot \vec{V}_0 - \vec{h}$
- if \vec{r}_0 is sufficiently small, then return (\vec{V}_0, \vec{P}_0) as the result
- $\vec{p}_0 = \vec{r}_0$

¹⁴https://en.wikipedia.org/wiki/Conjugate_gradient_method

- $k = 0$
- repeat
 - compute $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$
 - solve $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
 - compute $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k)$
 - $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$
 - $\vec{r}_{k+1} = \vec{r}_k - \alpha_k \mathbb{G}^T \cdot \vec{d}_k$
 - 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, 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

$$\vec{\mathcal{V}}_{k+1} = \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot \vec{\mathcal{P}}_{p+1}) \quad (305)$$

$$= \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot (\vec{\mathcal{P}}_k + \alpha_k \vec{p}_k)) \quad (306)$$

$$= \mathbb{K}^{-1} \cdot (f - \mathbb{G} \cdot \vec{\mathcal{P}}_k) - \alpha_k \mathbb{K}^{-1} \cdot \mathbb{G} \cdot \vec{p}_k \quad (307)$$

$$= \vec{\mathcal{V}}_k - \alpha_k \mathbb{K}^{-1} \cdot \tilde{\vec{p}}_k \quad (308)$$

$$= \vec{\mathcal{V}}_k - \alpha_k \vec{d}_k \quad (309)$$

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:

- compute $\vec{\mathcal{V}}_0 = \mathbb{K}^{-1} \cdot (\vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}_0)$
- $\vec{r}_0 = \mathbb{G}^T \cdot \vec{\mathcal{V}}_0 - \vec{h}$
- 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 $\alpha_k = (\vec{r}_k^T \cdot \vec{r}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k)$
 - $\vec{\mathcal{P}}_{k+1} = \vec{\mathcal{P}}_k + \alpha_k \vec{p}_k$
 - $\vec{\mathcal{V}}_{k+1} = \vec{\mathcal{V}}_k - \alpha_k \vec{d}_k$
 - $\vec{r}_{k+1} = \vec{r}_k - \alpha_k \mathbb{G}^T \cdot \vec{d}_k$
 - if \vec{r}_{k+1} is sufficiently small ($\|\vec{r}_{k+1}\|_2 / \|\vec{r}_0\|_2 < tol$), then exit loop
 - $\beta_k = (r_{k+1}^T r_{k+1}) / (r_k^T r_k)$

- $\vec{p}_{k+1} = \vec{r}_{k+1} + \beta_k \vec{p}_k$
- $k = k + 1$
- return $\vec{\mathcal{P}}_{k+1}$ as result

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 solver, iterative solver, ...)

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 a quick Google search yields:

```

r0 := b - Ax0
z0 := M-1r0
p0 := z0
k := 0
repeat
     $\alpha_k := \frac{\mathbf{r}_k^\top \mathbf{z}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k}$ 
     $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$ 
     $\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$ 
    if  $\mathbf{r}_{k+1}$  is sufficiently small then exit loop end if
     $\mathbf{z}_{k+1} := \mathbf{M}^{-1} \mathbf{r}_{k+1}$ 
     $\beta_k := \frac{\mathbf{z}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{z}_k^\top \mathbf{r}_k}$ 
     $\mathbf{p}_{k+1} := \mathbf{z}_{k+1} + \beta_k \mathbf{p}_k$ 
    k := k + 1
end repeat
The result is xk+1

```

Algorithm obtained from Wikipedia¹⁵.

Note that in the algorithm above the preconditioner matrix 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 M at each iteration, which means that M must be such that solving $M \cdot \vec{x} = \vec{f}$ 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:

```
- compute  $\mathcal{V}_0 = \mathbb{K}^{-1}(f - \mathbb{G}\mathcal{P}_0)$
  - $r_0 = \mathbb{G}^T \mathcal{V}_0 - h$
  - if  $\vec{r}_0$  is sufficiently small, then return  $(\vec{\mathcal{V}}_0, \vec{\mathcal{P}}_0)$  as the result
  - $\vec{z}_0 = M^{-1} \cdot \vec{r}_0$
  - $\vec{p}_0 = \vec{z}_0$
  - $k = 0$
  - repeat
    - compute  $\tilde{\vec{p}}_k = \mathbb{G} \cdot \vec{p}_k$

<sup>15</sup>[https://en.wikipedia.org/wiki/Conjugate\\_gradient\\_method](https://en.wikipedia.org/wiki/Conjugate_gradient_method)

- solve  $\mathbb{K} \cdot \vec{d}_k = \tilde{\vec{p}}_k$
- compute  $\alpha_k = (\vec{r}_k^T \cdot \vec{z}_k) / (\tilde{\vec{p}}_k^T \cdot \vec{d}_k)$
- $\vec{\mathcal{P}}_{k+1} = \mathcal{P}_k + \alpha_k \vec{p}_k$
- $\vec{\mathcal{V}}_{k+1} = \mathcal{V}_k - \alpha_k \vec{d}_k$
- $\vec{r}_{k+1} = \vec{r}_k - \alpha_k \mathbb{G}^T \cdot \vec{d}_k$
- if  $r_{k+1}$  is sufficiently small ( $\|r_{k+1}\|_2 / \|r_0\|_2 < tol$ ), then exit loop
- $\vec{z}_{k+1} = M^{-1} \cdot r_{k+1}$
- $\beta_k = (\tilde{\vec{z}}_{k+1}^T \cdot \vec{r}_{k+1}) / (\tilde{\vec{z}}_k^T \cdot \vec{r}_k)$
- $\vec{p}_{k+1} = \vec{z}_{k+1} + \beta_k \vec{p}_k$
- $k = k + 1$
- return  $\vec{\mathcal{P}}_{k+1}$  as result

Following [629] one can define the following matrix as preconditioner:

$$M = \text{diag} [\mathbb{G}^T (\text{diag}[\mathbb{K}])^{-1} \mathbb{G}]$$

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.

how to compute  $M$  for the Schur complement ?

#### 7.13.4 The Augmented Lagrangian approach

see LaCoDe paper.

We start from the saddle point Stokes system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{\mathcal{P}} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} \end{pmatrix} \quad (310)$$

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:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\lambda^{-1} \mathbb{M}_p \end{pmatrix} \cdot \begin{pmatrix} \vec{\mathcal{V}}^{k+1} \\ \vec{\mathcal{P}}^{k+1} \end{pmatrix} = \begin{pmatrix} \vec{f} \\ \vec{h} - \lambda^{-1} \mathbb{M}_p \cdot \vec{\mathcal{P}}^k \end{pmatrix} \quad (311)$$

where  $k$  is the iteration counter and  $\lambda$  is an artificial compressibility term which has the dimensions of dynamic viscosity. The choice of  $\lambda$  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.(311) converges to Eq.(310) and the velocity and pressure fields are solution of the unmodified system Eq.(310).

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. (311) 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:

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot \vec{\mathcal{P}}^{k+1}) \quad (312)$$

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k + \lambda \mathbb{M}_p^{-1} \cdot (\mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} - \vec{h})) \quad (313)$$

$$\mathbb{K} \cdot \vec{\mathcal{V}}^{k+1} + \lambda \mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T \cdot \vec{\mathcal{V}}^{k+1} = \vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k - \lambda \mathbb{M}_p^{-1} \vec{h}) \quad (314)$$

$$\underbrace{(\mathbb{K} + \lambda \mathbb{G} \cdot \mathbb{M}_p^{-1} \cdot \mathbb{G}^T)}_{\tilde{\mathbb{K}}} \cdot \vec{\mathcal{V}}^{k+1} = \underbrace{\vec{f} - \mathbb{G} \cdot (\vec{\mathcal{P}}^k - \lambda \mathbb{M}_p^{-1} \vec{h})}_{\vec{f}^{k+1}} \quad (315)$$

$$(316)$$

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  $\lambda >> \max_{\Omega} \eta$  then the matrix  $\tilde{\mathbb{K}}$  is ill-conditioned and an iterative solver must be used.

### 7.13.5 The GMRES approach

The Generalized Minimal Residual method [488] 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).

Resources: [177, p208] [487] [35]

finish GMRES algo description. not sure what to do, hard to explain, not easy to code.

## 7.14 The consistent boundary flux (CBF)

The Consistent Boundary Flux technique was devised to alleviate the problem of the accuracy of primary variables derivatives (mainly velocity and temperature) on boundaries, where basis function (nodal) derivatives do not exist. 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 [412] and later used in geodynamics [623]. It was finally implemented in the CitcomS code [625] 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 ??.

### 7.14.1 applied to the Stokes equation

We start from the strong form:

$$\nabla \cdot \boldsymbol{\sigma} = \mathbf{b}$$

and then write the weak form:

$$\int_{\Omega} N \nabla \cdot \boldsymbol{\sigma} dV = \int_{\Omega} N \mathbf{b} dV$$

where  $N$  is any test function. We then use the two equations:

$$\nabla \cdot (N \boldsymbol{\sigma}) = N \nabla \cdot \boldsymbol{\sigma} + \nabla N \cdot \boldsymbol{\sigma} \quad (\text{chain rule})$$

$$\int_{\Omega} (\nabla \cdot \mathbf{f}) dV = \int_{\Gamma} \mathbf{f} \cdot \mathbf{n} dS \quad (\text{divergence theorem})$$

Integrating the first equation over  $\Omega$  and using the second, we can write:

$$\int_{\Gamma} N \boldsymbol{\sigma} \cdot \mathbf{n} dS - \int_{\Omega} \nabla N \cdot \boldsymbol{\sigma} dV = \int_{\Omega} N \mathbf{b} dV$$

On  $\Gamma$ , the traction vector is given by  $\mathbf{t} = \boldsymbol{\sigma} \cdot \mathbf{n}$ :

$$\int_{\Gamma} N \mathbf{t} dS = \int_{\Omega} \nabla N \cdot \boldsymbol{\sigma} dV + \int_{\Omega} N \mathbf{b} dV$$

Considering the traction vector as an unknown living on the nodes on the boundary, we can expand (for  $Q_1$  elements)

$$t_x = \sum_{i=1}^2 t_{x|i} N_i \quad t_y = \sum_{i=1}^2 t_{y|i} N_i$$

on the boundary so that the left hand term yields a mass matrix  $M'$ . Finally, using our previous experience of discretising the weak form, we can write:

$$M' \cdot \mathcal{T} = -\mathbb{K}\mathcal{V} - \mathbb{G}\mathcal{P} + f$$

where  $\mathcal{T}$  is the vector of assembled tractions which we want to compute and  $\mathcal{V}$  and  $\mathcal{T}$  are the solutions of the Stokes problem. Note that the assembly only takes place on the elements along the boundary.

Note that the assembled mass matrix is tri-diagonal can be easily solved with a Conjugate Gradient method. With a trapezoidal integration rule (i.e. Gauss-Lobatto) the matrix can even be diagonalised and the resulting matrix is simply diagonal, which results in a very cheap solve [623].

### 7.14.2 applied to the heat equation

We start from the strong form of the heat transfer equation (without the source terms for simplicity):

$$\rho C_p \left( \frac{\partial T}{\partial t} + \mathbf{v} \cdot \nabla T \right) = \nabla \cdot k \nabla T$$

The weak form then writes:

$$\int_{\Omega} N \rho C_p \frac{\partial T}{\partial t} dV + \rho C_p \int_{\Omega} N \mathbf{v} \cdot \nabla T dV = \int_{\Omega} N \nabla \cdot k \nabla T dV$$

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 shape 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

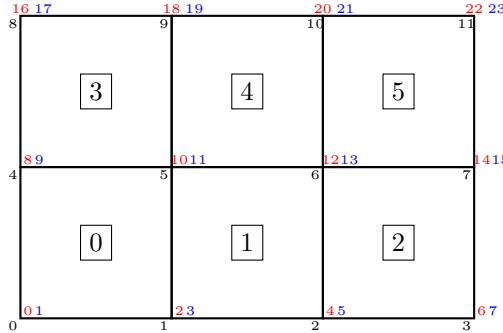
$$M' \cdot \mathcal{Q}_n = -M \cdot \frac{\partial \mathbf{T}}{\partial t} - K_a \cdot \mathbf{T} - K_d \cdot \mathbf{T}$$

where  $\mathcal{Q}_n$  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.

#### 7.14.3 implementation - Stokes equation

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, NfemV=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]
```

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 NfemTr is then the number of T in the `bc_fix` array. In our specific case, we have NfemTr=. This means that we need for each targeted dof to be able to find its identity/number between 0 and NfemTr-1. We therefore create the array `bc_nb` which is filled as follows:

[finish](#)

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

- A Prepare two arrays to store the matrix  $M_{cbf}$  and its right hand side  $rhs_{cbf}$
  - B Loop over all elements
  - C For each element touching a boundary, compute the residual vector  $R_{el} = -f_{el} + \mathbb{K}_{el}\mathcal{V}_{el} + \mathbb{G}_{el}\mathcal{P}_{el}$
  - D Loop over the four edges of the element using the connectivity array
  - E For each edge loop over the number of degrees of freedom (2 in 2D)
  - F 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)
rhs_cbf = np.zeros(NfemTr,np.float64)

for iel in range(0,nel):
 ... compute elemental residual ...
 # Boundary 0-1
 for i in range(0,ndofV):
 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] # I
 M_cbf[idofTr0,idofTr1]+=M_edge[0,1] # I
 M_cbf[idofTr1,idofTr0]+=M_edge[1,0] # I
 M_cbf[idofTr1,idofTr1]+=M_edge[1,1] # I

 # Boundary 1-2
 ...
 # Boundary 2-3
 ...
 # Boundary 3-0
 ...

```

## 7.15 The value of the timestep

The chosen time step  $\delta t$  used for time integration is chosen to comply with the Courant-Friedrichs-Lowy condition [17].

$$\delta t = C \min \left( \frac{h}{\max |\mathbf{v}|}, \frac{h^2}{\kappa} \right) \quad (317)$$

where  $h$  is a measure of the element size,  $\kappa = k/\rho C_p$  is the thermal diffusivity and  $C$  is the so-called CFL number chosen in  $[0, 1]$ .

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 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.

## 7.16 Mappings

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.

### 7.16.1 Linear mapping on a triangle

```

2
|\ s
| \ |_r
| \
3==1

```

Let us assume that the coordinates of the vertices are  $(x_1, y_1)$ ,  $(x_2, y_2)$ , and  $(x_3, y_3)$ . The coordinates inside the reference element are  $(r, s)$ . We then simply have the following relationship, i.e. any point of the reference element can be mapped to the physical triangle as follows:

$$x = rx_1 + sx_2 + (1 - r - s)x_3 \quad (318)$$

$$y = ry_1 + sy_2 + (1 - r - s)y_3 \quad (319)$$

There is also an inverse map, which is easily computed:

$$r = \frac{(y_2 - y_3)(x - x_3) - (x_2 - x_3)(y - y_3)}{(x_1 - x_3)(y_2 - y_3) - (y_1 - y_3)(x_2 - x_3)} \quad (320)$$

$$s = \frac{-(y_1 - y_3)(x - x_3) + (x_1 - x_3)(y - y_3)}{(x_1 - x_3)(y_2 - y_3) - (y_1 - y_3)(x_2 - x_3)} \quad (321)$$

**Remark.** The denominator will not vanish, because it is a multiple of the area of the triangle.

### 7.16.2 Linear mapping on a quadrilateral

```

4====3
| | s
| | |_r
|
1====2

```

The coordinates of the vertices are  $(x_1, y_1)$ ,  $(x_2, y_2)$ ,  $(x_3, y_3)$  and  $(x_4, y_4)$ . The coordinates inside the reference element are  $(r, s)$ . We then simply have the following relationship, i.e. any point of the reference element can be mapped to the physical quadrilateral as follows:

$$x = N_1(r, s)x_1 + N_2(r, s)x_2 + N_3(r, s)x_3 + N_4(r, s)x_4 \quad (322)$$

$$y = N_1(r, s)y_1 + N_2(r, s)y_2 + N_3(r, s)y_3 + N_4(r, s)y_4 \quad (323)$$

where the shape functions  $N_i(r, s)$  are defined in section ???. There is also an inverse map, which is not so easily computed (see section 7.19).

However, if the quadrilateral in the  $(x, y)$  space is a rectangle of size  $(h_x, h_y)$ , the inverse mapping is trivial:

$$r = \frac{x - x_1}{x_2 - x_1} \quad (324)$$

$$s = \frac{y - y_1}{y_4 - y_1} \quad (325)$$

Also in this case the shape functions can easily be written as functions of  $(x, y)$ :

$$\begin{aligned}N_1(x, y) &= \left( \frac{x_3 - x}{h_x} \right) \left( \frac{y_3 - y}{h_y} \right) \\N_2(x, y) &= \left( \frac{x - x_1}{h_x} \right) \left( \frac{y_3 - y}{h_y} \right) \\N_3(x, y) &= \left( \frac{x - x_1}{h_x} \right) \left( \frac{y - y_1}{h_y} \right) \\N_4(x, y) &= \left( \frac{x_3 - x}{h_x} \right) \left( \frac{y - y_1}{h_y} \right)\end{aligned}$$

## 7.17 Exporting data to vtk format

This format seems to be the universally accepted format for 2D and 3D visualisation in Computational Geodynamics. Such files can be opened with free softwares such as Paraview<sup>16</sup>, MayaVi<sup>17</sup> or Visit<sup>18</sup>.

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 format<sup>19</sup> but it delivers the information in a convoluted way. I therefore describe hereafter how **fieldstone** builds the vtk files.

I hereunder show vtk file corresponding to the 3x2 grid presented earlier 7.8. 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.vtu* 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
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
```

<sup>16</sup><https://www.paraview.org/>

<sup>17</sup><https://docs.enthought.com/mayavi/mayavi/>

<sup>18</sup><https://wci.llnl.gov/simulation/computer-codes/visit/>

<sup>19</sup><https://www.vtk.org/wp-content/uploads/2015/04/file-formats.pdf>

```

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
</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 said connectivity is 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 simply is the type of element as given in the vtk format document (9 corresponds to a generic quadrilateral with an internal numbering consistent with ours).

```

<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* file 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.

## 7.18 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 ODEs<sup>20</sup>. These methods are well documented in any numerical analysis textbook and the reader is referred to [241, 320]. Any Runge-Kutta method is uniquely identified by its Butcher tableau (REF?) which contains all necessary coefficients to build the algorithm.

The simplest RungeKutta method is the (forward) Euler method. Its tableau is:

0				
	1			

The standard second-order RK method method (also called midpoint method) is:

0				
1/2	1/2			
	0	1		

Another second-order RK method, called Heun's method<sup>21</sup> 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	1/6	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

The following method is called the Runge-Kutta-Fehlberg method and is commonly abbreviated RKF45<sup>22</sup>. Its Butcher tableau is as follows:

0					
1/4	1/4				
3/8	3/32	9/32			
12/13	1932/2197	-7200/2197	7296/2197		
1	439/216	-8	3680/513	-845/4104	
1/2	-8/27	2	-3544/2565	1859/4104	-11/40
	16/135	0	6656/12825	28561/56430	-9/50
	25/216	0	1408/2565	2197/4104	-1/5
					0

The first row of coefficients at the bottom of the table gives the fifth-order accurate method, and the second row gives the fourth-order accurate method.

### 7.18.1 Using RK methods to advect particles/markers

In the context of geodynamical modelling, one is usually confronted to the following problem: now that I have a velocity field on my FE mesh, how can I use it to advect the Lagrangian markers?

<sup>20</sup>[https://en.wikipedia.org/wiki/Runge-Kutta\\_methods](https://en.wikipedia.org/wiki/Runge-Kutta_methods)

<sup>21</sup>[https://en.wikipedia.org/wiki/Heun%27s\\_method](https://en.wikipedia.org/wiki/Heun%27s_method)

<sup>22</sup>[https://en.wikipedia.org/wiki/Runge-Kutta-Fehlberg\\_method](https://en.wikipedia.org/wiki/Runge-Kutta-Fehlberg_method)

Runge-Kutta methods are used to this effect but only their spatial component is used: the velocity solution is not recomputed at the intermediate fractional timesteps, i.e. only the coefficients of the right hand side of the tableaus is used.

The RK1 method is simple. Carry out a loop over markers and

1. interpolate velocity  $\vec{v}_m$  onto each marker  $m$
2. compute new position as follows:  $\vec{r}_m(t + \delta t) = \vec{r}_m(t) + \vec{v}_m \delta t$

The RK2 method is also simple but requires a bit more work. 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. 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 intermediate position as follows:  $\vec{r}_m^{(2)}(t + \delta t) = \vec{r}_m(t) + (2\vec{v}_m^{(1)} - \vec{v}_m) \delta t / 2$
5. compute velocity  $\vec{v}_m^{(2)}$  at position  $\vec{r}_m^{(2)}$
6. compute new position:  $\vec{r}_m(t + \delta t) = \vec{r}_m(t) + (\vec{v}_m + 4\vec{v}_m^{(1)} + \vec{v}_m^{(2)}) \delta t / 6$

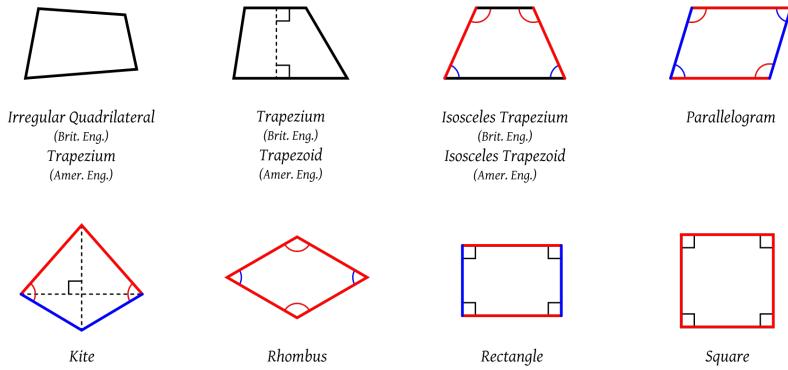
## 7.19 Am I in or not? - finding reduced coordinates

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?"

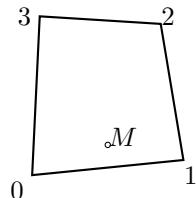
One typical occurrence of such a problem is linked to the use of the Particle-In-Cell technique: particles are advected and move through the mesh, and need to be localised at every time step. This question could arise in the context of a benchmark where certain quantities need to be measured at specific locations inside the domain.

### 7.19.1 Two-dimensional space

We shall first focus on quadrilaterals. There are many kinds of quadrilaterals as shown hereunder:



I wish to arrive at a single algorithm which is applicable to all quadrilaterals and therefore choose an irregular quadrilateral. For simplicity, let us consider a  $Q_1$  element, with a single node at each corner.



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.

### 7.19.2 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

$$\sum_{i=1}^4 N_i(r_M, s, t) x_i = x_M \quad \sum_{i=1}^4 N_i(r_M, s, t) y_i = y_M \quad \sum_{i=1}^4 N_i(r_M, s, t) z_i = z_M$$

This can be cast as a system of three equations and three unknowns. Unfortunately, each shape function  $N_i$  contains a term  $rst$  (as well as  $rs$ ,  $rt$ , and  $st$ ) so that it is not a linear system and standard techniques are not applicable. We must then use an iterative technique: the algorithm starts with a guess for values  $r, s, t$  and improves on their value iteration after iteration.

The classical way of solving nonlinear systems of equations is Newton's method. We can rewrite the equations above as  $\mathbf{F}(r, s, t) = 0$ :

$$\begin{aligned} \sum_{i=1}^8 N_i(r, s, t)x_i - x_M &= 0 \\ \sum_{i=1}^8 N_i(r, s, t)y_i - y_M &= 0 \\ \sum_{i=1}^8 N_i(r, s, t)z_i - z_M &= 0 \end{aligned} \quad (326)$$

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:

$$\begin{aligned} J_F(r_k, s_k, t_k) &= \begin{pmatrix} \frac{\partial F_r}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_r}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_r}{\partial t}(r_k, s_k, t_k) \\ \frac{\partial F_s}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_s}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_s}{\partial t}(r_k, s_k, t_k) \\ \frac{\partial F_t}{\partial r}(r_k, s_k, t_k) & \frac{\partial F_t}{\partial s}(r_k, s_k, t_k) & \frac{\partial F_t}{\partial t}(r_k, s_k, t_k) \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^8 \frac{\partial N_i}{\partial r}(r_k, s_k, t_k)x_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial s}(r_k, s_k, t_k)x_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial t}(r_k, s_k, t_k)x_i \\ \sum_{i=1}^8 \frac{\partial N_i}{\partial r}(r_k, s_k, t_k)y_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial s}(r_k, s_k, t_k)y_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial t}(r_k, s_k, t_k)y_i \\ \sum_{i=1}^8 \frac{\partial N_i}{\partial r}(r_k, s_k, t_k)z_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial s}(r_k, s_k, t_k)z_i & \sum_{i=1}^8 \frac{\partial N_i}{\partial t}(r_k, s_k, t_k)z_i \end{pmatrix} \end{aligned}$$

In practice, we solve the following system:

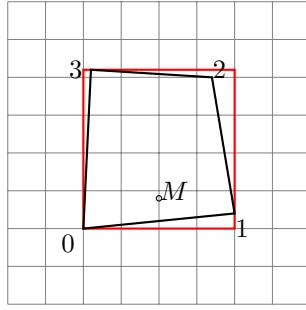
$$J_F(r_k, s_k, t_k) \left[ \begin{pmatrix} r_{k+1} \\ s_{k+1} \\ t_{k+1} \end{pmatrix} - \begin{pmatrix} r_k \\ s_k \\ t_k \end{pmatrix} \right] = - \begin{pmatrix} F_r(r_k, s_k, t_k) \\ F_s(r_k, s_k, t_k) \\ F_t(r_k, s_k, t_k) \end{pmatrix}$$

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  $3 \times 3$  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.

write about case when element is rectangle/cuboid

## 7.20 Error measurements and convergence rates

What follows is written in the case of a two-dimensional model. Generalisation to 3D is trivial. What follows is mostly borrowed from [541].

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  $\Psi$ , the  $L_1$  and  $L_2$  norms are computed as

$$\|\Psi\|_1 = \int_V |\Psi| dV \quad \|\Psi\|_2 = \sqrt{\int_V \Psi^2 dV} \quad (327)$$

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:

$$\|\vec{k}\|_1 = \int_V (|k_x| + |k_y|) dV \quad \|\vec{k}\|_2 = \sqrt{\int_V (k_x^2 + k_y^2) dV} \quad (328)$$

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

$$e_p^h|_1 = \sum_{i=1}^{n_e} \sum_{q=1}^{n_q} |e_p^h(\vec{r}_q)| w_q |J_q| \quad e_p^h|_2 = \sqrt{\sum_{i=1}^{n_e} \sum_{q=1}^{n_q} |e_p^h(\vec{r}_q)|^2 w_q |J_q|} \quad (329)$$

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

$$e_{\vec{v}}^h|_1 = \sum_{i=1}^{n_e} \sum_{q=1}^{n_q} [|e_u^h(\vec{r}_q)| + |e_v^h(\vec{r}_q)|] w_q |J_q| \quad e_{\vec{v}}^h|_2 = \sqrt{\sum_{i=1}^{n_e} \sum_{q=1}^{n_q} [|e_u^h(\vec{r}_q)|^2 + |e_v^h(\vec{r}_q)|^2] w_q |J_q|} \quad (330)$$

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 [165, 330, 308]. This norm is expressed as follows for a function  $f$  such that  $f, |\nabla f| \in L^2(\Omega)$ <sup>23</sup>

$$\|f\|_{H^1} = \left( \int_{\Omega} (|f|^2 + |\nabla f|^2) d\Omega \right)^{1/2} \quad (331)$$

We then have

$$e_{\vec{v}}^h|_{H^1} = \|\vec{v}^h - \vec{v}\|_{H^1} = \sqrt{\sum_{i=1}^d \int_{\Omega} [(v_i^h - v_i)^2 + \vec{\nabla}(v_i^h - v_i) \cdot \vec{\nabla}(v_i^h - v_i)] d\Omega} \quad (332)$$

where  $d$  is the number of dimensions. Note that sometimes the following semi-norm is used [164, 65]:

$$e_{\vec{v}}^h|_{H^1} = \|\vec{v}^h - \vec{v}\|_{H^1} = \sqrt{\sum_{i=1}^d \int_{\Omega} [\vec{\nabla}(v_i^h - v_i) \cdot \vec{\nabla}(v_i^h - v_i)] d\Omega} \quad (333)$$

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:

$$e_{\vec{v}}^h|_1 = Ch^{rvL_1} \quad e_{\vec{v}}^h|_2 = Ch^{rvL_2} \quad e_{\vec{v}}^h|_{H^1} = Ch^{rvH^1} \quad (334)$$

---

<sup>23</sup>[https://en.wikipedia.org/wiki/Sobolev\\_space](https://en.wikipedia.org/wiki/Sobolev_space)

$$e_p^h|_1 = Ch^{rpL_1} \quad e_p^h|_2 = Ch^{rpL_2} \quad (335)$$

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 [164], when finite element solutions converge at the same rates as the interpolants we say that the method is optimal, i.e.:

$$e_v^h|_{L_2} = \mathcal{O}(h^3) \quad e_v^h|_{H^1} = \mathcal{O}(h^2) \quad e_p^h|_{L_2} = \mathcal{O}(h^2) \quad (336)$$

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.

### 7.20.1 About extrapolation

*Section contributed by W. Bangerth and part of Thieulot & Bangerth [in prep.]*

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  $\vec{u}_h, p_h$  and the corresponding value  $X$  computed from the exact solution  $\vec{u}, p$ . Examples of such quantities  $X$  are the root mean square velocity  $v_{rms}$ , 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

$$e_h = |X_h - X| \approx Ch^r. \quad (337)$$

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 (337) we also have

$$e_{h/2} = |X_{h/2} - X| \approx C \left( \frac{h}{2} \right)^r, \quad (338)$$

$$e_{h/4} = |X_{h/4} - X| \approx C \left( \frac{h}{4} \right)^r. \quad (339)$$

Taking ratios of equations (337)–(339), 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  $r$ :

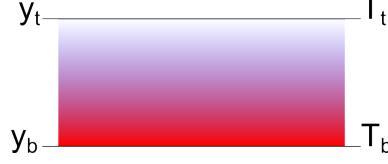
$$X^* = \frac{X_h X_{h/2} - X_{h/2}^2}{X_h - 2X_{h/2} + X_{h/4}}, \quad r = \log_2 \frac{X_{h/2} - X^*}{X_{h/4} - X^*}.$$

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 (337) are less visible on finer meshes.

## 7.21 The initial temperature field

### 7.21.1 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

$$\rho C_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \vec{\nabla} \cdot (k \vec{\nabla} T) + \rho H$$

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

$$\Delta T = 0$$

Assuming the layer to be parallel to the  $x$ -axis, the temperature is

$$T(x, y) = T(y) = \alpha T + \beta$$

In order to specify the constants  $\alpha$  and  $\beta$ , 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

$$T(y) = \boxed{\frac{T_t - T_b}{y_t - y_b}(y - y_b) + T_b}$$

If now the heat production coefficient is not zero, the differential equation reads

$$k \Delta T + H = 0$$

The temperature field is then expected to be of the form

$$T(y) = -\frac{H}{2k}y^2 + \alpha y + \beta$$

Supplied again with the same boundary conditions, this leads to

$$\beta = T_b + \frac{H}{2k}y_b^2 - \alpha y_b$$

ie,

$$T(y) = -\frac{H}{2k}(y^2 - y_b^2) + \alpha(y - y_b) + T_b$$

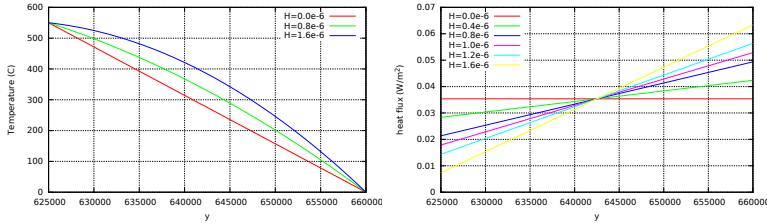
and finally

$$\alpha = \frac{T_t - T_b}{y_t - y_b} + \frac{H}{2k}(y_b + y_t)$$

or,

$$T(y) = -\frac{H}{2k}(y^2 - y_b^2) + \left( \frac{T_t - T_b}{y_t - y_b} + \frac{H}{2k}(y_b + y_t) \right)(y - y_b) + T_b$$

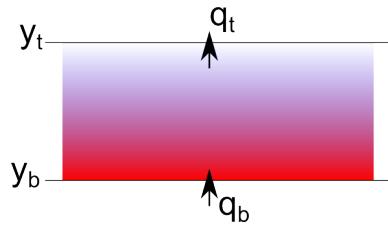
Taking  $H = 0$  in this equation obviously yields the temperature field obtained previously. Taking  $k = 2.25$ ,  $T_t = 0C$ ,  $T_b = 550C$ ,  $y_t = 660km$ ,  $y_b = 630km$  yields the following temperature profiles and heat fluxes when the heat production  $H$  varies:



Looking at the values at the top, which are somewhat estimated to be about  $55 - 65 \text{ mW/m}^2$  [329, 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  $0.03 \text{ W/m}^2$  which is somewhat problematic since the heat flux at the Moho is reported to be somewhere between 10 and  $20 \text{ mW/m}^2$  in [329, table 7.1].

### 7.21.2 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

$$k\Delta T + H = 0$$

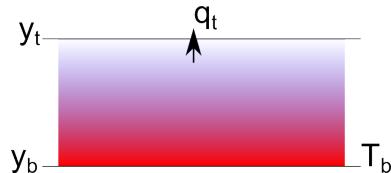
but only integrate it once:

$$k \frac{dT}{dy} + Hy + \alpha = 0$$

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](#)

### 7.21.3 Single layer with imposed heat flux and temperature b.c.



[to finish](#)

### 7.21.4 Half cooling space

### 7.21.5 Plate model

### 7.21.6 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 [409] derived such approximate initial field from the steady-state energy equation in two dimensions:

$$\rho C_p \vec{v} \cdot \vec{\nabla} T = k \vec{\nabla}^2 T$$

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 [409]):

$$\frac{\partial^2 T'}{\partial x'^2} - 2R \frac{\partial T'}{\partial x'} + \frac{\partial^2 T'}{\partial z'^2} = 0$$

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

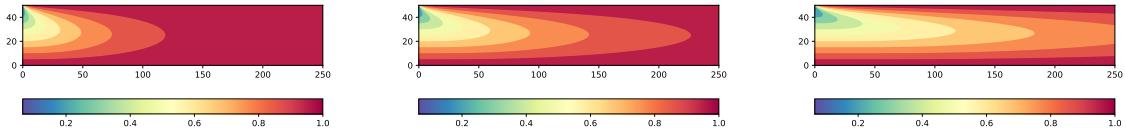
$$T'(x', z') = 1 + \sum_n C_n \exp \left[ \left( R - (R^2 + n^2 \pi^2)^{1/2} \right) x' \right] \sin(n\pi z')$$

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

$$T'(x', z') = 1 + 2 \sum_n \frac{(-1)^n}{n\pi} \exp \left[ \left( R - (R^2 + n^2 \pi^2)^{1/2} \right) x' \right] \sin(n\pi z')$$

(340)

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`.

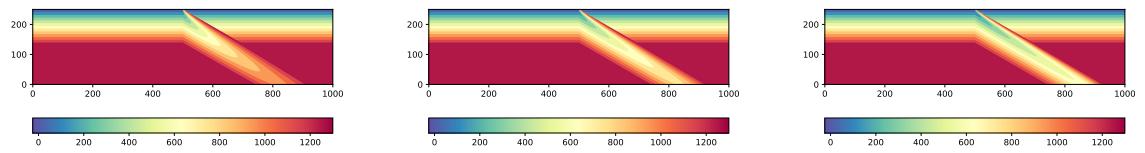


Left to right: Dimensionless temperature  $T'$  in a  $250\text{km} \times 50\text{km}$  slab for  $v_x = 0.5, 1, 2\text{cm/year}$

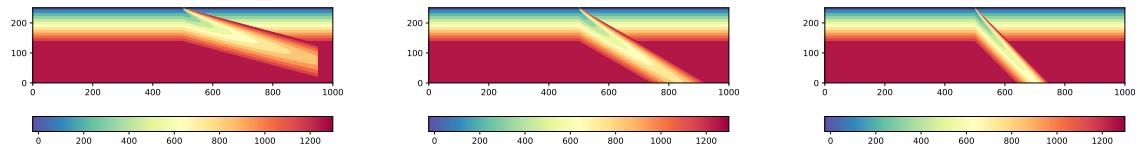
We logically recover the fact that the slower the slab penetrates the mantle the more temperature diffusion dominates over temperature advection. For  $v = 0.5\text{cm/year}$  we see that that the slab assumes a constant temperature  $T' = 1$  at all depths  $0 \leq z' \leq 1$  for  $x' \geq 125\text{km}$ .

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  $1300^\circ$ . 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  $\phi = 15, 30, 45^\circ$ .

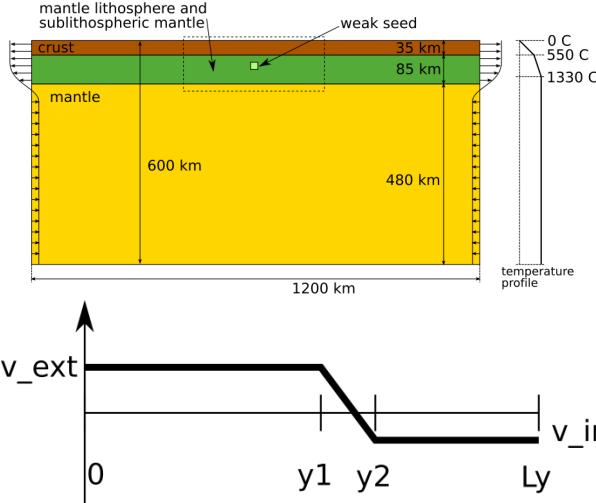
## 7.22 Kinematic boundary conditions

Boundary conditions come in two basic flavors: essential and natural.

- Essential bcs directly affect DOFs, and are imposed on the FEM matrix.
- Natural bcs do not directly affect DOFs and are imposed on the right-hand side vector.

### 7.22.1 In-out flux boundary conditions for lithospheric models

C. Thieulot / Physics of the Earth and Planetary Interiors 188 (2011) 47–68



The velocity on the side is given by

$$\begin{aligned} u(y) &= v_{ext} \quad y < L_y \\ u(y) &= \frac{v_{in} - v_{ext}}{y_2 - y_1}(y - y_1) + v_{ext} \quad y_1 < y < y_2 \\ u(y) &= v_{in} \quad y > y_2 \end{aligned}$$

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  $y_2$ .

$$\begin{aligned} \Phi &= \int_0^{y_1} u(y) dy + \int_{y_1}^{y_2} u(y) dy + \int_{y_2}^{L_y} u(y) dy \\ &= v_{ext}y_1 + \frac{1}{2}(v_{in} + v_{ext})(y_2 - y_1) + (L_y - y_2)v_{in} \\ &= v_{ext}[y_1 + \frac{1}{2}(y_2 - y_1)] + v_{in}[\frac{1}{2}(y_2 - y_1) + (L_y - y_2)] \\ &= v_{ext}\frac{1}{2}(y_1 + y_2) + v_{in}[L_y - \frac{1}{2}(y_1 + y_2)] \end{aligned}$$

and finally

$$v_{ext} = -v_{in} \frac{L_y - \frac{1}{2}(y_1 + y_2)}{\frac{1}{2}(y_1 + y_2)}$$

## 7.23 Computing gradients - the recovery process

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  $Q_1$  representation of the field  $\vec{\nabla}f^h$ . Since the derivative of the shape function does not exist on the nodes we need to design an algorithm do do so. This problem is well known and has been investigated [?]

refs!

- . The main standard techniques are listed hereafter.

### 7.23.1 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  $\phi$  by all the shape functions  $N_i$  and since we have  $g^h = \sum_j N_i g_i$  we then obtain

$$\sum_j \int N_i N_j d\Omega g_i = \int N_i \vec{\nabla} f^h \, d\Omega$$

or,

$$\mathbb{M} \cdot \vec{\mathcal{G}} = \vec{f}$$

### 7.23.2 Local recovery - centroid average over patch

### 7.23.3 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  $(\vec{\nabla} f^h)_e(\vec{r}_j)$  is the gradient of  $f$  as obtained with the shape functions inside element  $e$  and computed at location  $\vec{r}_j$ .

### 7.23.4 Local recovery - least squares over patch

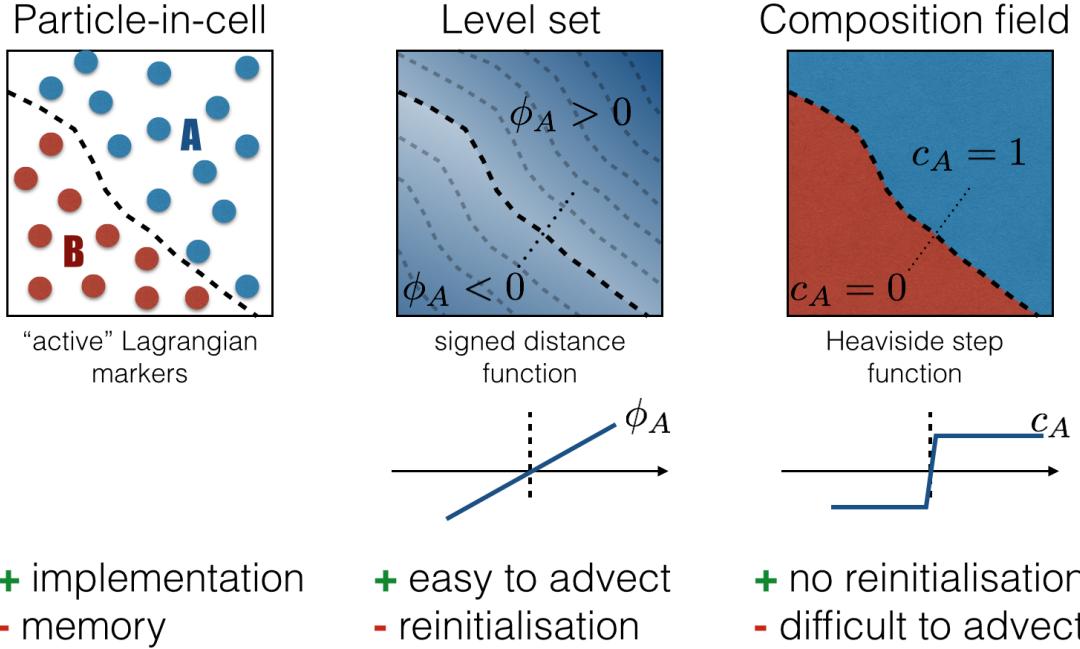
### 7.23.5 Link to pressure smoothing

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 6.3, 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)

tie to fieldstone 12

## 7.24 Tracking materials and/or interfaces

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 ...

### 7.24.1 The Particle-in-cell technique

**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 [293]. For more information on MAC see the review paper by McKee et al [408].

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 [454] in non-parallel codes such as SOPALE [218], I2VIS [250] or CITCOM [410] (Appendix B). It has been implemented in ASPECT [?] and the inherent load balancing issues arising from the parallel implementation as well as from the use of Adaptive Mesh Refinement are discussed.

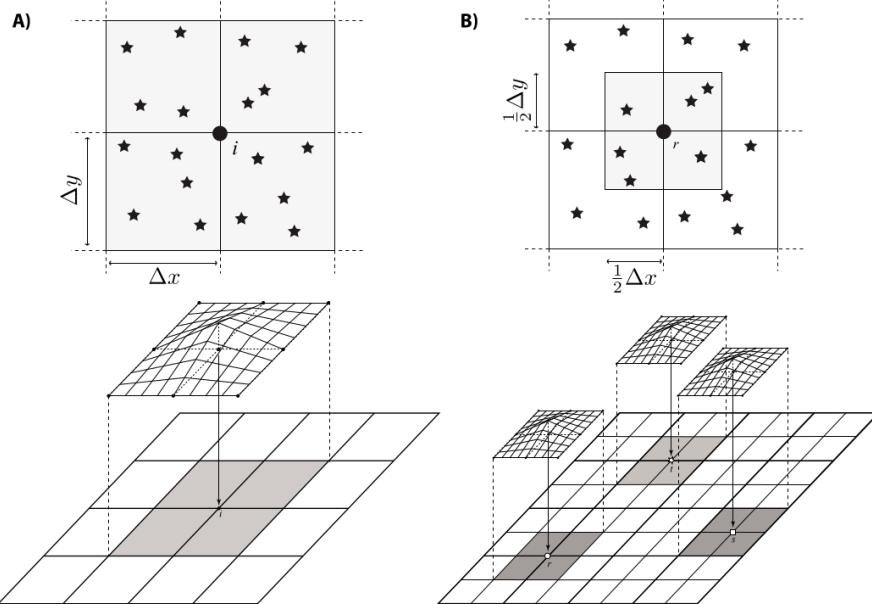
The basic methodology goes as follows:

1. distribute particles in the domain
2. assign a material identity (and/or any other quantity) to each of them
3. project particle quantities of the Eulerian 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 fieldstone 13 (Section 20).

**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 quantity 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 [171] as illustrated in the following figure:



Taken from [171]. 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, only the particles inside the element will contribute to the average values assigned to the quadrature points. 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. In this case the FDM approach applies.

Finally, in both FDM and FEM bi/trilinear shape functions are used for the interpolation as they can be interpreted as weighing functions. Higher order shape functions could also be used but the standard

$Q_2$  shape functions (Section 4.5) are 2-nd order polynomials which can take negative values (as opposed to the  $Q_1$  shape 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. Q1 projection PUCKETT

it would be nice to have a Q1 and Q2 drawing of a 1D element and show that indeed negative values arise

Assuming that we have established a list of particles, all tracking a field  $f(\vec{r})$  and that each particle has an associated weight  $N_i$  (function of the location where the average is to be computed or not), we must now compute their average value  $\langle f \rangle$ . The simplest approach which comes to mind is the (weighed) arithmetic mean (*am*):

$$\langle f \rangle_{am} = \frac{\sum_{i=1}^n N_i f_i}{\sum_{i=1}^n N_i}$$

In the case where  $f$  is the (mass) density  $\rho$ , it is indeed what should be used. However, turning now to viscosity  $\eta$ , 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  $10^{18}\text{Pa s}$  and  $10^{25}\text{Pa s}$ . As explained in [497] 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 (assuming  $N_i = 1$  for simplicity) 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:

$$\langle f \rangle_{gm} = \left( \prod_i f_i^{N_i} \right)^{1/\sum_i N_i} \quad \text{or,} \quad \log_{10} \langle f \rangle_{gm} = \frac{\sum_i N_i \log_{10} f_i}{\sum_i N_i}$$

and

$$\langle f \rangle_{hm} = \left( \frac{\sum_{i=1}^n N_i \frac{1}{f_i}}{\sum_i N_i} \right)^{-1} \quad \text{or,} \quad \frac{1}{\langle f \rangle_{hm}} = \frac{\sum_{i=1}^n N_i \frac{1}{f_i}}{\sum_i N_i}$$

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:

$$\langle f \rangle_{am} \geq \langle f \rangle_{gm} \geq \langle f \rangle_{hm} \tag{341}$$

When all  $f_i$  are equal to  $f_0$  their computed average should also be equal to  $f_0$ . As a consequence the weights  $N_i$  should fulfill the condition  $\sum_{i=1}^n N_i = 1$ . If all weights are equal, then  $N_i = 1/n$  and the averagings become:

$$\langle f \rangle_{am} = \frac{1}{n} \sum_{i=1}^n f_i \quad \langle f \rangle_{gm} = \prod_i f_i^{1/n} \quad \langle f \rangle_{hm} = \left( \frac{1}{n} \sum_i \frac{1}{f_i} \right)^{-1} \tag{342}$$

There are many papers which have looked at particle averagings and projections. I will for now simply point to the following ones: [497] [162] [171] [416] [455] [541] [224].

write more about particle averaging and projection

### Interpolation of the velocity onto particles .

Once the particle  $i$  has been localised inside a given element (Section 7.19) and its reduced coordinates  $(r, s, t)$  determined, the velocity at this location can be computed through the shape 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 that<sup>24</sup>  $\vec{\nabla} \cdot \vec{v} = 0$  (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 [587].

**Moving the particles** This is discussed in the context of the Runge-Kutta Methods, see Section 7.18.1.

#### 7.24.2 The level set function technique

This method was developed in the 80's by Stanley Osher and James Sethian []

The Level-set Method (LSM), as it is commonly used in Computational Fluid Dynamics – and especially in Computational Geodynamics – represents a close curve  $\Gamma$  (say, in our case, the interface between two fluids or layers) by means of a function  $\phi$  (called the level-set function, or LSF).  $\Gamma$  is then the zero level-set of  $\phi$ :

$$\Gamma = \{(x, y) \mid \phi(x, y) = 0\} \quad (343)$$

The convention is that  $\phi > 0$  inside the region delimited by  $\Gamma$  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  $\Gamma$  moves with a velocity  $\vec{v}$ , then it satisfies the following equation:

$$\frac{\partial \phi}{\partial t} + \vec{v} \cdot \vec{\nabla} \phi = 0 \quad (344)$$

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.

As explained in [303], the level-set function  $\phi$  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 [438, 489].

The level set method has not often been used in the geodynamics community with some notable exceptions [71, 72, 289, 281, 637, 290, 526, 525, 303] An overview of the method and applications can be found in [437].

#### 7.24.3 The field/composition technique

This is the approach taken by the ASPECT developers [358, 299]. Each material  $i$  is represented by a compositional field  $c_i$ , which takes values between 0 and 1. 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 Heavyside 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 ( $c_i > 1$ ), 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 7.3) and the entropy viscosity method add small amounts of diffusion to dampen the under- and overshoots. This means that at a given point two or more compositions may

---

<sup>24</sup>for incompressible flows, of course

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).

write about DG approach

#### 7.24.4 Hybrid methods

In Braun et al. [80] 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 B) 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 [2]. Note that another very similar approach is used in [489].

## 7.25 Static condensation

The idea behind 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 ??). In the following,  $\vec{V}^*$  corresponds to the list of such dofs inside an element. The discretised Stokes equations on any element looks like:

$$\begin{pmatrix} \mathbb{K} & L & \mathbb{G} \\ L^T & \mathbb{K}^* & H \\ \mathbb{G}^T & H^T & 0 \end{pmatrix}_e \begin{pmatrix} \vec{\mathcal{V}} \\ \vec{V}^* \\ \vec{\mathcal{P}} \end{pmatrix}_e = \begin{pmatrix} \vec{f} \\ \vec{f}^* \\ \vec{h} \end{pmatrix}_e \quad (345)$$

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  $\mathbb{K}^*$  is diagonal.

This can also be re-written in non-matrix form:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + L \cdot \vec{V}^* + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (346)$$

$$L^T V + K^* \cdot \vec{V}^* + H \cdot \vec{\mathcal{P}} = \vec{f}^* \quad (347)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} + H^T \vec{V}^* = \vec{h} \quad (348)$$

The  $V^*$  in the second equation can be isolated:

$$\vec{V}^* = \mathbb{K}^{-*} \cdot (\vec{f}^* - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}})$$

and inserted in the first and third equations:

$$\mathbb{K} \cdot \vec{\mathcal{V}} + L \left[ \mathbb{K}^{-*} (\vec{f}^* - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}}) \right] + \mathbb{G} \cdot \vec{\mathcal{P}} = \vec{f} \quad (349)$$

$$\mathbb{G}^T \cdot \vec{\mathcal{V}} + H^T \left[ \mathbb{K}^{-*} (\vec{f}^* - L^T \cdot \vec{\mathcal{V}} - H \cdot \vec{\mathcal{P}}) \right] = \vec{h} \quad (350)$$

or,

$$(\mathbb{K} - L \cdot \mathbb{K}^{-*} \cdot L^T) \cdot \vec{\mathcal{V}} + (G - L \cdot \mathbb{K}^{-*} \cdot H) \cdot \vec{\mathcal{P}} = \vec{f} - L \cdot \mathbb{K}^{-*} \cdot \vec{f}^* \quad (351)$$

$$(G^T - H^T \cdot \mathbb{K}^{-*} \cdot L^T) \cdot \vec{\mathcal{V}} - (H^T \cdot \mathbb{K}^{-*} \cdot H) \cdot \vec{\mathcal{P}} = \vec{h} - H^T \cdot \mathbb{K}^{-*} \cdot \vec{f}^* \quad (352)$$

i.e.

$$\underline{\mathbb{K}} \cdot \vec{\mathcal{V}} + \underline{\mathbb{G}} \cdot \vec{\mathcal{P}} = \underline{\vec{f}} \quad (353)$$

$$\underline{\mathbb{G}}^T \cdot \vec{\mathcal{V}} - \underline{\mathbb{C}} \cdot \vec{\mathcal{P}} = \underline{\vec{h}} \quad (354)$$

with

$$\underline{\mathbb{K}} = K - L \cdot \mathbb{K}^{-*} \cdot L^T \quad (355)$$

$$\underline{\mathbb{G}} = G - L \cdot \mathbb{K}^{-*} \cdot H \quad (356)$$

$$\underline{\mathbb{C}} = H^T \cdot \mathbb{K}^{-*} \cdot H \quad (357)$$

$$\underline{\vec{f}} = \vec{f} - L \cdot \mathbb{K}^{-*} \cdot \vec{f}^* \quad (358)$$

$$\underline{\vec{h}} = \vec{h} - H^T \cdot \mathbb{K}^{-*} \cdot \vec{f}^* \quad (359)$$

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. However, it is then important to note that static condensation introduces a pressure-pressure term which was not there in the original formulation.

## 7.26 Measuring incompressibility

The velocity divergence error integrated over the whole element is given by

$$e_{div} = \int_{\Omega} (\vec{\nabla} \cdot \vec{v}^h - \underbrace{\vec{\nabla} \cdot \vec{v}}_{=0}) d\Omega = \int_{\Omega} (\vec{\nabla} \cdot \vec{v}^h) d\Omega \quad (360)$$

where  $\Gamma_e$  is the boundary of element  $e$  and  $\vec{n}$  is the unit outward normal of  $\Gamma_e$ .

Furthermore, one can show that [164]:

$$e_{div} = \int_{\Gamma_e} \vec{v}^h \cdot \vec{n} d\Gamma$$

The reason is as follows and is called the divergence theorem<sup>25</sup>: 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 [277].

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.

[implement and report](#)

---

<sup>25</sup>[https://en.wikipedia.org/wiki/Divergence\\_theorem](https://en.wikipedia.org/wiki/Divergence_theorem)

## 7.27 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.*

## 7.28 Removing rotational nullspace

July 10, 2019 - C.T.

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 7.12), the solution is simple:

1. fix the tangential velocity at *one* node on a boundary, and solve the system (the nullspace has been removed)
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 order to remove the angular momentum, we search for a rotation vector  $\vec{\omega}$  such that

$$\int_{\Omega} \rho[\vec{r} \times (\vec{v} - \vec{\omega} \times \vec{r})] dV = \vec{0} \quad (361)$$

The angular momentum of a rigid body can be obtained from the sum of the angular momentums of the particles forming the body<sup>26</sup>:

$$\vec{H} = \sum_i \vec{L}_i \quad (362)$$

$$= \sum_i \vec{r}_i \times m_i \vec{v}_i \quad (363)$$

$$= \sum_i \vec{r}_i \times m_i (\vec{\omega}_i \times \vec{r}_i) \quad (364)$$

$$= \sum_i m_i \begin{pmatrix} \sum_i m_i(y_i^2 + z_i^2) & -\sum_i m_i x_i y_i & -\sum_i m_i x_i z_i \\ -\sum_i m_i x_i y_i & \sum_i m_i(x_i^2 + z_i^2) & -\sum_i m_i y_i z_i \\ -\sum_i m_i x_i z_i & -\sum_i m_i y_i z_i & \sum_i m_i(x_i^2 + y_i^2) \end{pmatrix} \cdot \begin{pmatrix} \omega_x \\ \omega_y \\ \omega_z \end{pmatrix} \quad (365)$$

In the continuum limit, we have:

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \vec{r} \times \vec{v} dV \quad (366)$$

and the  $3 \times 3$  moment of inertia tensor  $\mathbf{I}$  (also called inertia tensor) is given by<sup>27</sup>

$$\mathbf{I} = \int_{\Omega} \rho(\vec{r}) [\vec{r} \cdot \vec{r} \mathbf{1} - \vec{r} \times \vec{r}] dV \quad (367)$$

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 [625, eq. 26]:

$$\vec{v}_{new} = \vec{v}_{old} - \vec{\omega} \times \vec{r} \quad (368)$$

**Angular velocity approach** The angular velocity<sup>28</sup> vector is given by  $\vec{\omega} = \frac{\vec{r} \times \vec{v}}{r^2}$  so that the volume-averaged angular velocity of the cylindrical shell is:

$$\vec{\omega} = \frac{1}{|\Omega|} \int_{\Omega} \frac{\vec{r} \times \vec{v}}{r^2} dV \quad (369)$$

<sup>26</sup><http://www.kwon3d.com/theory/moi/iten.html>

<sup>27</sup>[https://en.wikipedia.org/wiki/Moment\\_of\\_inertia](https://en.wikipedia.org/wiki/Moment_of_inertia)

<sup>28</sup>[https://en.wikipedia.org/wiki/angular\\_velocity](https://en.wikipedia.org/wiki/angular_velocity)

### 7.28.1 Three dimensions

The angular momentum vector is given by:

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} yw - zv \\ zu - xw \\ xv - yu \end{pmatrix} d\vec{r} = \begin{pmatrix} \int_{\Omega} \rho(\vec{r})(yw - zv) d\vec{r} \\ \int_{\Omega} \rho(\vec{r})(zu - xw) d\vec{r} \\ \int_{\Omega} \rho(\vec{r})(xv - yu) d\vec{r} \end{pmatrix} = \begin{pmatrix} H_x \\ H_y \\ H_z \end{pmatrix} \quad (370)$$

while the inertia tensor for a continuous body is given by

$$\mathbf{I} = \int_{\Omega} \rho(\vec{r}) [\vec{r} \cdot \vec{r} \mathbf{1} - \vec{r} \times \vec{r}] d\vec{r} \quad (371)$$

$$= \int_{\Omega} \rho(\vec{r}) \left[ \begin{pmatrix} x^2 + y^2 + z^2 & 0 & 0 \\ 0 & x^2 + y^2 + z^2 & 0 \\ 0 & 0 & x^2 + y^2 + z^2 \end{pmatrix} - \begin{pmatrix} xx & xy & xz \\ yx & yy & yz \\ zx & zy & zz \end{pmatrix} \right] d\vec{r} \quad (372)$$

$$= \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} y^2 + z^2 & -xy & -xz \\ -yx & x^2 + z^2 & -yz \\ -zx & -zy & x^2 + y^2 \end{pmatrix} d\vec{r} \quad (373)$$

$$= \begin{pmatrix} \int_{\Omega} \rho(\vec{r})(y^2 + z^2) d\vec{r} & -\int_{\Omega} \rho(\vec{r})xy d\vec{r} & -\int_{\Omega} \rho(\vec{r})xz d\vec{r} \\ -\int_{\Omega} \rho(\vec{r})yx d\vec{r} & \int_{\Omega} \rho(\vec{r})(x^2 + z^2) d\vec{r} & -\int_{\Omega} \rho(\vec{r})yz d\vec{r} \\ -\int_{\Omega} \rho(\vec{r})zx d\vec{r} & -\int_{\Omega} \rho(\vec{r})zy d\vec{r} & \int_{\Omega} \rho(\vec{r})(x^2 + y^2) d\vec{r} \end{pmatrix} \quad (374)$$

$$= \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} \quad (375)$$

### 7.28.2 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

$$\vec{H} = \int_{\Omega} \rho(\vec{r}) \begin{pmatrix} 0 \\ 0 \\ xv - yu \end{pmatrix} d\vec{r} = \begin{pmatrix} 0 \\ 0 \\ \int_{\Omega} \rho(\vec{r})(xv - yu) d\vec{r} \end{pmatrix} \quad (376)$$

and

$$\mathbf{I} = \begin{pmatrix} I_{xx} & I_{xy} & I_{xz} \\ I_{yx} & I_{yy} & I_{yz} \\ I_{zx} & I_{zy} & I_{zz} \end{pmatrix} = \begin{pmatrix} I_{xx} & I_{xy} & 0 \\ I_{yx} & I_{yy} & 0 \\ 0 & 0 & I_{zz} \end{pmatrix} \quad (377)$$

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 G.2):

$$\omega_x = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} 0 & I_{xy} & 0 \\ 0 & I_{yy} & 0 \\ H_3 & 0 & I_{zz} \end{vmatrix} = 0 \quad (378)$$

$$\omega_y = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} I_{xx} & 0 & 0 \\ I_{yx} & 0 & 0 \\ 0 & H_z & I_{zz} \end{vmatrix} = 0 \quad (379)$$

$$\omega_z = \frac{1}{\det(\mathbf{I})} \begin{vmatrix} I_{xx} & I_{xy} & 0 \\ I_{yx} & I_{yy} & 0 \\ 0 & 0 & H_z \end{vmatrix} \quad (380)$$

$$= \frac{1}{\det(\mathbf{I})} (I_{xx}I_{yy}H_z - I_{yx}I_{xy}H_z) \quad (381)$$

$$= \frac{1}{\det(\mathbf{I})} (I_{xx}I_{yy} - I_{yx}I_{xy}) H_z \quad (382)$$

with  $\det(\mathbf{I}) = I_{xx}I_{yy}I_{zz} - I_{yx}I_{xy}I_{zz}$ . Concretely, this means that in 2D one does not need to solve the system  $\mathbf{I} \cdot \vec{\omega} = \vec{H}$  since only  $\omega_z$  is not zero.

Likewise, the volume-averaged angular velocity is then simply:

$$\omega_z = \frac{1}{|\Omega|} \int_{\Omega} \frac{xv - yu}{r^2} d\vec{r} \quad (383)$$

## 7.29 Picard and Newton

explain why our eqs are nonlinear

### 7.29.1 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  $\vec{T}^{new}$
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  $\epsilon$ . 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  $\epsilon = 10^{-3}$ K) we could consider that iterations have converged but looking now at velocities which are of the order of a cm/year (i.e.  $\sim 3 \cdot 10^{-11}$ 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  $\epsilon$  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) [477]:

$$\left( \frac{(\vec{X}^{old} - \vec{X}^{new}) \cdot (\vec{X}^{old} - \vec{X}^{new})}{\vec{X}^{new} \cdot \vec{X}^{new}} \right)^{1/2} < \epsilon$$

Yet another convergence criterion is used in [542]: the means  $\langle \vec{X}^{old} \rangle$ ,  $\langle \vec{X}^{new} \rangle$  as well as the variances  $\sigma^{old}$  and  $\sigma^{new}$  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

$$\mathbb{A}(\vec{X}^n) \cdot \vec{X}^* = \vec{b}(\vec{X}^n)$$

and then updates  $\vec{X}^n$  as follows:

$$\vec{X}^n = \gamma \vec{X}^n + (1 - \gamma) \vec{X}^* \quad 0 < \gamma \leq 1$$

When  $\gamma = 1$  we recover the standard Picard iterations formula above.

### 7.30 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 [209]) 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:

$$\begin{pmatrix} \mathbb{K}_k & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \delta \vec{\mathcal{V}}_k \\ \delta \vec{\mathcal{P}}_k \end{pmatrix} = \begin{pmatrix} -\vec{F}_k^{\mathcal{V}} \\ -\vec{F}_k^{\mathcal{P}} \end{pmatrix}$$

### 7.31 Parallel or not?

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 [544]). 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

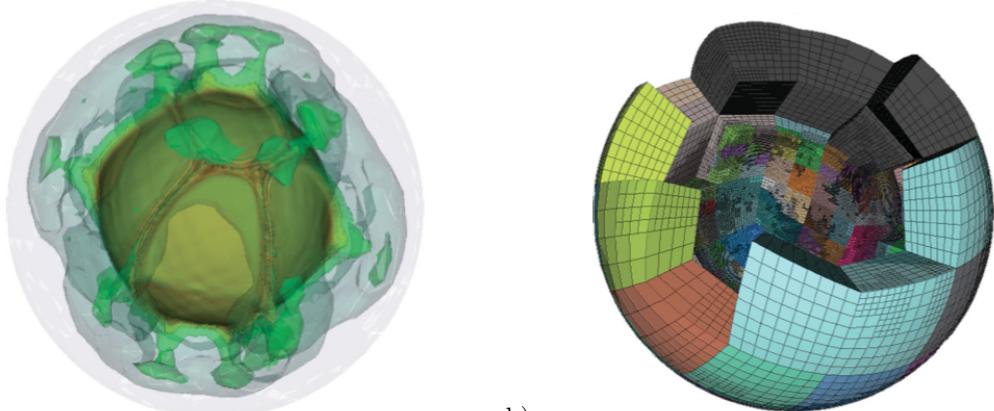
$$N_{particles} \geq 10^{10}$$

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.

An example of such a large parallel calculation of 3D convection with domain decomposition in a spherical shell can be found in [358]:



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.)

## 7.32 Stream function

### 7.32.1 In Cartesian coordinates

The Stream function (commonly denoted by  $\Phi$  or  $\Psi$ ) approach is a useful approach in fluid dynamics as it can provide relatively quick solutions to 2D incompressible flow problems. Lines of constant  $\Phi$  are called stream lines and give a useful representation of the flow. The definition of the stream function is such that

$$u = -\frac{\partial \Phi}{\partial y} \quad (384)$$

$$v = \frac{\partial \Phi}{\partial x} \quad (385)$$

It then follows that the velocity field based on the above equations automatically fulfills the continuity equation:

$$\vec{\nabla} \cdot \vec{v} = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = -\frac{\partial^2 \Phi}{\partial x \partial y} + \frac{\partial^2 \Phi}{\partial y \partial x} = 0$$

The stream function can also be substituted into the (constant viscosity) Stokes equation  $-\vec{\nabla} p + \eta \Delta \vec{v} = \vec{0}$ :

$$-\frac{\partial p}{\partial x} - \eta \left( \frac{\partial^3 \Phi}{\partial^2 x \partial y} + \frac{\partial^3 \Phi}{\partial^3 y} \right) = 0 \quad (386)$$

$$-\frac{\partial p}{\partial y} - \eta \left( \frac{\partial^3 \Phi}{\partial^3 x} + \frac{\partial^3 \Phi}{\partial x \partial^2 y} \right) = 0 \quad (387)$$

We can now eliminate the pressure term by taking the partial derivative of the first equation with respect to  $y$  and the partial derivative of the second one with respect to  $x$ , and substracting both. We get:

$$\frac{\partial^4 \Phi}{\partial x^4} + \frac{\partial^4 \Phi}{\partial x^2 \partial y^2} + \frac{\partial^4 \Phi}{\partial y^4} = 0 \quad (388)$$

or,

$$\left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \Phi = 0 \quad (389)$$

or,

$$\vec{\nabla}^2 \vec{\nabla}^2 \Phi = \vec{\nabla}^4 \Phi = 0$$

which is known as the Biharmonic operator.

### 7.32.2 In Cylindrical coordinates

TODO

VERIFY THOSE! minus signs ?

$$\begin{aligned} v_r &= \frac{1}{r} \frac{\partial \Phi}{\partial \theta} \\ v_\theta &= -\frac{\partial \Phi}{\partial r} \end{aligned}$$

### 7.33 Corner flow

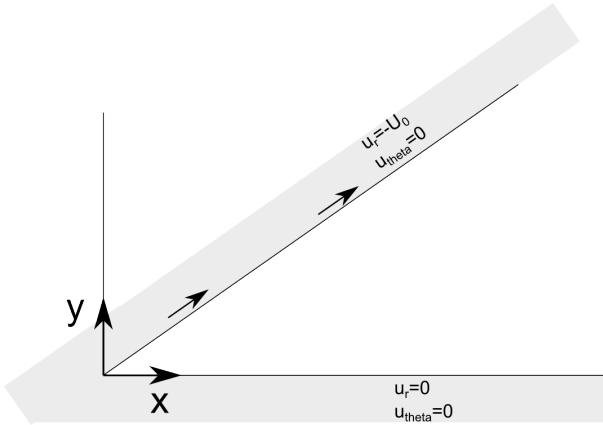
The mantle wedge comprised between the downgoing slab and the overriding plate has been extensively studied since very important geodynamical processes take place in it or right above it (slab dehydration and water transport, melting, over-riding plate deformation, vulcanism, ...).

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 [36].

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  $\Phi$ , 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 7.32, the stream function  $\Phi$  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

$$v_r = \frac{1}{r} \frac{\partial \Phi}{\partial \theta} = f'(\theta)$$

$$v_\theta = -\frac{\partial \Phi}{\partial r} = -f(\theta)$$

$A, B, C$  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,

$$f'(0) = A + D = 0 \quad (390)$$

$$f(0) = B = 0 \quad (391)$$

$$f'(\theta_0) = -U_0 \quad (392)$$

$$f(\theta_0) = 0 \quad (393)$$

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

$$\mathbf{e}_r = \cos \theta \mathbf{e}_x + \sin \theta \mathbf{e}_y \quad (394)$$

$$\mathbf{e}_\theta = -\sin \theta \mathbf{e}_x + \cos \theta \mathbf{e}_y \quad (395)$$

so that the velocity field can be expressed in cartesian coordinates:

$$\begin{aligned}\mathbf{v} &= \mathbf{v}_r \mathbf{e}_r + \mathbf{v}_\theta \mathbf{e}_\theta \\ &= \mathbf{v}_r (\cos \theta \mathbf{u}_x + \sin \theta \mathbf{u}_y) + \mathbf{v}_\theta (-\sin \theta \mathbf{u}_x + \cos \theta \mathbf{u}_y) \\ &= (\mathbf{v}_r \cos \theta - \mathbf{v}_\theta \sin \theta) \mathbf{e}_x + (\mathbf{v}_r \sin \theta + \mathbf{v}_\theta \cos \theta) \mathbf{e}_y\end{aligned} \quad (396)$$

## 8 fieldstone\_01: simple analytical solution

This benchmark was developed in collaboration with Job Mos.

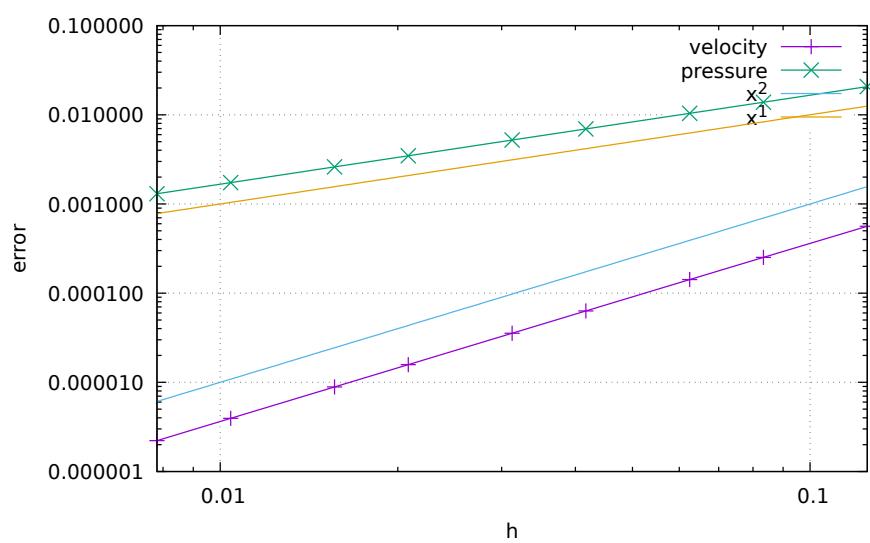
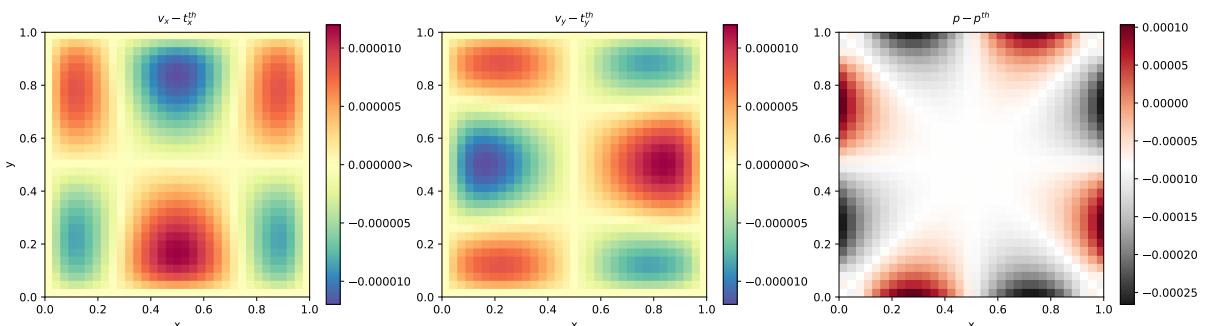
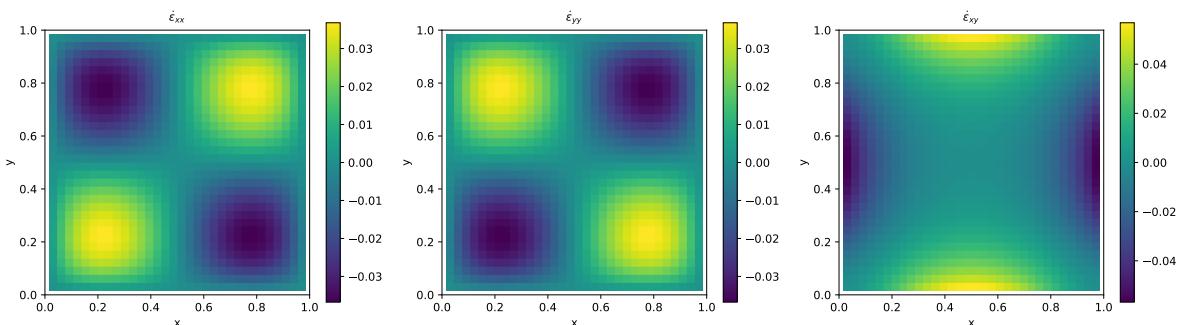
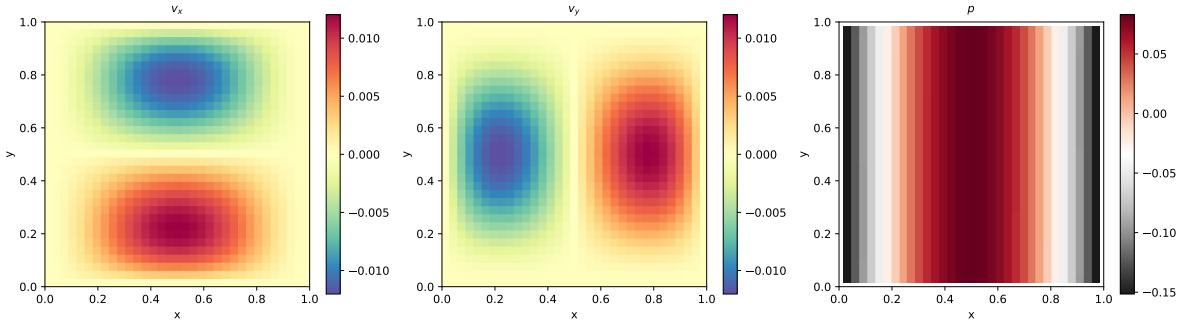
This benchmark is taken from [165] and is described fully in section 7.4. In order to illustrate the behavior of selected mixed finite elements in the solution of stationary Stokes flow, we consider a two-dimensional problem in the square domain  $\Omega = [0, 1] \times [0, 1]$ , which possesses a closed-form analytical solution. The problem consists of determining the velocity field  $\mathbf{v} = (u, v)$  and the pressure  $p$  such that

$$\begin{aligned}\eta \Delta \vec{v} - \vec{\nabla} p + \vec{b} &= \vec{0} && \text{in } \Omega \\ \vec{\nabla} \cdot \vec{v} &= 0 && \text{in } \Omega \\ \vec{v} &= \vec{0} && \text{on } \Gamma_D\end{aligned}$$

where the fluid viscosity is taken as  $\eta = 1$ .

### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (no-slip)
- direct solver
- isothermal
- isoviscous
- analytical solution



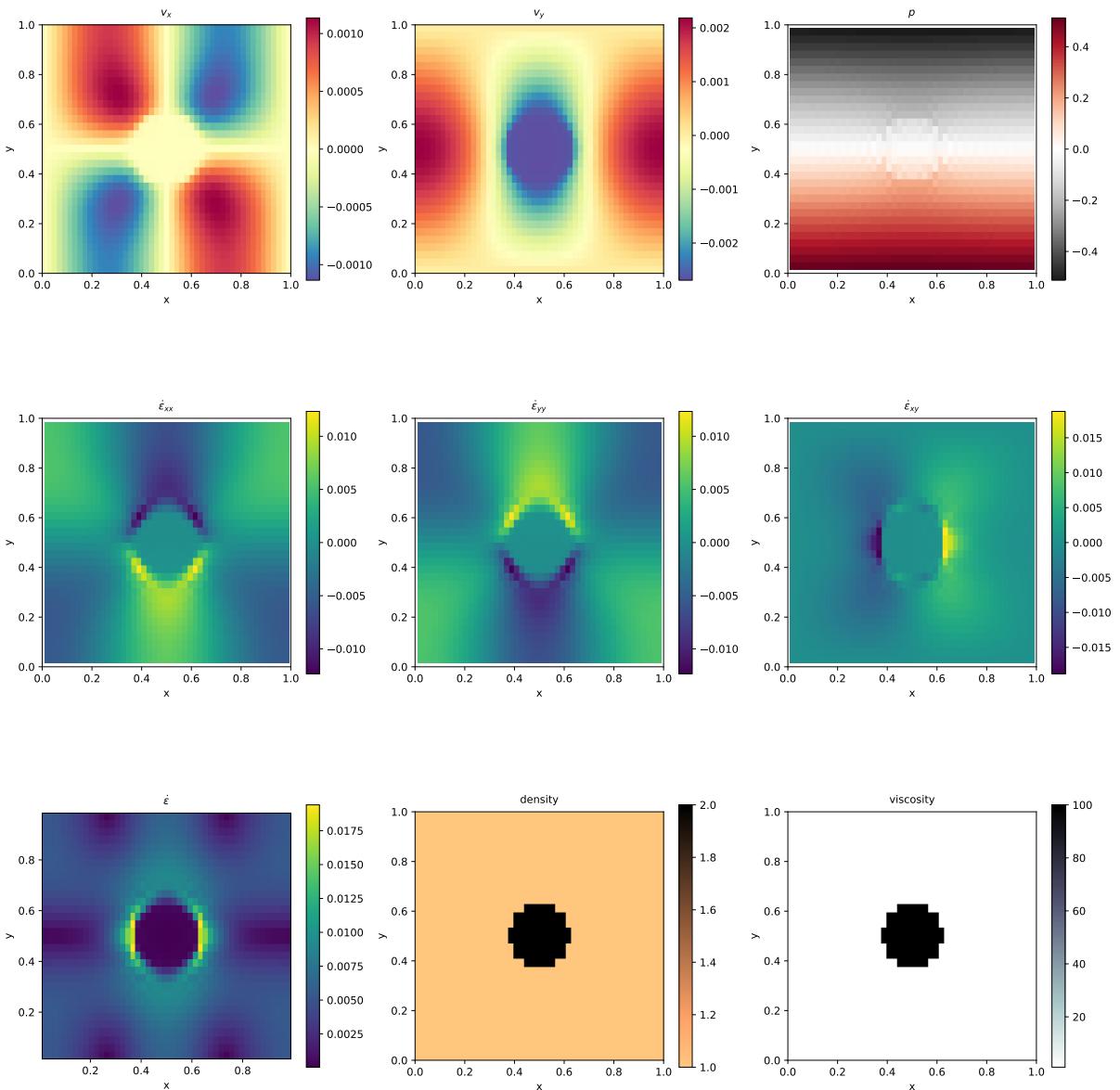
Quadratic convergence for velocity error, linear convergence for pressure error, as expected.

## 9 fieldstone\_02: Stokes sphere

Viscosity and density directly computed at the quadrature points.

### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (free-slip)
- isothermal
- non-isoviscous
- buoyancy-driven flow
- Stokes sphere



## 10 fieldstone\_03: Convection in a 2D 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, I carry out the case 1a, 1b, and 1c experiments as shown in [64]: 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

$$Ra = \frac{\alpha g_y \Delta T h^3}{\kappa \nu} = \frac{\alpha g_y \Delta T h^3 \rho^2 c_p}{k \mu} \quad (397)$$

In what follows, I use the following parameter values:  $L_x = L_y = 1, \rho_0 = c_P = k = \mu = 1, T_0 = 0, \alpha = 10^{-2}, g = 10^2 Ra$  and I run the model with  $Ra = 10^4, 10^5$  and  $10^6$ .

The initial temperature field is given by

$$T(x, y) = (1 - y) - 0.01 \cos(\pi x) \sin(\pi y) \quad (398)$$

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 Nusselt number (i.e. the mean surface temperature gradient over mean bottom temperature) is computed as follows [64]:

$$Nu = L_y \frac{\int \frac{\partial T}{\partial y} (y = L_y) dx}{\int T(y = 0) dx} \quad (399)$$

Note that in our case the denominator is equal to 1 since  $L_x = 1$  and the temperature at the bottom is prescribed to be 1.

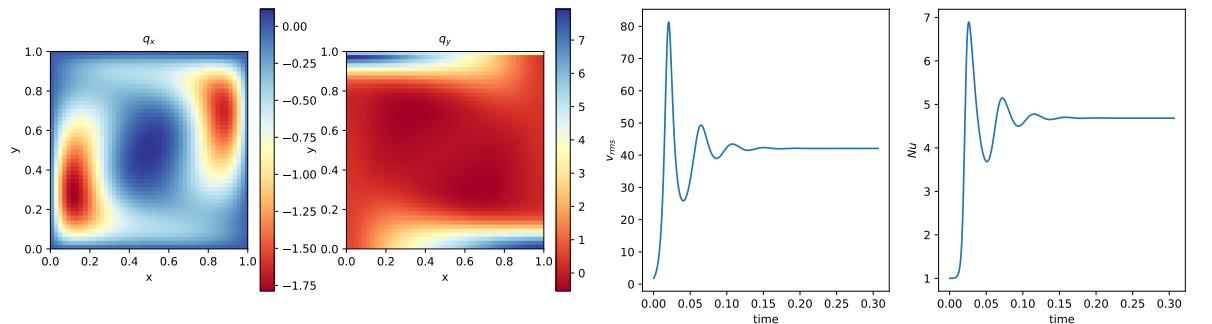
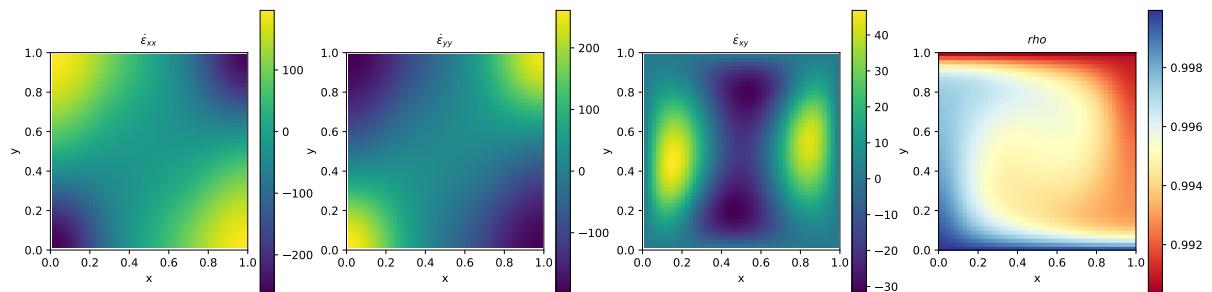
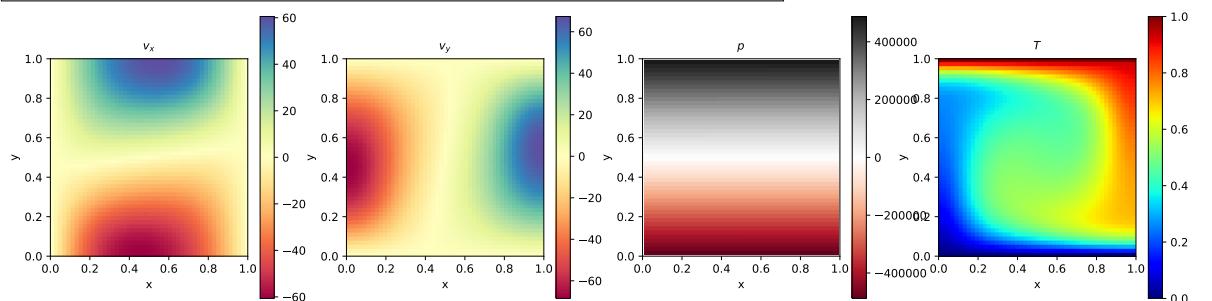
Finally, the steady state root mean square velocity and Nusselt number measurements are indicated in Table ?? alongside those of [64] and [528]. (Note that this benchmark was also carried out and published in other publications [553, 8, 241, 157, 370] but since they did not provide a complete set of measurement values, they are not included in the table.)

		Blankenbach et al	Tackley [528]
$Ra = 10^4$	$V_{rms}$	$42.864947 \pm 0.000020$	42.775
	$Nu$	$4.884409 \pm 0.000010$	4.878
$Ra = 10^5$	$V_{rms}$	$193.21454 \pm 0.00010$	193.11
	$Nu$	$10.534095 \pm 0.000010$	10.531
$Ra = 10^6$	$V_{rms}$	$833.98977 \pm 0.00020$	833.55
	$Nu$	$21.972465 \pm 0.000020$	21.998

Steady state Nusselt number  $Nu$  and  $V_{rms}$  measurements as reported in the literature.

## features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (free-slip)
- Boussinesq approximation
- direct solver
- non-isothermal
- buoyancy-driven flow
- isoviscous
- CFL-condition



ToDo:

- implement steady state criterion
- reach steady state
- do  $\text{Ra}=1\text{e}4, 1\text{e}5, 1\text{e}6$
- plot against blankenbach paper and aspect
- look at critical Ra number

## 11 fieldstone\_04: The lid driven cavity

The lid driven cavity is a famous Computational Fluid Dynamics test case [340, 256, 441, 68, 89, 278, 507] and has been studied in countless publications with a wealth of numerical techniques (see [188] for a succinct review) and also in the laboratory [355].

It models a plane flow of an isothermal isoviscous fluid in a rectangular (usually square) lid-driven cavity. The boundary conditions are indicated in the Fig. ??a. The gravity is set to zero.

### 11.1 the lid driven cavity problem (ldc=0)

In the standard case, the upper side of the cavity moves in its own plane at unit speed, while the other sides are fixed. This thereby introduces a discontinuity in the boundary conditions at the two upper corners of the cavity and yields an uncertainty as to which boundary (side or top) the corner points belong to. In this version of the code the top corner nodes are considered to be part of the lid. If these are excluded the recovered pressure showcases an extremely large checkboard pattern.

This benchmark is usually discussed in the context of low to very high Reynolds number with the full Navier-Stokes equations being solved (with the noticeable exception of [491, 492, 125, 176] which focus on the Stokes equation). In the case of the incompressible Stokes flow, the absence of inertia renders this problem instantaneous so that only one time step is needed.

### 11.2 the lid driven cavity problem - regularisation I (ldc=1)

We avoid the top corner nodes issue altogether by prescribing the horizontal velocity of the lid as follows:

$$u(x) = x^2(1-x)^2. \quad (400)$$

In this case the velocity and its first derivative is continuous at the corners. This is the so-called regularised lid-driven cavity problem [448].

### 11.3 the lid driven cavity problem - regularisation II (ldc=2)

Another regularisation was presented in [160]. Also in Appendix D.4 of [330]. Here, a regularized lid driven cavity is studied which is consistent in the sense that  $\nabla \cdot \mathbf{v} = 0$  holds also at the corners of the domain. There are no-slip conditions at the boundaries  $x = 0$ ,  $x = 1$ , and  $y = 0$ .

The velocity at  $y = 1$  is given by

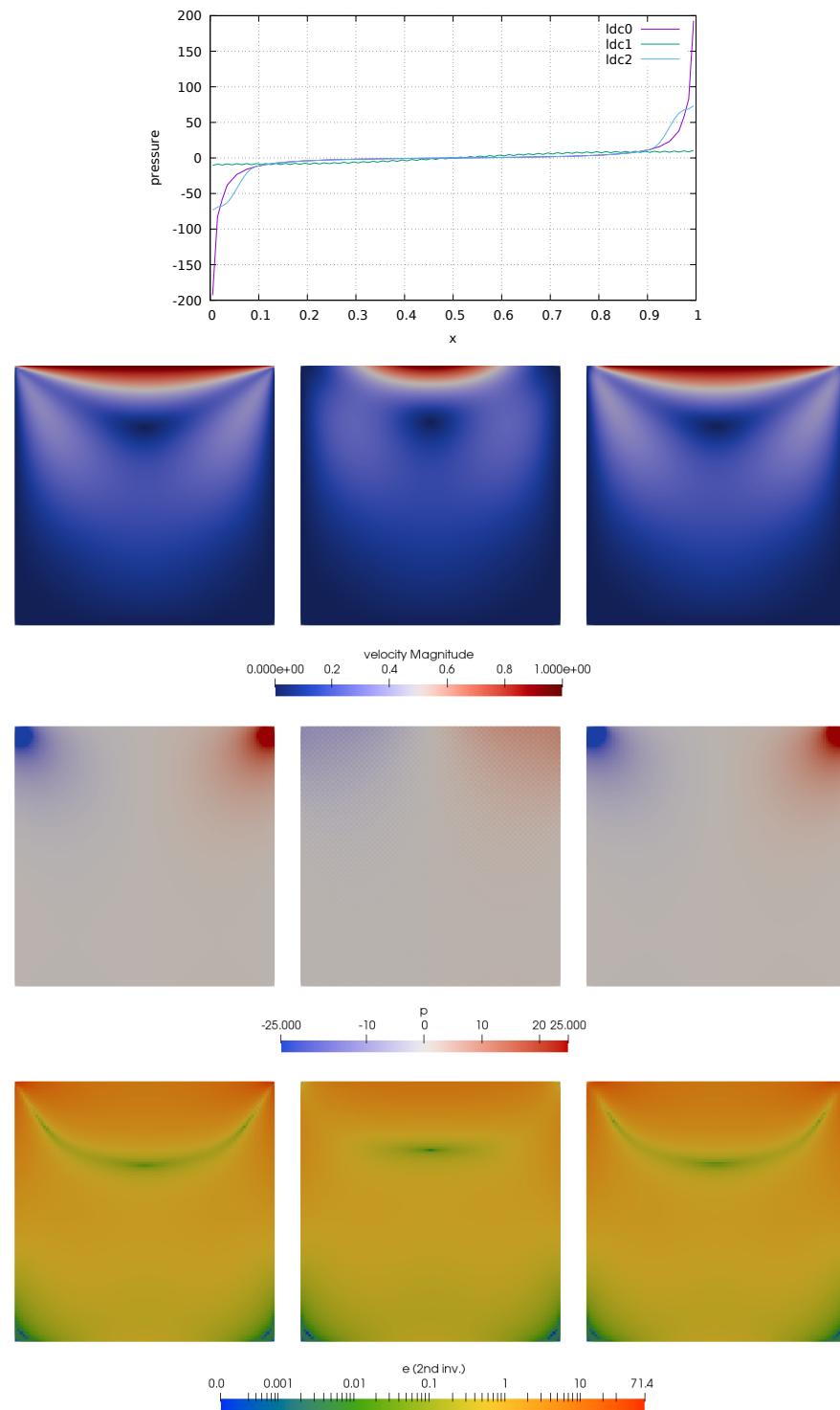
$$\begin{aligned} u(x) &= 1 - \frac{1}{4} \left( 1 - \cos\left(\frac{x_1 - x}{x_1}\pi\right) \right)^2 & x \in [0, x_1] \\ u(x) &= 1 & x \in [x_1, 1 - x_1] \\ u(x) &= 1 - \frac{1}{4} \left( 1 - \cos\left(\frac{x - (1 - x_1)}{x_1}\pi\right) \right)^2 & x \in [1 - x_1, 1] \end{aligned} \quad (401)$$

Results are obtained with  $x_1 = 0.1$ .

#### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- isothermal
- isoviscous

A 100x100 element grid is used. No-slip boundary conditions are prescribed on sides and bottom. A zero vertical velocity is prescribed at the top and the exact form of the prescribed horizontal velocity is controlled by the `ldc` parameter.



## 12 fieldstone\_05: SolCx benchmark

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  $\mathbf{g} = (0, -1)^T$  and the density is given by

$$\rho(x, y) = \sin(\pi y) \cos(\pi x) \quad (402)$$

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:

$$\mu(x, y) = \begin{cases} 1 & \text{for } x < 0.5 \\ 10^6 & \text{for } x > 0.5 \end{cases} \quad (403)$$

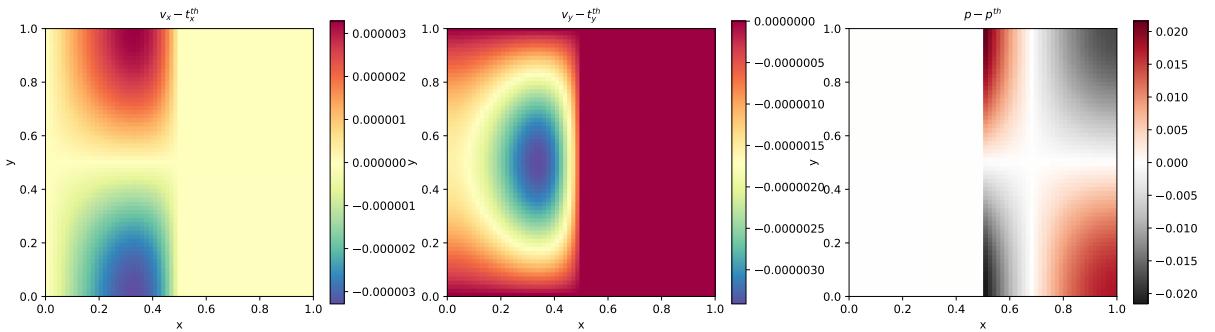
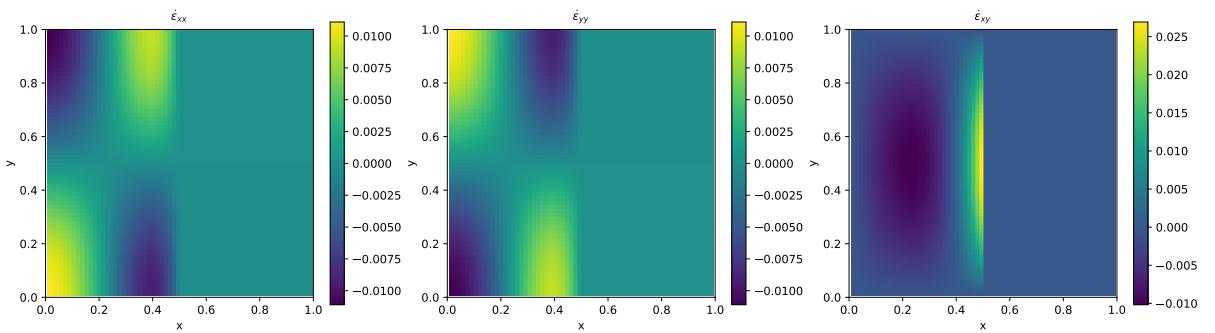
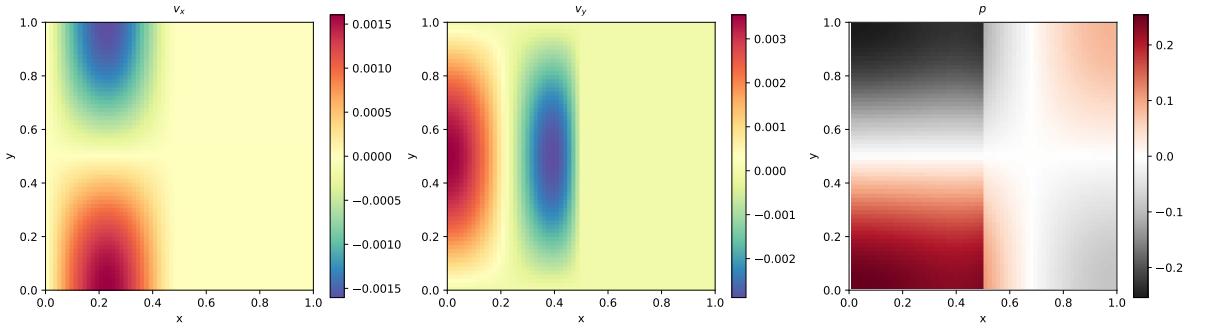
Note the strongly discontinuous viscosity field yields a stagnant flow in the right half of the domain and thereby yields a pressure discontinuity along the interface.

The SolCx benchmark was previously used in [171] (references to earlier uses of the benchmark are available there) and its analytic solution is given in [628]. It has been carried out in [358] and [245]. 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 ([417], <http://underworldproject.org>).

In this particular example, the viscosity is computed analytically at the quadrature points (i.e. tracers are not used to attribute a viscosity to the element). If the number of elements is even in any direction, all elements (and their associated quadrature points) have a constant viscosity(1 or  $10^6$ ). If it is odd, then the elements situated at the viscosity jump have half their integration points with  $\mu = 1$  and half with  $\mu = 10^6$  (which is a pathological case since the used quadrature rule inside elements cannot represent accurately such a jump).

### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- analytical solution



**What we learn from this**

### 13 fieldstone\_06: SolKz benchmark

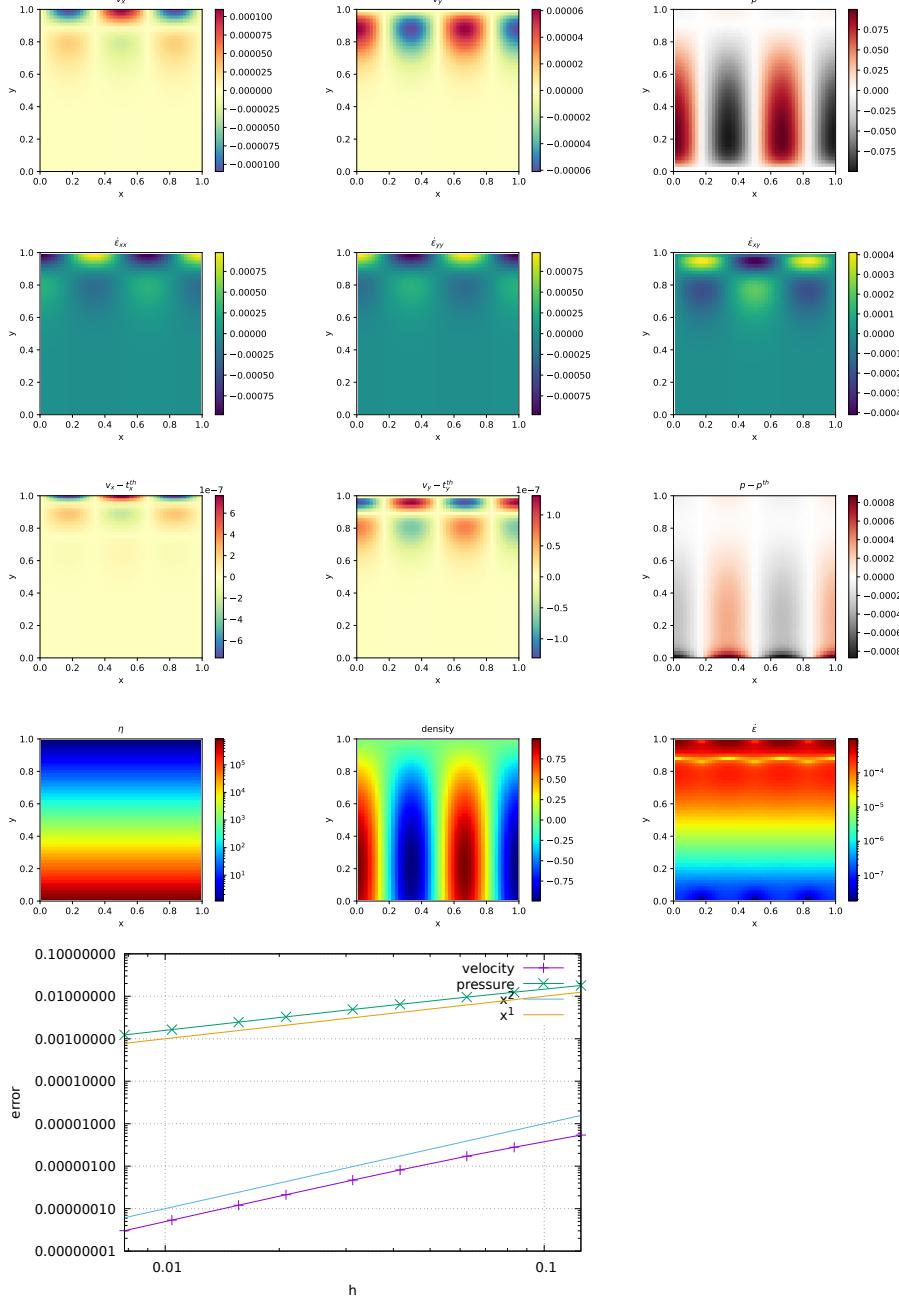
The SolKz benchmark [478] is similar to the SolCx benchmark, but the viscosity is now a function of the space coordinates:

$$\mu(y) = \exp(By) \quad \text{with} \quad B = 13.8155 \quad (404)$$

It is however 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:

$$\rho(x, y) = \sin(2y) \cos(3\pi x) \quad (405)$$

Free slip boundary conditions are imposed on all sides of the domain. This benchmark is presented in [628] as well and is studied in [171] and [245].



## 14 fieldstone\_07: SolVi benchmark

Following SolCx and SolKz, the SolVi inclusion benchmark solves a problem with a discontinuous viscosity field, but in this case the viscosity field is chosen in such a way that the discontinuity is along a circle. Given the regular nature of the grid used by a majority of codes and the present one, 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. [498] derived a simple analytic solution for the pressure and velocity fields for a circular inclusion under simple shear and it was used in [162], [526], [171], [358] and [245].

Because of the symmetry of the problem, we only have to solve over the top right quarter of the domain (see Fig. ??a).

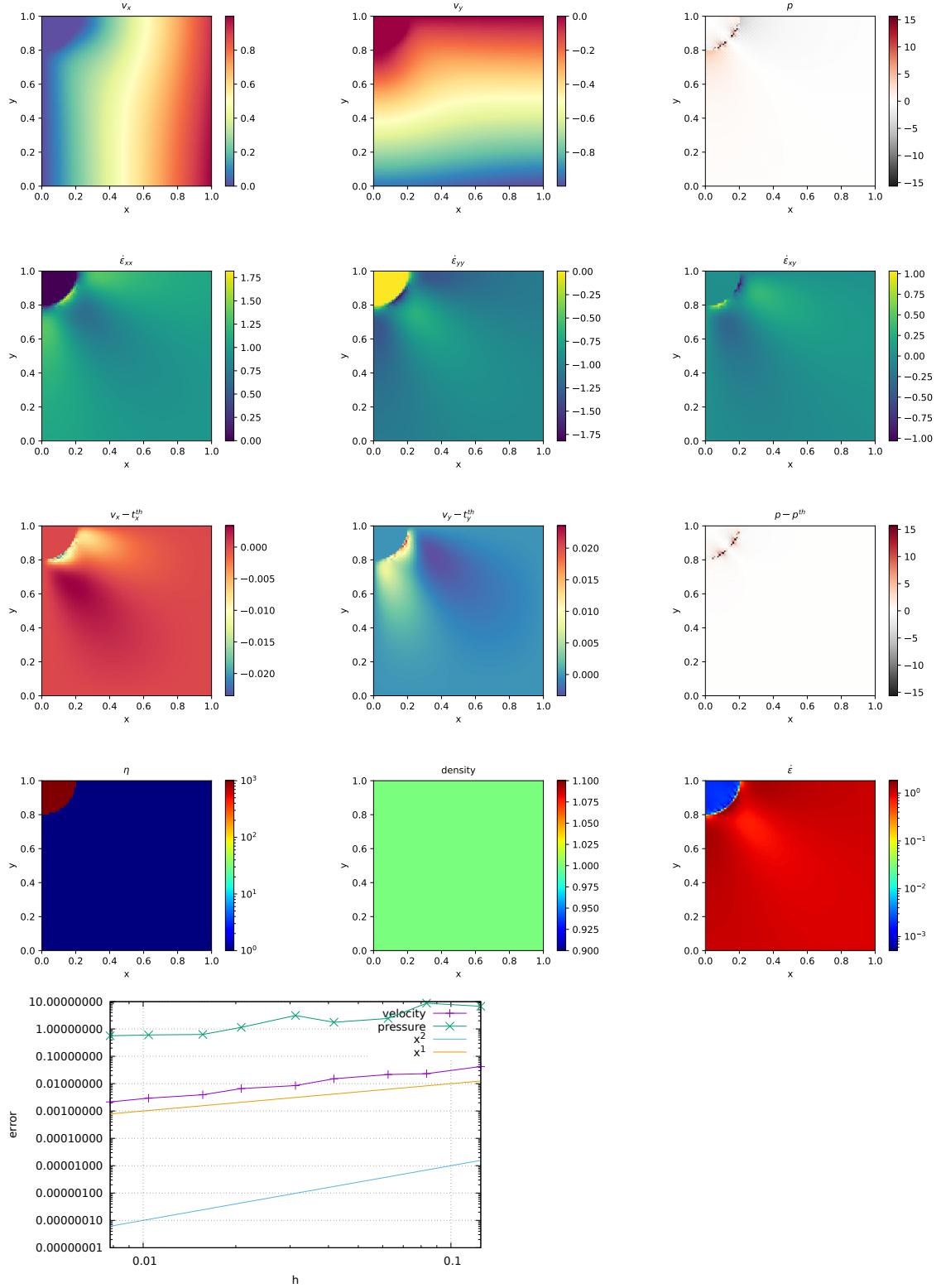
The analytical solution requires a strain rate boundary condition (e.g., pure shear) to be applied far away from the inclusion. In order to avoid using very large domains and/or dealing with this type of boundary condition altogether, the analytical solution is evaluated and imposed on the boundaries of the domain. By doing so, the truncation error introduced while discretizing the strain rate boundary condition is removed.

A characteristic of the analytic solution is that the pressure is zero inside the inclusion, while outside it follows the relation

$$p_m = 4\dot{\epsilon} \frac{\mu_m(\mu_i - \mu_m)}{\mu_i + \mu_m} \frac{r_i^2}{r^2} \cos(2\theta) \quad (406)$$

where  $\mu_i = 10^3$  is the viscosity of the inclusion and  $\mu_m = 1$  is the viscosity of the background media,  $\theta = \tan^{-1}(y/x)$ , and  $\dot{\epsilon} = 1$  is the applied strain rate.

[162] thoroughly investigated this problem with various numerical methods (FEM, FDM), with and without tracers, and conclusively showed how various averagings lead to different results. [171] obtained a first order convergence for both pressure and velocity, while [358] and [245] showed that the use of adaptive mesh refinement in respectively the FEM and FDM yields convergence rates which depend on refinement strategies.



## 15 fieldstone\_08: the indentor benchmark

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) [334].

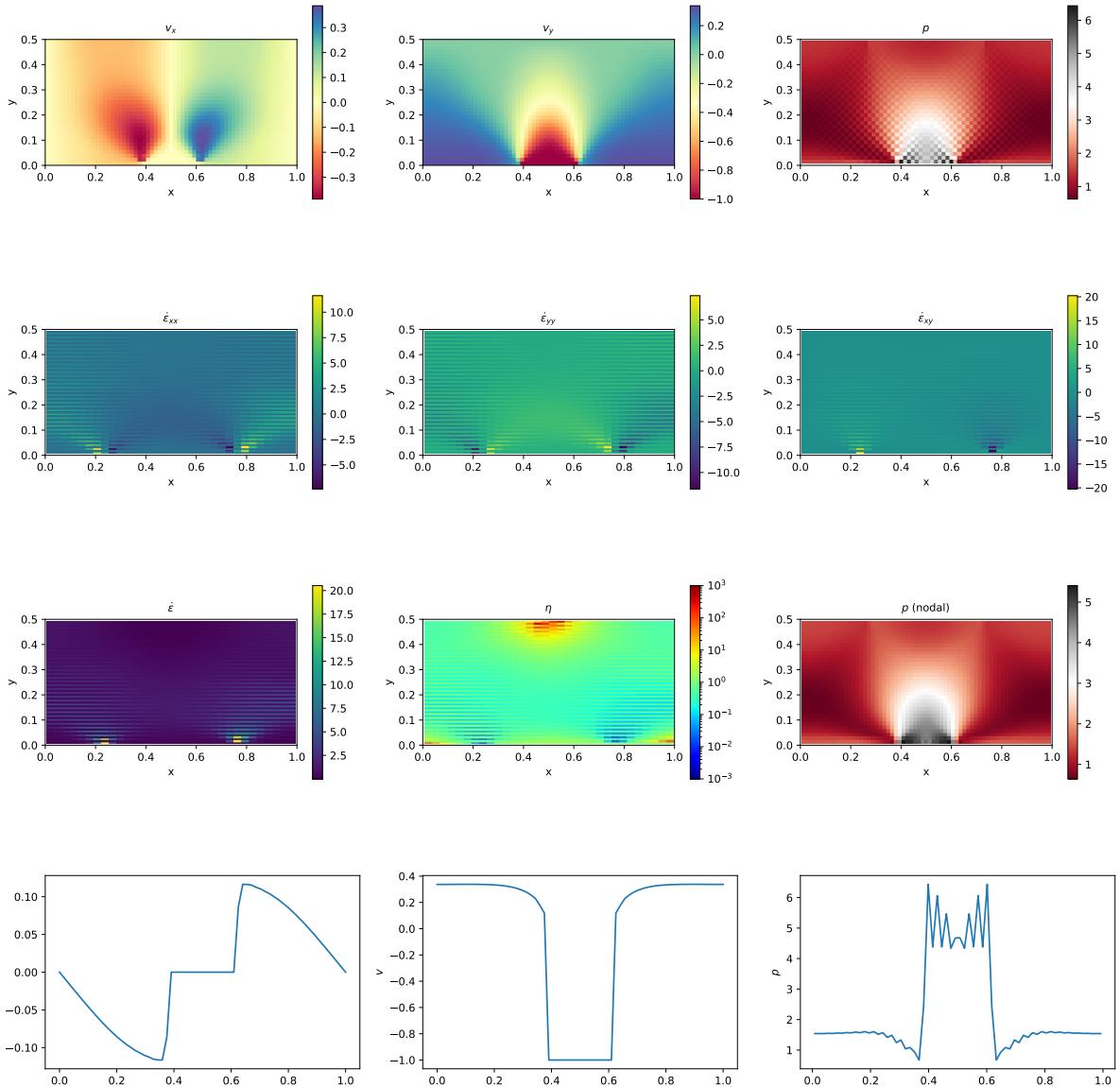
In this experiment, a rigid punch indents a rigid plastic half space; the slip line field theory gives exact solutions as shown in Fig. ??a. The plane strain formulation of the equations and the detailed solution to the problem were derived in the Appendix of [546] and are also presented in [235].

The two dimensional punch problem has been extensively studied numerically for the past 40 years [635, 634, ?, 133, 312, 616, 90, 475] and has been used to draw a parallel with the tectonics of eastern China in the context of the India-Eurasia collision [534, 413]. 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 [442, 158, 331].

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 following parameters are used:  $L_x = 1$ ,  $L_y = 0.5$ ,  $\mu_{min} = 10^{-3}$ ,  $\mu_{max} = 10^3$ ,  $v_p = 1$ ,  $\delta = 0.123456789$  and the yield value of the material is set to  $k = 1$ .

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).



ToDo: smooth punch

#### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (no-slip)
- isothermal
- non-isoviscous
- nonlinear rheology

## 16 fieldstone\_09: the annulus benchmark

This fieldstone was developed in collaboration with Prof. E.G.P. Puckett.

This benchmark is based on Thieulot & Puckett [Subm.] in which an analytical solution to the isoviscous incompressible Stokes equations is derived in an annulus geometry. The velocity and pressure fields are as follows:

$$v_r(r, \theta) = g(r)k \sin(k\theta), \quad (407)$$

$$v_\theta(r, \theta) = f(r) \cos(k\theta), \quad (408)$$

$$p(r, \theta) = kh(r) \sin(k\theta), \quad (409)$$

$$\rho(r, \theta) = \aleph(r)k \sin(k\theta), \quad (410)$$

with

$$f(r) = Ar + B/r, \quad (411)$$

$$g(r) = \frac{A}{2}r + \frac{B}{r} \ln r + \frac{C}{r}, \quad (412)$$

$$h(r) = \frac{2g(r) - f(r)}{r}, \quad (413)$$

$$\aleph(r) = g'' - \frac{g'}{r} - \frac{g}{r^2}(k^2 - 1) + \frac{f}{r^2} + \frac{f'}{r}, \quad (414)$$

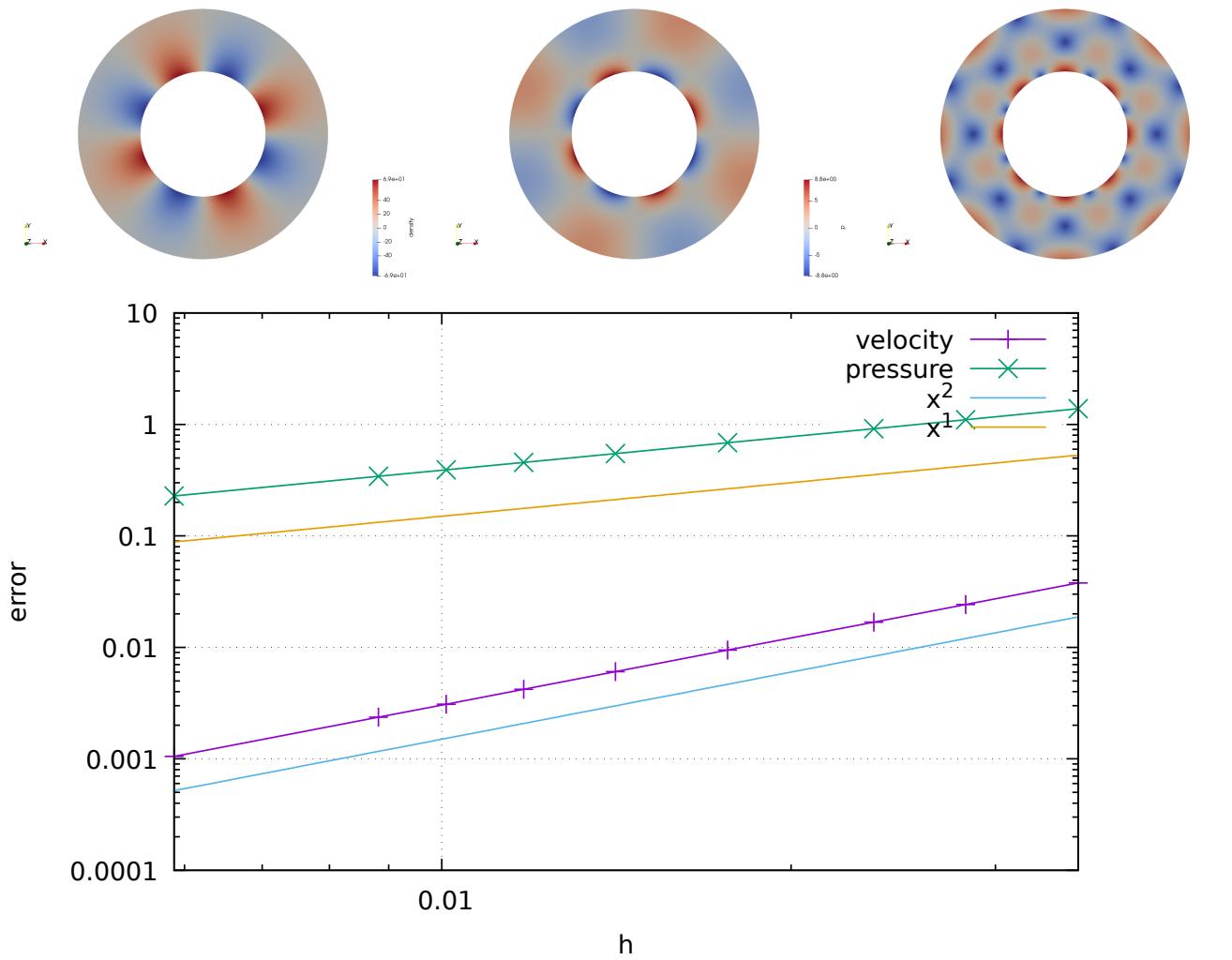
$$A = -C \frac{2(\ln R_1 - \ln R_2)}{R_2^2 \ln R_1 - R_1^2 \ln R_2}, \quad (415)$$

$$B = -C \frac{R_2^2 - R_1^2}{R_2^2 \ln R_1 - R_1^2 \ln R_2}. \quad (416)$$

The parameters  $A$  and  $B$  are chosen so that  $v_r(R_1) = v_r(R_2) = 0$ , i.e. the velocity is tangential to both inner and outer surfaces. The gravity vector is radial and of unit length. In the present case, we set  $R_1 = 1$ ,  $R_2 = 2$  and  $C = -1$ .

### features

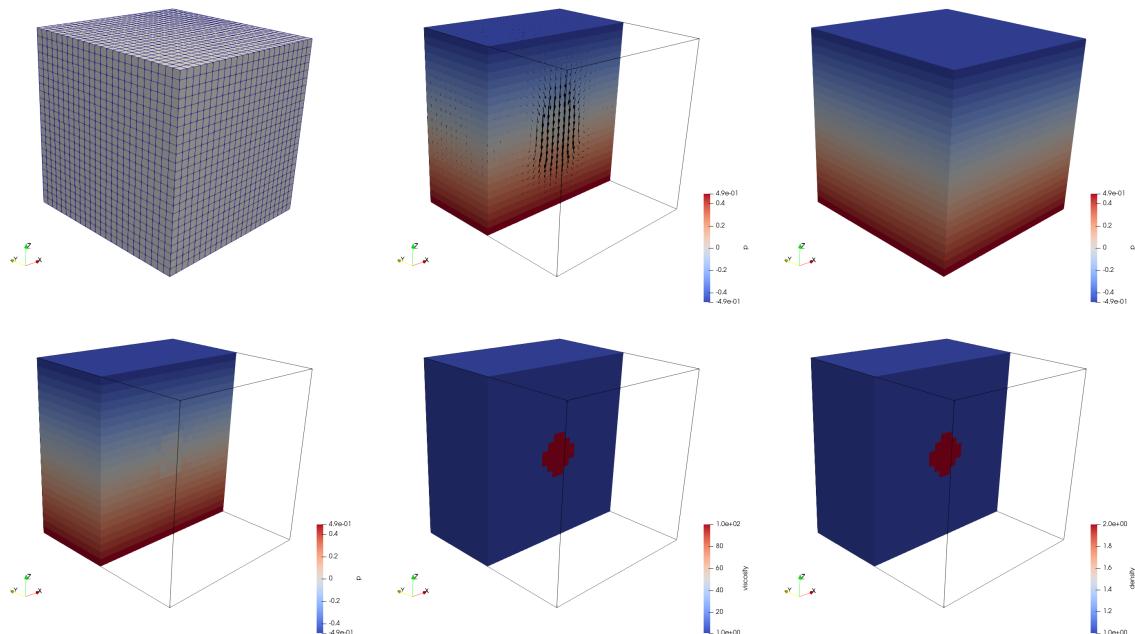
- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions
- direct solver
- isothermal
- isoviscous
- analytical solution
- annulus geometry
- elemental boundary conditions



## 17 fieldstone\_10: Stokes sphere (3D) - penalty

### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- 3D
- elemental b.c.
- buoyancy driven



## 18 fieldstone\_11: stokes sphere (3D) - mixed formulation

This is the same setup as Section 17.

### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- 3D
- elemental b.c.
- buoyancy driven

## 19 fieldstone\_12: consistent pressure recovery

What follows is presented in [633]. The second part of their paper wishes 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 with:

$$p = -\lambda \nabla \cdot \mathbf{v}$$

which lead to

$$(q, p) = -\lambda(q, \nabla \cdot \mathbf{v})$$

and then

$$\left( \int \mathbf{N} \mathbf{N} d\Omega \right) \cdot \mathbf{P} = - \left( \lambda \int \mathbf{N} \nabla \mathbf{N} d\Omega \right) \cdot \mathbf{V}$$

or,

$$\mathbf{M} \cdot \mathbf{P} = -\mathbf{D} \cdot \mathbf{V}$$

and finally

$$\mathbf{P} = -\mathbf{M}^{-1} \cdot \mathbf{D} \cdot \mathbf{V}$$

with  $\mathbf{M}$  of size  $(np \times np)$ ,  $\mathbf{D}$  of size  $(np * ndof \times np * ndof)$  and  $\mathbf{V}$  of size  $(np * ndof)$ . The vector  $\mathbf{P}$  contains the  $np$  nodal pressure values directly, with no need for a smoothing scheme. 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.

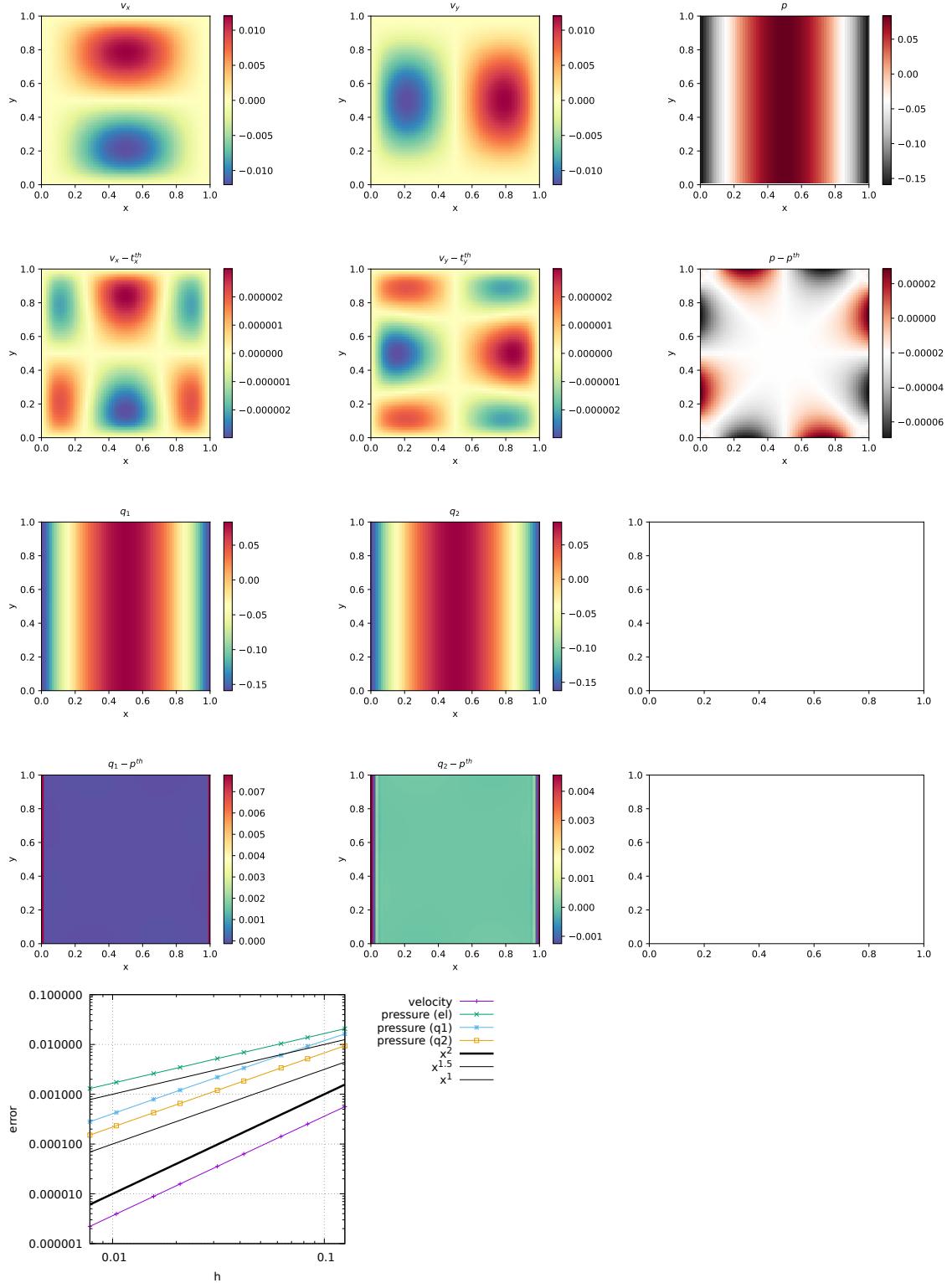
As noted by [633], it is interesting to note that when linear elements are used and the lumped matrices are used for the  $\mathbf{M}$  the resulting algebraic equation is identical to the smoothing scheme based on the averaging method only if the uniform square finite element mesh is used. In this respect this method is expected to yield different results when elements are not square or even rectangular.

---

$q_1$  is smoothed pressure obtained with the center-to-node approach.

$q_2$  is recovered pressure obtained with [633].

All three fulfill the zero average condition:  $\int p d\Omega = 0$ .



In terms of pressure error,  $q_2$  is better than  $q_1$  which is better than elemental.

QUESTION: why are the averages exactly zero ?!

TODO:

- add randomness to internal node positions.
- look at elefant algorithms

## 20 fieldstone\_13: the Particle in Cell technique (1) - the effect of averaging

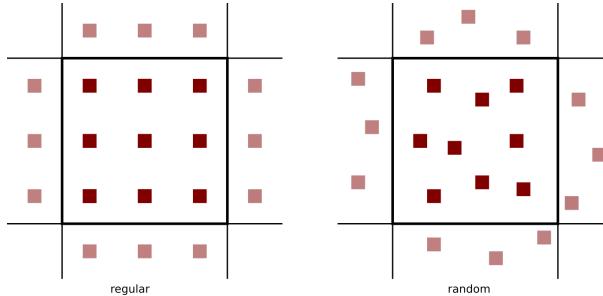
This fieldstone is being developed in collaboration with BSc student Eric Hoogen.

### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (no-slip)
- isothermal
- non-isoviscous
- particle-in-cell

After the initial setup of the grid, markers can then be generated and placed in the domain. One could simply randomly generate the marker positions in the whole domain but unless a *very* large number of markers is used, the chance that an element does not contain any marker exists and this will prove problematic. In order to get a better control over the markers spatial distribution, one usually generates the marker per element, so that the total number of markers in the domain is the product of the number of elements times the user-chosen initial number of markers per element.

Our next concern is how to actually place the markers inside an element. Two methods come to mind: on a regular grid, or in a random manner, as shown on the following figure:



In both cases we make use of the basis shape functions: we generate the positions of the markers (random or regular) in the reference element first ( $r_{im}, s_{im}$ ), and then map those out to the real element as follows:

$$x_{im} = \sum_i^m N_i(r_{im}, s_{im}) x_i \quad y_{im} = \sum_i^m N_i(r_{im}, s_{im}) y_i \quad (417)$$

where  $x_i, y_i$  are the coordinates of the vertices of the element.

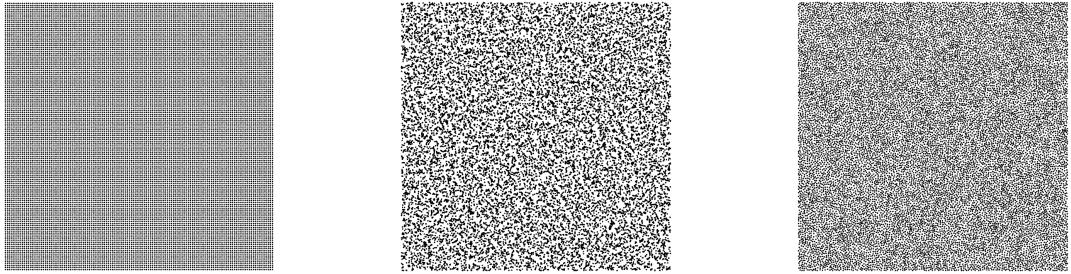
A third option consists in the use of the so-called Poisson-disc sampling which produces points that are tightly-packed, but no closer to each other than a specified minimum distance, resulting in a more natural pattern<sup>29</sup>. Note that the Poisson-disc algorithm fills the whole domain at once, not element after element.

say smthg about avrg dist

insert here theory and link about Poisson disc

---

<sup>29</sup><https://en.wikipedia.org/wiki/SuperSampling>



Left: regular distribution, middle: random, right: Poisson disc.  
16384 markers (32x32 grid, 16 markers per element).

When using *active* markers, one is faced with the problem of transferring the properties they carry to the mesh on which the PDEs are to be solved. As we have seen, building the FE matrix involves a loop over all elements, so one simple approach consists of assigning each element a single property computed as the average of the values carried by the markers in that element. Often in colloquial language "average" refers to the arithmetic mean:

$$\langle \phi \rangle_{am} = \frac{1}{n} \sum_k^n \phi_i \quad (418)$$

where  $\langle \phi \rangle_{am}$  is the arithmetic average of the  $n$  numbers  $\phi_i$ . However, in mathematics other means are commonly used, such as the geometric mean:

$$\langle \phi \rangle_{gm} = \left( \prod_i^n \phi_i \right) \quad (419)$$

PROBLEM with this formula!!!! and the harmonic mean:

$$\langle \phi \rangle_{hm} = \left( \frac{1}{n} \sum_i^n \frac{1}{\phi_i} \right)^{-1} \quad (420)$$

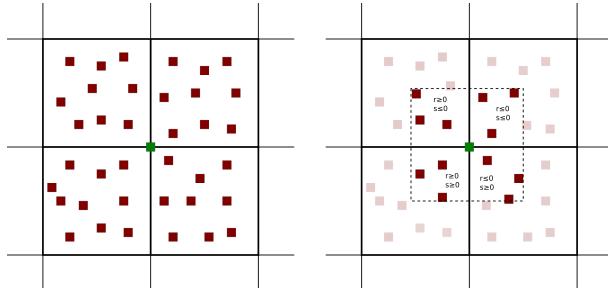
Furthermore, there is a well known inequality for any set of positive numbers,

$$\langle \phi \rangle_{am} \geq \langle \phi \rangle_{gm} \geq \langle \phi \rangle_{hm} \quad (421)$$

which will prove to be important later on.

Let us now turn to a simple concrete example: the 2D Stokes sphere. There are two materials in the domain, so that markers carry the label "mat=1" or "mat=2". For each element an average density and viscosity need to be computed. The majority of elements contains markers with a single material label so that the choice of averaging does not matter (it is trivial to verify that if  $\phi_i = \phi_0$  then  $\langle \phi \rangle_{am} = \langle \phi \rangle_{gm} = \langle \phi \rangle_{hm} = \phi_0$ ). Remain the elements crossed by the interface between the two materials: they contain markers of both materials and the average density and viscosity inside those depends on 1) the total number of markers inside the element, 2) the ratio of markers 1 to markers 2, 3) the type of averaging.

This averaging problem has been studied and documented in the literature [497, 162, 541, 459]



Nodal projection. Left: all markers inside elements to which the green node belongs to are taken into account.  
Right: only the markers closest to the green node count.

Let  $k$  be the green node of the figures above. Let  $(r, s)$  denote the coordinates of a marker inside its element. For clarity, we define the follow three nodal averaging schemes:

- nodal type A:

$$f_k = \frac{\text{sum of values carried by markers in 4 neighbour elements}}{\text{number of markers in 4 neighbour elements}}$$

- nodal type B:

$$f_k = \frac{\text{sum of values carried by markers inside dashed line}}{\text{number of markers in area delimited by the dashed line}}$$

- nodal type C

$$f_k = \frac{\text{sum of values carried by markers in 4 neighbour elements} * N_p(r, s)}{\text{sum of } N_p(r, s)}$$

where  $N_p$  is the  $Q_1$  basis function corresponding to node  $p$  defined on each element. Since these functions are 1 on node  $k$  and then linearly decrease and become zero on the neighbouring nodes, this effectively gives more weight to those markers closest to node  $k$ .

This strategy is adopted in [3, 407] (although it is used to interpolate onto the nodes of  $Q_2P_{-1}$  elements. It is formulated as follows:

"We assume that an arbitrary material point property  $f$ , is discretized via  $f(\mathbf{x}) \simeq \delta(\mathbf{x} - \mathbf{x}_p)f_p$ . We then utilize an approximate local  $L_2$  projection of  $f_p$  onto a continuous  $Q_1$  finite element space. The corner vertices of each  $Q_2$  finite element define the mesh  $f_p$  is projected onto. The local reconstruction for a node  $i$  is defined by

$$\hat{f}_i = \frac{\int_{\Omega_i} N_i(\mathbf{x}) f(\mathbf{x})}{\int_{\Omega_i} N_i(\mathbf{x})} \simeq \frac{\sum_p N_i(\mathbf{x}_p) f_p}{\sum_p N_i(\mathbf{x}_p)}$$

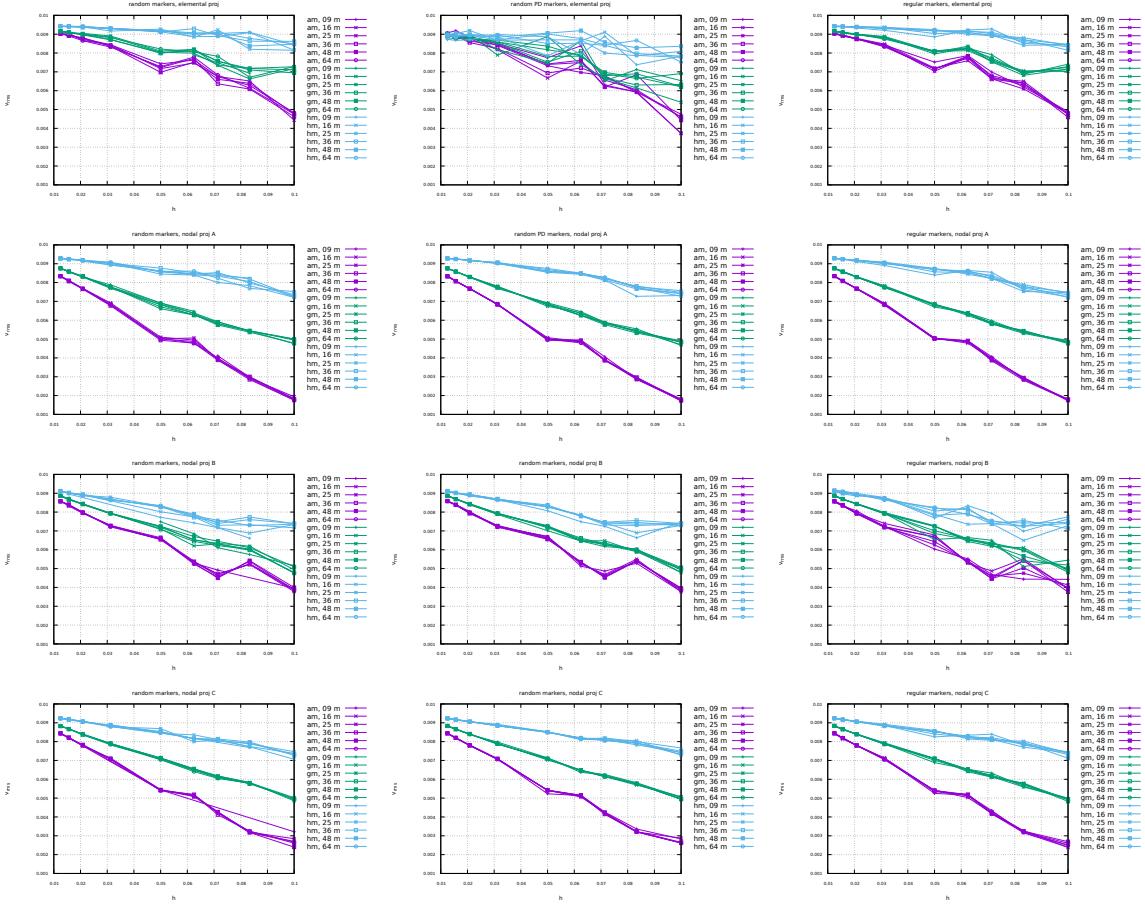
where the summation over  $p$  includes all material points contained within the support  $\Omega_i$  of the trilinear interpolant  $N_i$ ".

The setup is identical to the Stokes sphere experiment. The bash script `script_runall` runs the code for many resolutions, both initial marker distribution and all four averaging types. The viscosity of the sphere has been set to  $10^3$  while the viscosity of the surrounding fluid is 1. The average density is always computed with an arithmetic mean.

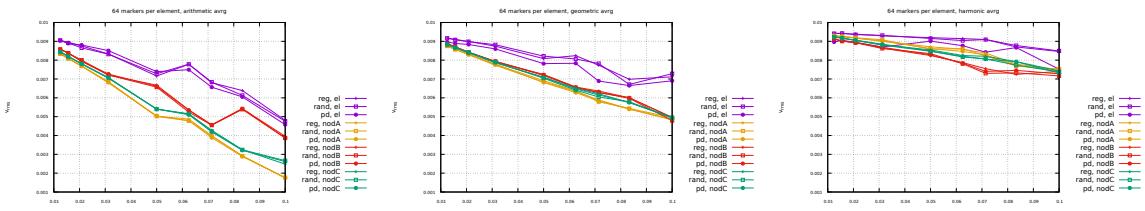
Conclusions:

- With increasing resolution ( $h \rightarrow 0$ ) vrms values seem to converge towards a single value, irrespective of the number of markers.
- At low resolution, say 32x32 (i.e.  $h=0.03125$ ), vrms values for the three averagings differ by about 10%. At higher resolution, say 128x128, vrms values are still not converged.
- The number of markers per element plays a role at low resolution, but less and less with increasing resolution.
- Results for random and regular marker distributions are not identical but follow a similar trend and seem to converge to the same value.
- elemental values yield better results (espcecially at low resolutions)
- harmonic mean yields overal the best results

Root mean square velocity results are shown hereunder:



Left column: random markers, middle column: Poisson disc, right column: regular markers. First row: elemental projection, second row: nodal 1 projection, third row: nodal 2 projection, fourth row: nodal 3 projection.



Left to right: arithmetic, geometric, harmonic averaging for viscosity.

## 21 fieldstone\_f14: solving the full saddle point problem

The details of the numerical setup are presented in Section ??.

The main difference is that we no longer use the penalty formulation and therefore keep both velocity and pressure as unknowns. Therefore we end up having to solve the following system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} V \\ P \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix} \quad \text{or,} \quad \mathbb{A} \cdot X = rhs$$

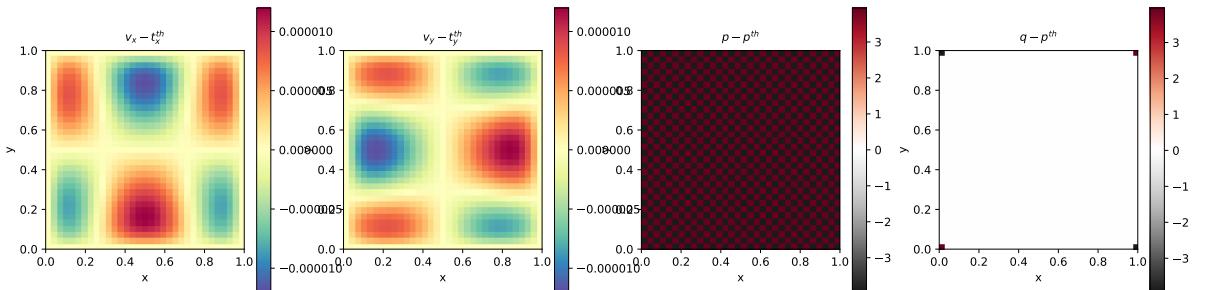
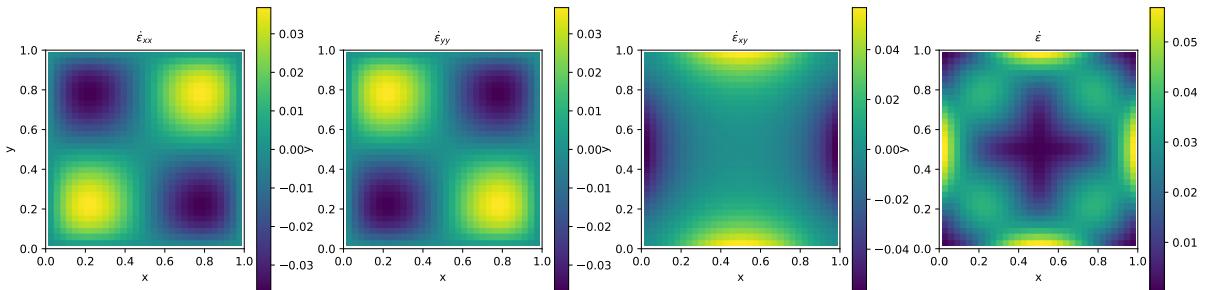
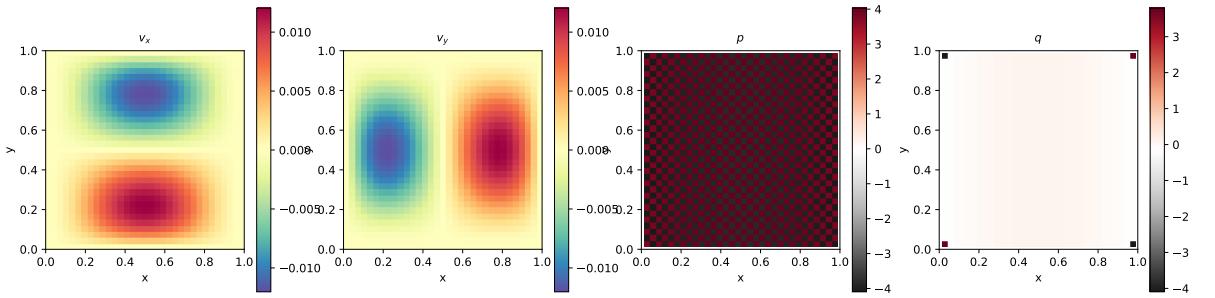
Each block  $\mathbb{K}$ ,  $\mathbb{G}$  and vector  $f$ ,  $h$  are built separately in the code and assembled into the matrix  $\mathbb{A}$  and vector  $rhs$  afterwards.  $\mathbb{A}$  and  $rhs$  are then passed to the solver. We will see later that there are alternatives to solve this approach which do not require to build the full Stokes matrix  $\mathbb{A}$ .

Each element has  $m = 4$  vertices so in total  $ndofV \times m = 8$  velocity dofs and a single pressure dof, commonly situated in the center of the element. The total number of velocity dofs is therefore  $NfemV = nnp \times ndofV$  while the total number of pressure dofs is  $NfemP = nel$ . The total number of dofs is then  $Nfem = NfemV + NfemP$ .

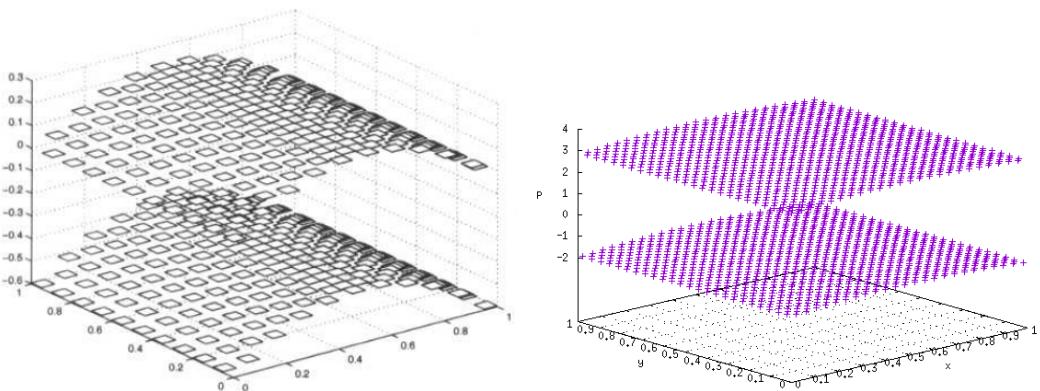
As a consequence, matrix  $\mathbb{K}$  has size  $NfemV, NfemV$  and matrix  $\mathbb{G}$  has size  $NfemV, NfemP$ . Vector  $f$  is of size  $NfemV$  and vector  $h$  is of size  $NfemP$ .

### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Dirichlet boundary conditions (no-slip)
- direct solver (?)
- isothermal
- isoviscous
- analytical solution
- pressure smoothing



Unlike the results obtained with the penalty formulation (see Section ??), the pressure showcases a very strong checkerboard pattern, similar to the one in [165].



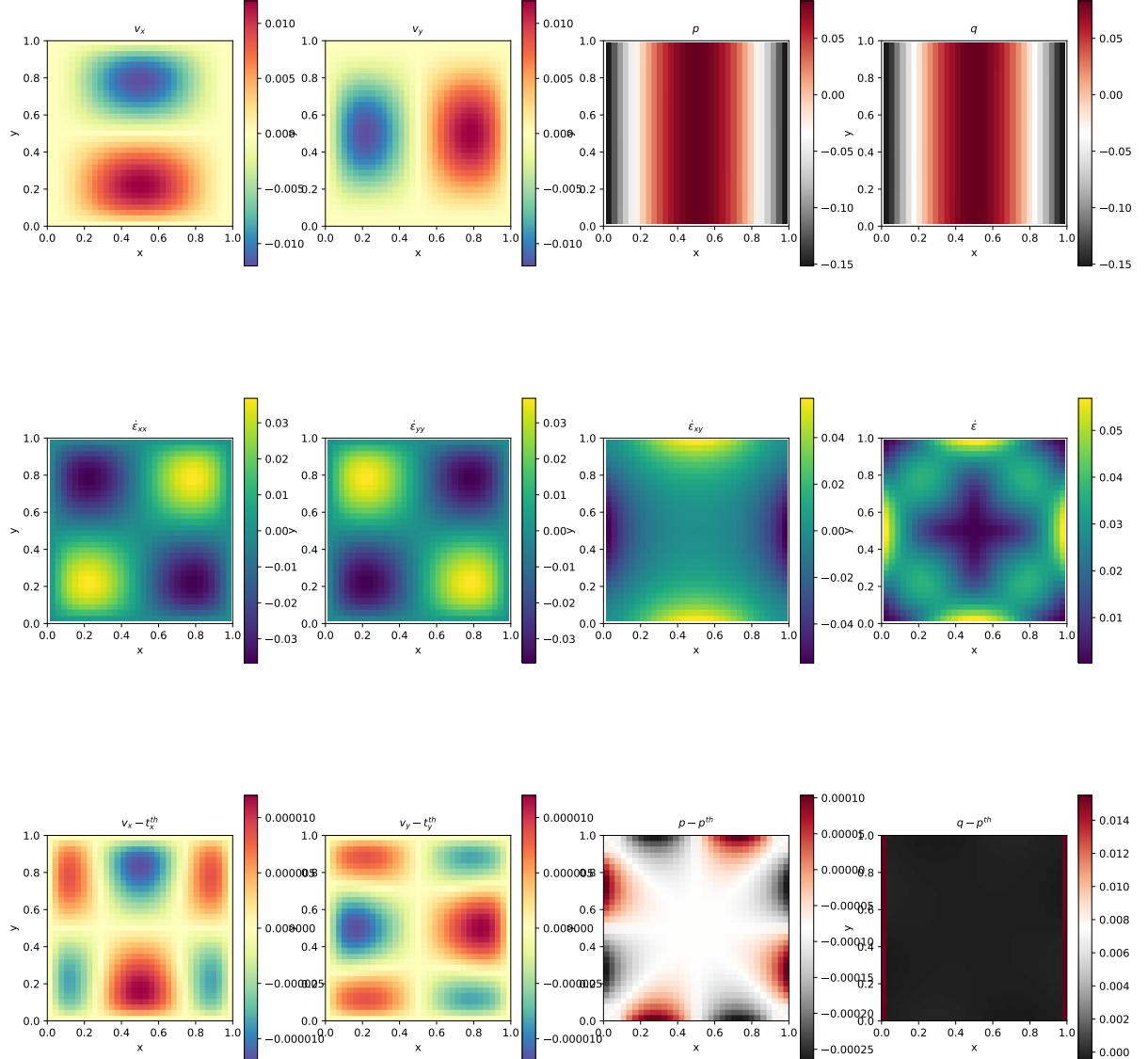
Left: pressure solution as shown in [165]; Right: pressure solution obtained with fieldstone.

Rather interestingly, the nodal pressure (obtained with a simple center-to-node algorithm) fails to recover a correct pressure at the four corners.

Note that the umfpack solver complains a lot about the matrix condition number, even at (very) low resolutions. I believe it does not like the zeros on the (2,2) block of the assembled Stokes matrix.

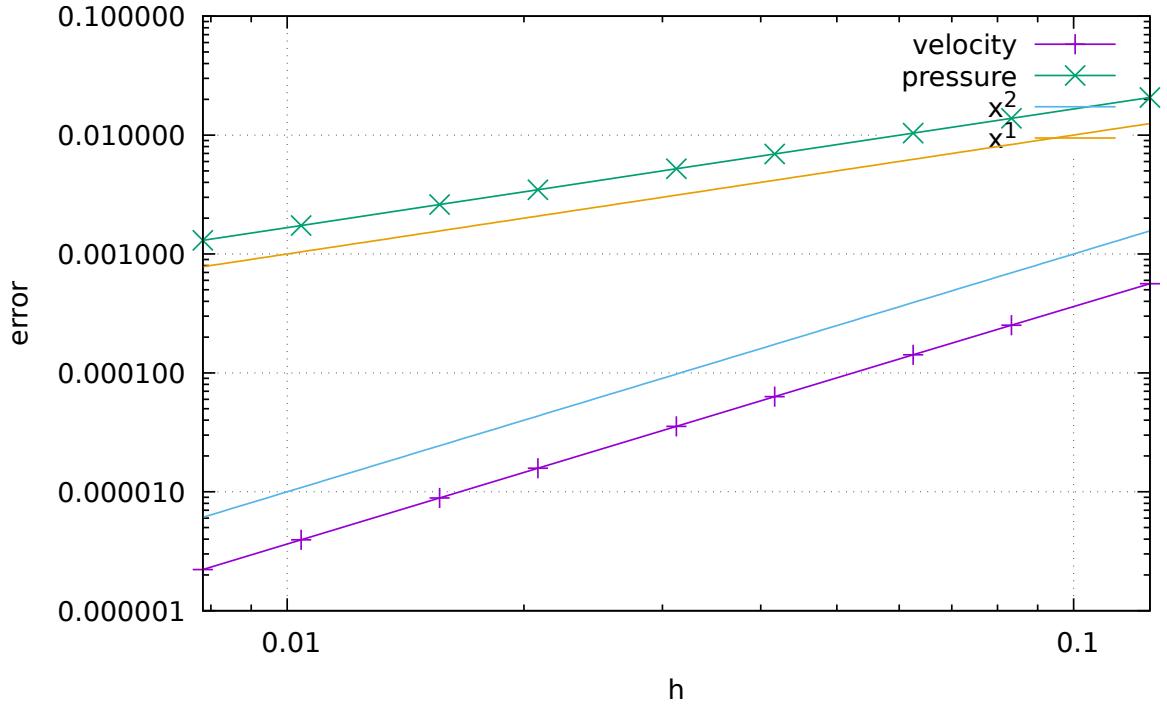
## 22 fieldstone\_f15: saddle point problem with Schur complement approach - benchmark

The details of the numerical setup are presented in Section ???. The main difference resides in the Schur complement approach to solve the Stokes system, as presented in Section ??? (see `solver_cg`). This iterative solver is very easy to implement once the blocks  $\mathbb{K}$  and  $\mathbb{G}$ , as well as the rhs vectors  $f$  and  $h$  have been built.

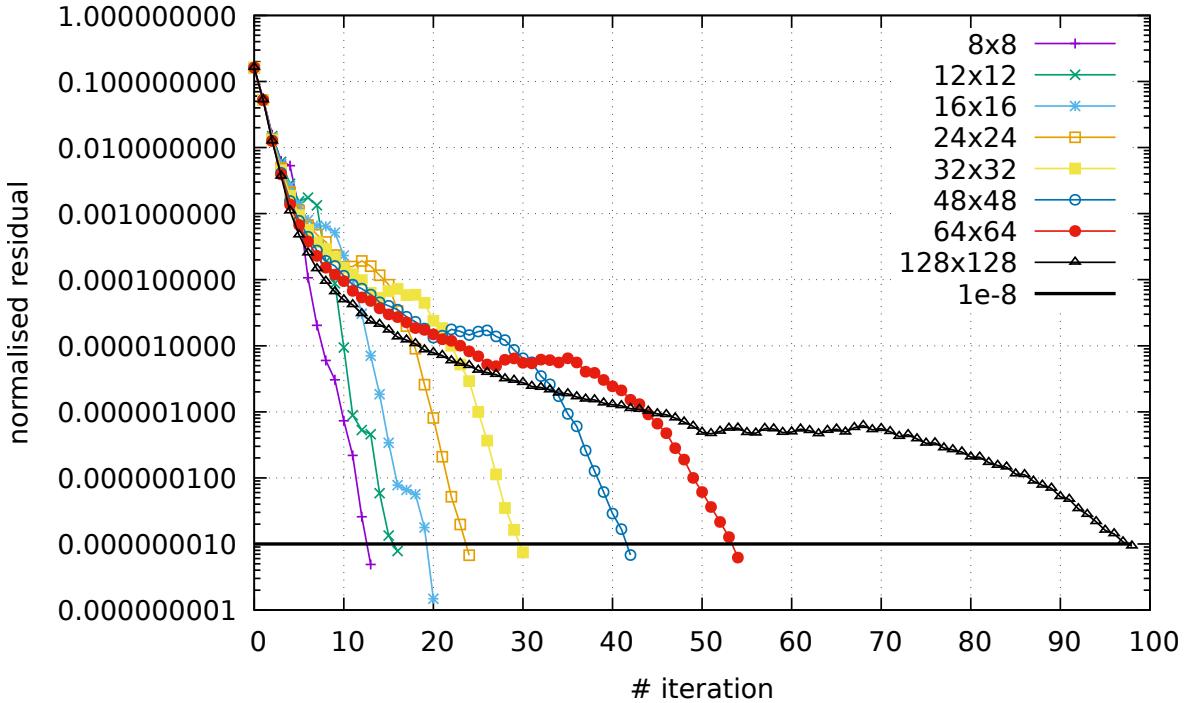


Rather interestingly the pressure checkerboard modes are not nearly as present as in Section ??? which uses a full matrix approach.

Looking at the discretisation errors for velocity and pressure, we of course recover the same rates and values as in the full matrix case.



Finally, for each experiment the normalised residual (see `solver_cg`) was recorded. We see that all things equal the resolution has a strong influence on the number of iterations the solver must perform to reach the required tolerance. This is one of the manifestations of the fact that the  $Q_1 \times P_0$  element is not a stable element: the condition number of the matrix increases with resolution. We will see that this is not the case of stable elements such as  $Q_2 \times Q_1$ .



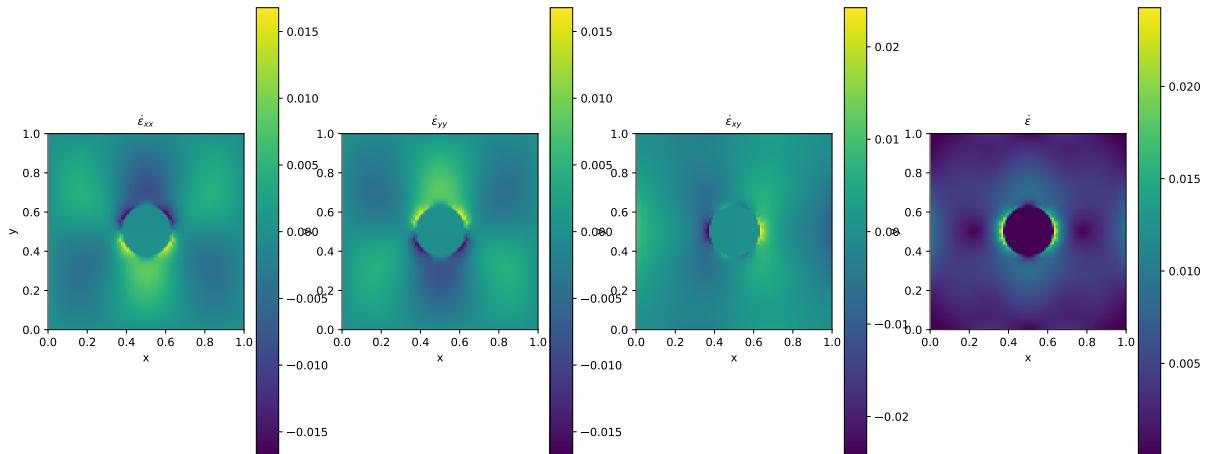
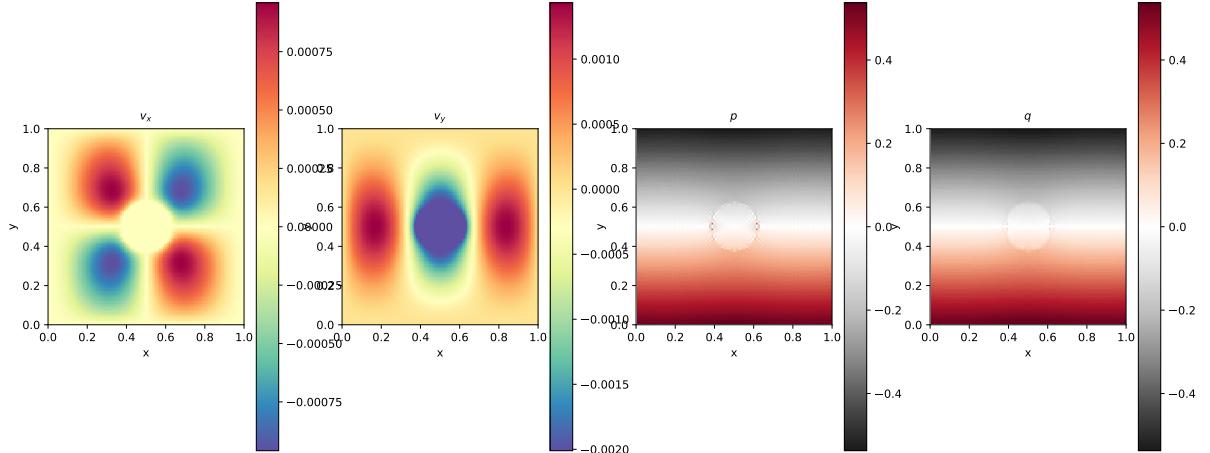
**features**

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Schur complement approach
- isothermal
- isoviscous
- analytical solution

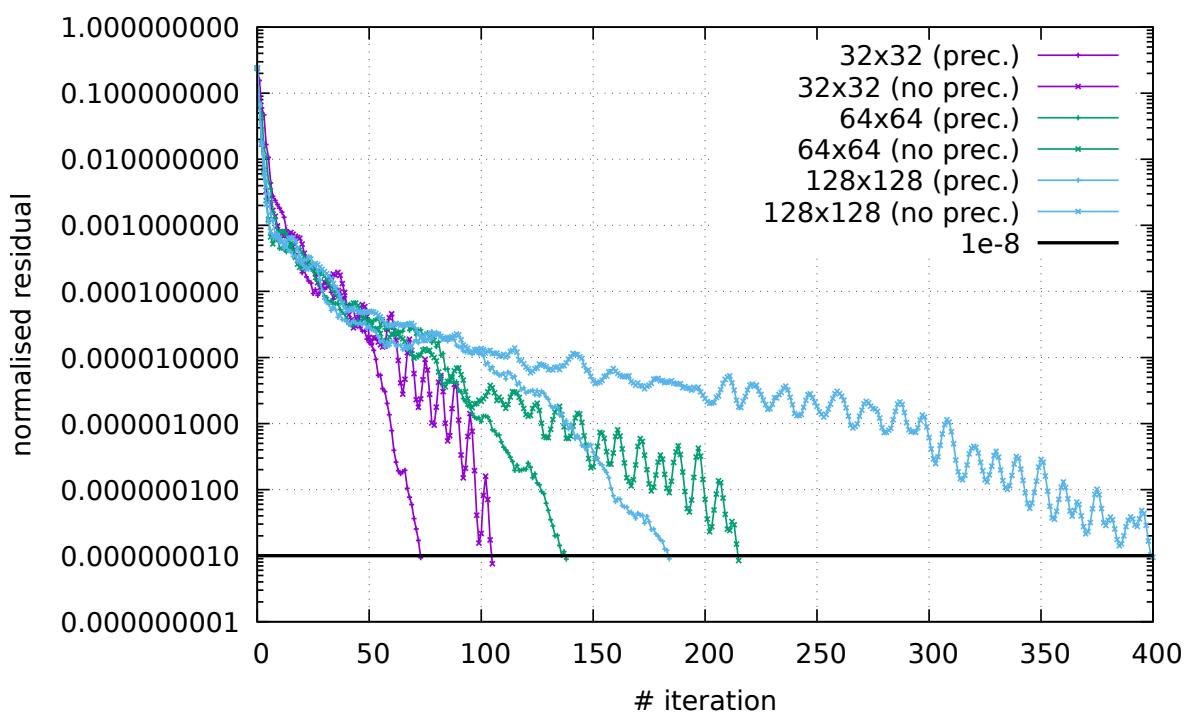
build S and have python compute its smallest and largest eigenvalues as a function of resolution?

## 23 fieldstone\_f16: saddle point problem with Schur complement approach - Stokes sphere

We are revisiting the 2D Stokes sphere problem, but this time we use the Schur complement approach to solve the Stokes system. Because there are viscosity contrasts in the domain, it is advisable to use the Preconditioned Conjugate Gradient as presented in Section ?? (see `solver_pcg`).



The normalised residual (see `solver_pcg`) was recorded. We see that all things equal the resolution has a strong influence on the number of iterations the solver must perform to reach the required tolerance. However, we see that the use of the preconditioner can substantially reduce the number of iterations inside the Stokes solver. At resolution 128x128, this number is halved.



#### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Schur complement approach
- isothermal
- non-isoviscous
- Stokes sphere

## 24 fieldstone\_17: solving the full saddle point problem in 3D

When using  $Q_1 \times P_0$  elements, this benchmark fails because of the Dirichlet b.c. on all 6 sides and all three components. However, as we will see, it does work well with  $Q_2 \times Q_1$  elements. .

This benchmark begins by postulating a polynomial solution to the 3D Stokes equation [164]:

$$\mathbf{v} = \begin{pmatrix} x + x^2 + xy + x^3y \\ y + xy + y^2 + x^2y^2 \\ -2z - 3xz - 3yz - 5x^2yz \end{pmatrix} \quad (422)$$

and

$$p = xyz + x^3y^3z - 5/32 \quad (423)$$

While it is then trivial to verify that this velocity field is divergence-free, the corresponding body force of the Stokes equation can be computed by inserting this solution into the momentum equation with a given viscosity  $\mu$  (constant or position/velocity/strain rate dependent). The domain is a unit cube and velocity boundary conditions simply use Eq. (422). Following [113], the viscosity is given by the smoothly varying function

$$\mu = \exp(1 - \beta(x(1-x) + y(1-y) + z(1-z))) \quad (424)$$

One can easily show that the ratio of viscosities  $\mu^*$  in the system follows  $\mu^* = \exp(-3\beta/4)$  so that choosing  $\beta = 10$  yields  $\mu^* \simeq 1808$  and  $\beta = 20$  yields  $\mu^* \simeq 3.269 \times 10^6$ .

We start from the momentum conservation equation:

$$-\nabla p + \nabla \cdot (2\mu \dot{\epsilon}) = \mathbf{f}$$

The  $x$ -component of this equation writes

$$f_x = -\frac{\partial p}{\partial x} + \frac{\partial}{\partial x}(2\mu \dot{\epsilon}_{xx}) + \frac{\partial}{\partial y}(2\mu \dot{\epsilon}_{xy}) + \frac{\partial}{\partial z}(2\mu \dot{\epsilon}_{xz}) \quad (425)$$

$$= -\frac{\partial p}{\partial x} + 2\mu \frac{\partial}{\partial x} \dot{\epsilon}_{xx} + 2\mu \frac{\partial}{\partial y} \dot{\epsilon}_{xy} + 2\mu \frac{\partial}{\partial z} \dot{\epsilon}_{xz} + 2 \frac{\partial \mu}{\partial x} \dot{\epsilon}_{xx} + 2 \frac{\partial \mu}{\partial y} \dot{\epsilon}_{xy} + 2 \frac{\partial \mu}{\partial z} \dot{\epsilon}_{xz} \quad (426)$$

Let us compute all the block separately:

$$\begin{aligned} \dot{\epsilon}_{xx} &= 1 + 2x + y + 3x^2y \\ \dot{\epsilon}_{yy} &= 1 + x + 2y + 2x^2y \\ \dot{\epsilon}_{zz} &= -2 - 3x - 3y - 5x^2y \\ 2\dot{\epsilon}_{xy} &= (x + x^3) + (y + 2xy^2) = x + y + 2xy^2 + x^3 \\ 2\dot{\epsilon}_{xz} &= (0) + (-3z - 10xyz) = -3z - 10xyz \\ 2\dot{\epsilon}_{yz} &= (0) + (-3z - 5x^2z) = -3z - 5x^2z \end{aligned}$$

In passing, one can verify that  $\dot{\epsilon}_{xx} + \dot{\epsilon}_{yy} + \dot{\epsilon}_{zz} = 0$ . We further have

$$\begin{aligned}\frac{\partial}{\partial x} 2\dot{\epsilon}_{xx} &= 2(2 + 6xy) \\ \frac{\partial}{\partial y} 2\dot{\epsilon}_{xy} &= 1 + 4xy \\ \frac{\partial}{\partial z} 2\dot{\epsilon}_{xz} &= -3 - 10xy \\ \frac{\partial}{\partial x} 2\dot{\epsilon}_{xy} &= 1 + 2y^2 + 3x^2 \\ \frac{\partial}{\partial y} 2\dot{\epsilon}_{yy} &= 2(2 + 2x^2) \\ \frac{\partial}{\partial z} 2\dot{\epsilon}_{yz} &= -3 - 5x^2 \\ \frac{\partial}{\partial x} 2\dot{\epsilon}_{xz} &= -10yz \\ \frac{\partial}{\partial y} 2\dot{\epsilon}_{yz} &= 0 \\ \frac{\partial}{\partial z} 2\dot{\epsilon}_{zz} &= 2(0)\end{aligned}$$

$$\frac{\partial p}{\partial x} = yz + 3x^2y^3z \quad (427)$$

$$\frac{\partial p}{\partial y} = xz + 3x^3y^2z \quad (428)$$

$$\frac{\partial p}{\partial z} = xy + x^3y^3 \quad (429)$$

**Pressure normalisation** Here again, because Dirichlet boundary conditions are prescribed on all sides the pressure is known up to an arbitrary constant. This constant can be determined by (arbitrarily) choosing to normalised the pressure field as follows:

$$\int_{\Omega} p \, d\Omega = 0 \quad (430)$$

This is a single constraint associated to a single Lagrange multiplier  $\lambda$  and the global Stokes system takes the form

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} & 0 \\ \mathbb{G}^T & 0 & \mathcal{C} \\ 0 & \mathcal{C}^T & 0 \end{pmatrix} \begin{pmatrix} V \\ P \\ \lambda \end{pmatrix}$$

In this particular case the constraint matrix  $\mathcal{C}$  is a vector and it only acts on the pressure degrees of freedom because of Eq.(430). Its exact expression is as follows:

$$\int_{\Omega} p \, d\Omega = \sum_e \int_{\Omega_e} p \, d\Omega = \sum_e \int_{\Omega_e} \sum_i N_i^p p_i \, d\Omega = \sum_e \sum_i \left( \int_{\Omega_e} N_i^p \, d\Omega \right) p_i = \sum_e \mathcal{C}_e \cdot \mathbf{p}_e$$

where  $\mathbf{p}_e$  is the list of pressure dofs of element  $e$ . The elemental constraint vector contains the corresponding pressure basis functions integrated over the element. These elemental constraints are then assembled into the vector  $\mathcal{C}$ .

#### 24.0.1 Constant viscosity

Choosing  $\beta = 0$  yields a constant velocity  $\mu(x, y, z) = \exp(1) \simeq 2.718$  (and greatly simplifies the right-hand side) so that

$$\frac{\partial}{\partial x} \mu(x, y, z) = 0 \quad (431)$$

$$\frac{\partial}{\partial y} \mu(x, y, z) = 0 \quad (432)$$

$$\frac{\partial}{\partial z} \mu(x, y, z) = 0 \quad (433)$$

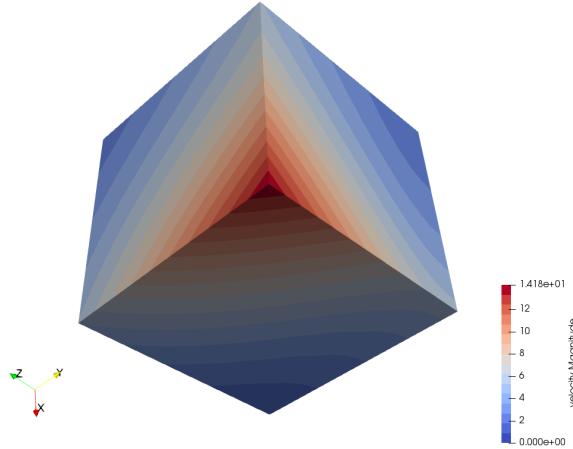
and

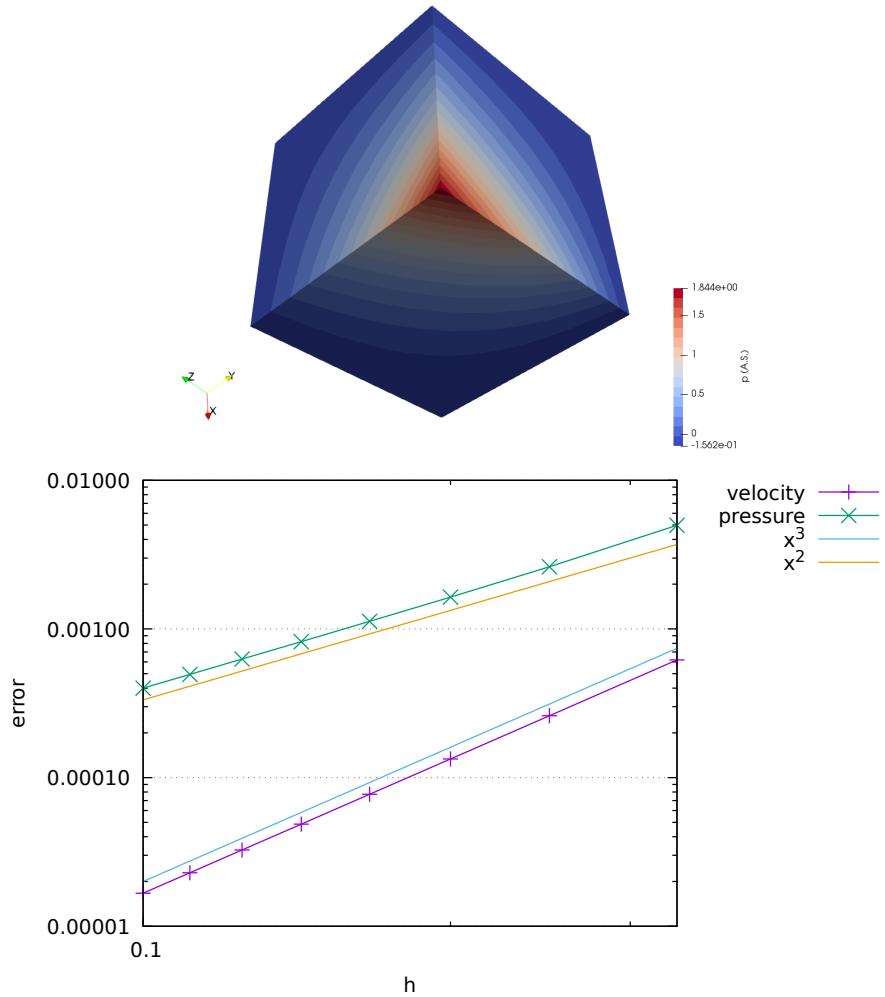
$$\begin{aligned} f_x &= -\frac{\partial p}{\partial x} + 2\mu \frac{\partial}{\partial x} \dot{\epsilon}_{xx} + 2\mu \frac{\partial}{\partial y} \dot{\epsilon}_{xy} + 2\mu \frac{\partial}{\partial z} \dot{\epsilon}_{xz} \\ &= -(yz + 3x^2y^3z) + 2(2 + 6xy) + (1 + 4xy) + (-3 - 10xy) \\ &= -(yz + 3x^2y^3z) + \mu(2 + 6xy) \\ f_y &= -\frac{\partial p}{\partial y} + 2\mu \frac{\partial}{\partial x} \dot{\epsilon}_{xy} + 2\mu \frac{\partial}{\partial y} \dot{\epsilon}_{yy} + 2\mu \frac{\partial}{\partial z} \dot{\epsilon}_{yz} \\ &= -(xz + 3x^3y^2z) + \mu(1 + 2y^2 + 3x^2) + \mu(2(2 + 2x^2) + \mu(-3 - 5x^2)) \\ &= -(xz + 3x^3y^2z) + \mu(2 + 2x^2 + 2y^2) \\ f_z &= -\frac{\partial p}{\partial z} + 2\mu \frac{\partial}{\partial x} \dot{\epsilon}_{xz} + 2\mu \frac{\partial}{\partial y} \dot{\epsilon}_{yz} + 2\mu \frac{\partial}{\partial z} \dot{\epsilon}_{zz} \\ &= -(xy + x^3y^3) + \mu(-10yz) + 0 + 0 \\ &= -(xy + x^3y^3) + \mu(-10yz) \end{aligned}$$

Finally

$$\mathbf{f} = - \begin{pmatrix} yz + 3x^2y^3z \\ xz + 3x^3y^2z \\ xy + x^3y^3 \end{pmatrix} + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix}$$

Note that there seems to be a sign problem with Eq.(26) in [113].





#### 24.0.2 Variable viscosity

The spatial derivatives of the viscosity are then given by

$$\begin{aligned}\frac{\partial}{\partial x} \mu(x, y, z) &= -(1 - 2x)\beta\mu(x, y, z) \\ \frac{\partial}{\partial y} \mu(x, y, z) &= -(1 - 2y)\beta\mu(x, y, z) \\ \frac{\partial}{\partial z} \mu(x, y, z) &= -(1 - 2z)\beta\mu(x, y, z)\end{aligned}$$

and the right-hand side by

$$\begin{aligned}
\mathbf{f} &= - \begin{pmatrix} yz + 3x^2y^3z \\ xz + 3x^3y^2z \\ xy + x^3y^3 \end{pmatrix} + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix} \\
&\quad - (1 - 2x)\beta\mu(x, y, z) \begin{pmatrix} 2\dot{\epsilon}_{xx} \\ 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{xz} \end{pmatrix} - (1 - 2y)\beta\mu(x, y, z) \begin{pmatrix} 2\dot{\epsilon}_{xy} \\ 2\dot{\epsilon}_{yy} \\ 2\dot{\epsilon}_{yz} \end{pmatrix} - (1 - 2z)\beta\mu(x, y, z) \begin{pmatrix} 2\dot{\epsilon}_{xz} \\ 2\dot{\epsilon}_{yz} \\ 2\dot{\epsilon}_{zz} \end{pmatrix} \\
&= - \begin{pmatrix} yz + 3x^2y^3z \\ xz + 3x^3y^2z \\ xy + x^3y^3 \end{pmatrix} + \mu \begin{pmatrix} 2 + 6xy \\ 2 + 2x^2 + 2y^2 \\ -10yz \end{pmatrix} \\
&\quad - (1 - 2x)\beta\mu \begin{pmatrix} 2 + 4x + 2y + 6x^2y \\ x + y + 2xy^2 + x^3 \\ -3z - 10xyz \end{pmatrix} - (1 - 2y)\beta\mu \begin{pmatrix} x + y + 2xy^2 + x^3 \\ 2 + 2x + 4y + 4x^2y \\ -3z - 5x^2z \end{pmatrix} - (1 - 2z)\beta\mu \begin{pmatrix} -3z - 10xyz \\ -3z - 5x^2z \\ -4 - 6x - 6y - 10xz \end{pmatrix}
\end{aligned}$$

Note that at  $(x, y, z) = (0, 0, 0)$ ,  $\mu = \exp(1)$ , and at  $(x, y, z) = (0.5, 0.5, 0.5)$ ,  $\mu = \exp(1 - 3\beta/4)$  so that the maximum viscosity ratio is given by

$$\mu^* = \frac{\exp(1 - 3\beta/4)}{\exp(1)} = \exp(-3\beta/4)$$

By varying  $\beta$  between 1 and 22 we can get up to 7 orders of magnitude viscosity difference.

#### features

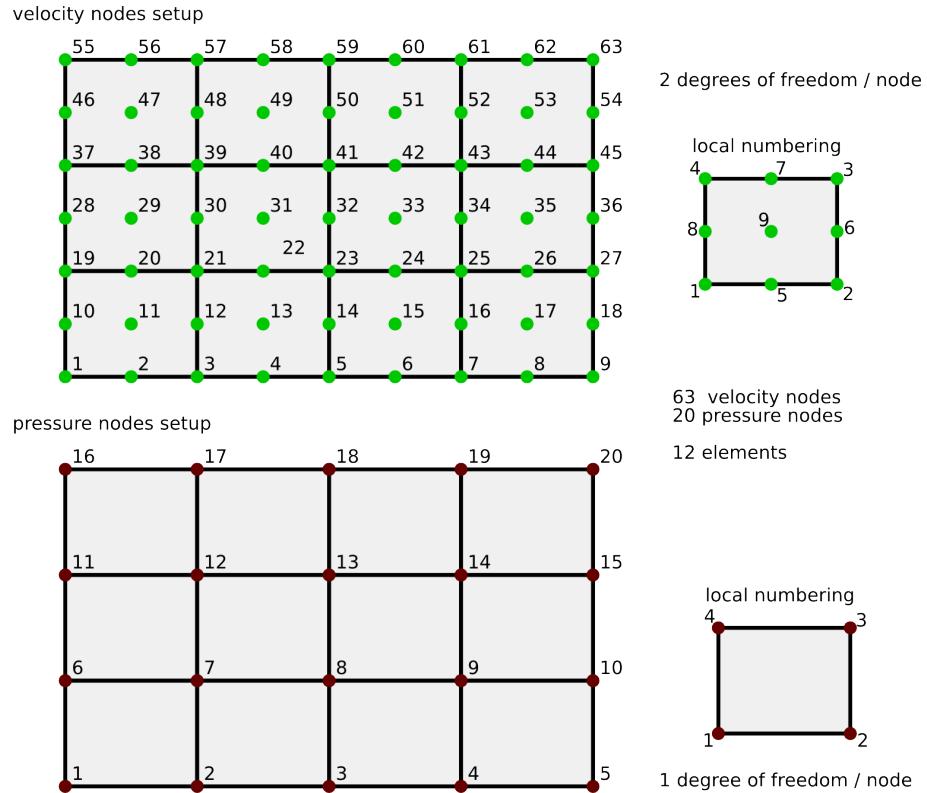
- $Q_1 \times P_0$  element
- incompressible flow
- saddle point system
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- 3D
- elemental b.c.
- analytical solution

## 25 fieldstone\_18: solving the full saddle point problem with $Q_2 \times Q_1$ elements

The details of the numerical setup are presented in Section ??.

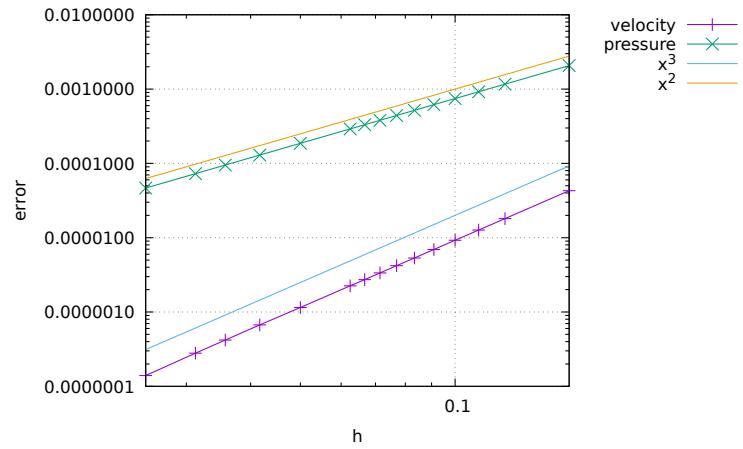
Each element has  $m_V = 9$  vertices so in total  $ndof_V \times m_V = 18$  velocity dofs and  $ndof_P \times m_P = 4$  pressure dofs. The total number of velocity dofs is therefore  $NfemV = nnp \times ndofV$  while the total number of pressure dofs is  $NfemP = nel$ . The total number of dofs is then  $Nfem = NfemV + NfemP$ .

As a consequence, matrix  $\mathbb{K}$  has size  $NfemV, NfemV$  and matrix  $\mathbb{G}$  has size  $NfemV, NfemP$ . Vector  $f$  is of size  $NfemV$  and vector  $h$  is of size  $NfemP$ .



### features

- $Q_2 \times Q_1$  element
- incompressible flow
- mixed formulation
- Dirichlet boundary conditions (no-slip)
- isothermal
- isoviscous
- analytical solution



## 26 fieldstone\_19: solving the full saddle point problem with $Q_3 \times Q_2$ elements

The details of the numerical setup are presented in Section ??.

Each element has  $m_V = 16$  vertices so in total  $ndof_V \times m_V = 32$  velocity dofs and  $ndof_P \times m_P = 9$  pressure dofs. The total number of velocity dofs is therefore  $NfemV = nnp \times ndofV$  while the total number of pressure dofs is  $NfemP = nel$ . The total number of dofs is then  $Nfem = NfemV + NfemP$ .

As a consequence, matrix  $\mathbb{K}$  has size  $NfemV, NfemV$  and matrix  $\mathbb{G}$  has size  $NfemV, NfemP$ . Vector  $f$  is of size  $NfemV$  and vector  $h$  is of size  $NfemP$ .

```

60====61====62====63====64====65====66====67====68====70
|| || || || ||
50 51 52 53 54 55 56 57 58 59
|| || || || ||
40 41 42 43 44 45 46 47 48 49
|| || || || ||
30====31====32====33====34====35====36====37====38====39
|| || || || ||
20 21 22 23 24 25 26 27 28 29
|| || || || ||
10 11 12 13 14 15 16 17 18 19
|| || || || ||
00====01====02====03====04====05====06====07====08====09

```

Example of 3x2 mesh.  $n_{nx}=10$ ,  $n_{ny}=7$ ,  $nnp=70$ ,  $nelx=3$ ,  $nely=2$ ,  $nel=6$

```

12====13====14====15 06=====07=====08
|| || || || || || || || ||
08====09====10====11 || || || ||
|| || || || 03=====04=====05
04====05====06====07 || || || ||
|| || || || 00=====01=====02
00====01====02====03

```

Velocity (Q3)

```

(r,s)_{00}=(-1,-1) (r,s)_{00}=(-1,-1)
(r,s)_{01}=(-1/3,-1) (r,s)_{01}=(0,-1)
(r,s)_{02}=(+1/3,-1) (r,s)_{02}=(+1,-1)
(r,s)_{03}=(+1,-1) (r,s)_{03}=(-1,0)
(r,s)_{04}=(-1,-1/3) (r,s)_{04}=(0,0)
(r,s)_{05}=(-1/3,-1/3) (r,s)_{05}=(+1,0)
(r,s)_{06}=(+1/3,-1/3) (r,s)_{06}=(-1,+1)
(r,s)_{07}=(+1,-1/3) (r,s)_{07}=(0,+1)
(r,s)_{08}=(-1,+1/3) (r,s)_{08}=(+1,+1)
(r,s)_{09}=(-1/3,+1/3)
(r,s)_{10}=(+1/3,+1/3)
(r,s)_{11}=(+1,+1/3)
(r,s)_{12}=(-1,+1)
(r,s)_{13}=(-1/3,+1)
(r,s)_{14}=(+1/3,+1)
(r,s)_{15}=(+1,+1)

```

Pressure (Q2)

```

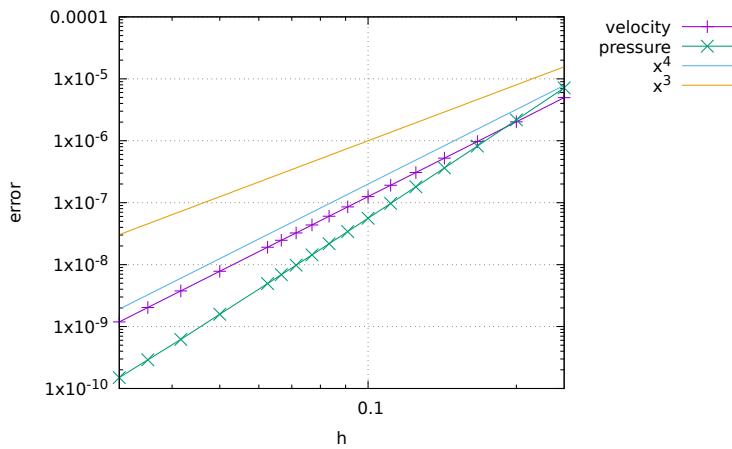
(r,s)_{00}=(-1,-1)
(r,s)_{01}=(0,-1)
(r,s)_{02}=(+1,-1)
(r,s)_{03}=(-1,0)
(r,s)_{04}=(0,0)
(r,s)_{05}=(+1,0)
(r,s)_{06}=(-1,+1)
(r,s)_{07}=(0,+1)
(r,s)_{08}=(+1,+1)

```

Write about 4 point quadrature.

**features**

- $Q_3 \times Q_2$  element
- incompressible flow
- mixed formulation
- isothermal
- isoviscous
- analytical solution



velocity error rate is cubic, pressure superconvergent since the pressure field is quadratic and therefore lies into the  $Q_2$  space.

## 27 fieldstone\_20: the Busse benchmark

This three-dimensional benchmark was first proposed by [114]. It has been subsequently presented in [528, 553, 8, 434, 157, 358]. We here focus on Case 1 of [114]: an isoviscous bimodal convection experiment at  $Ra = 3 \times 10^5$ .

The domain is of size  $a \times b \times h$  with  $a = 1.0079h$ ,  $b = 0.6283h$  with  $h = 2700\text{km}$ . It is filled with a Newtonian fluid characterised by  $\rho_0 = 3300\text{kg.m}^{-3}$ ,  $\alpha = 10^{-5}\text{K}^{-1}$ ,  $\mu = 8.0198 \times 10^{23}\text{Pa.s}$ ,  $k = 3.564\text{W.m}^{-1}\text{.K}^{-1}$ ,  $c_p = 1080\text{J.K}^{-1}\text{.kg}^{-1}$ . The gravity vector is set to  $\mathbf{g} = (0, 0, -10)^T$ . The temperature is imposed at the bottom ( $T = 3700^\circ\text{C}$ ) and at the top ( $T = 0^\circ\text{C}$ ).

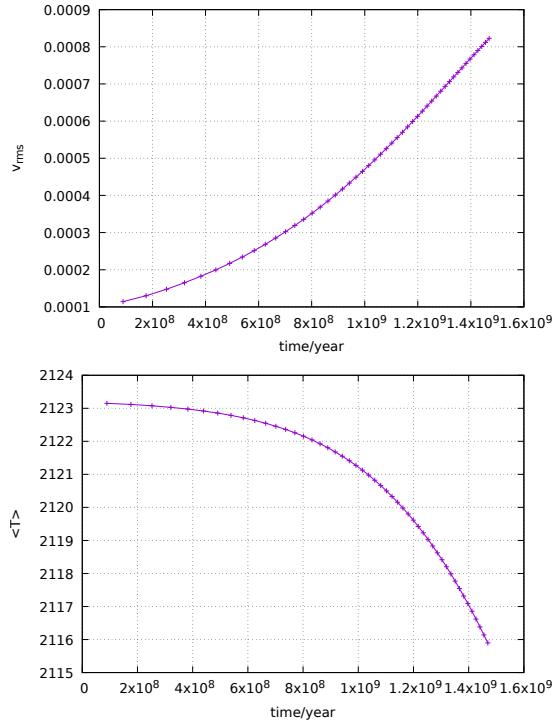
The various measurements presented in [114] are listed hereafter:

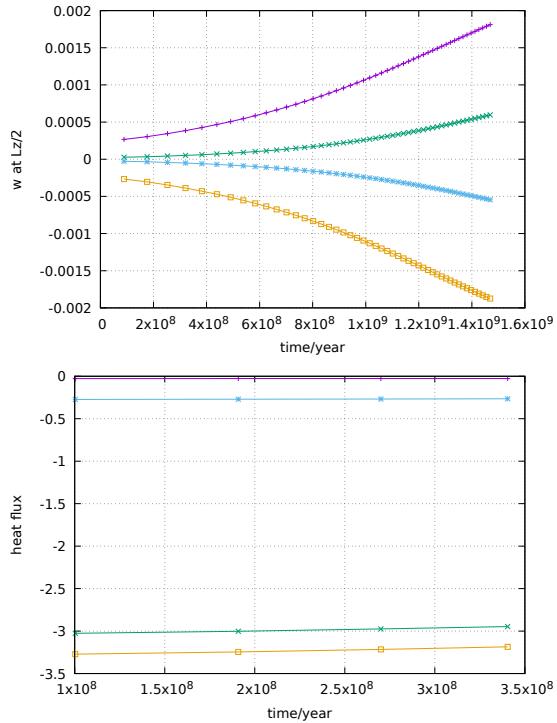
- The Nusselt number  $Nu$  computed at the top surface following Eq. (399):

$$Nu = L_z \frac{\int \int_{z=L_z} \frac{\partial T}{\partial y} dx dy}{\int \int_{z=0} T dx dy}$$

- the root mean square velocity  $v_{rms}$  and the temperature mean square velocity  $T_{rms}$
- The vertical velocity  $w$  and temperature  $T$  at points  $\mathbf{x}_1 = (0, 0, L_z/2)$ ,  $\mathbf{x}_2 = (L_x, 0, L_z/2)$ ,  $\mathbf{x}_3 = (0, L_y, L_z/2)$  and  $\mathbf{x}_4 = (L_x, L_y, L_z/2)$ ;
- the vertical component of the heat flux  $Q$  at the top surface at all four corners.

The values plotted hereunder are adimensionalised by means of a reference temperature (3700K), a reference lengthscale 2700km, and a reference time  $L_z^2/\kappa$ .





### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- 3D
- elemental b.c.
- buoyancy driven

**ToDo:** look at energy conservation. run to steady state and make sure the expected values are retrieved.

## 28 fieldstone\_21: The non-conforming $Q_1 \times P_0$ element

### features

- Non-conforming  $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- isothermal
- non-isoviscous
- analytical solution
- pressure smoothing

try Q1 mapping instead of isoparametric.

## 29 fieldstone\_22: The stabilised $Q_1 \times Q_1$ element

The details of the numerical setup are presented in Section 7.4.

We wish to use  $Q_1 \times Q_1$  element, which, unless stabilised, violates the LBB stability condition and therefore is unusable. Stabilisation can be of two types: least-squares [165, 539, 342, 285], or by means of an additional term in the weak form as first introduced in [164, 65], which is appealing since there is no explicit stabilisation parameter. It is further analysed in [431, 373, 310, 499, 276]. Note that an equal-order velocity-pressure formulation that does not exhibit spurious pressure modes (without stabilisation) has been presented in [479].

This element corresponds to bilinear velocities, bilinear pressure (equal order interpolation for both velocity and pressure) which is very convenient in terms of data structures since all dofs are colocated.

In geodynamics, it is used in the Rhea code [518, 113] and in Gale [23]. It is also used in [370] in its stabilised form, in conjunction with AMR. This element is quickly discussed at page 217 of Volker John's book [330].

The stabilisation term  $\mathbb{C}$  enters the Stokes matrix in the (2,2) position:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T & -\mathbb{C} \end{pmatrix} \cdot \begin{pmatrix} \mathcal{V} \\ \mathcal{P} \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix}$$

The purpose of the  $\mathbb{C}$  term is to stabilise the linear system. It is given by:

$$\mathbb{C}(p, q) = \sum_e \int_{\Omega_e} \frac{1}{\eta} (p - \Pi p)(q - \Pi q) d\Omega$$

where  $\Pi$  is the  $L^2$ -projection onto the space of element-wise constant functions:

$$\Pi p = \frac{1}{|\Omega_e|} \int_{\Omega_e} p d\Omega$$

Because of the stabilisation matrix  $\mathbb{C}$ , the numerical solution satisfies the incompressibility condition only approximately. Local mesh refinement helps to control these unwanted effects [112, 113]. Since  $\mathbb{K}$  and  $\mathbb{C}$  are symmetric matrices, the Stokes system is then an indefinite symmetric system. The Schur complement matrix  $\mathbb{S}$  is then given by

$$\mathbb{S} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}$$

One can further expand the above expression for the  $\mathbb{C}$  term:

$$\begin{aligned} \mathbb{C}(p, q) &= \sum_e \int_{\Omega_e} \frac{1}{\eta} (p - \Pi p)(q - \Pi q) d\Omega \\ &= \sum_e \int_{\Omega_e} \frac{1}{\eta} [pq - (\Pi p)q - (\Pi q)p + (\Pi p)(\Pi q)] d\Omega \\ &= \sum_e \frac{1}{\eta_e} \left[ \int_{\Omega_e} pq d\Omega - \int_{\Omega_e} (\Pi p)qd\Omega - \int_{\Omega_e} (\Pi q)pd\Omega + \int_{\Omega_e} (\Pi p)(\Pi q)d\Omega \right] \\ &= \sum_e \frac{1}{\eta_e} \left[ \int_{\Omega_e} pq d\Omega - (\Pi p) \int_{\Omega_e} q d\Omega - (\Pi q) \int_{\Omega_e} p d\Omega + (\Pi p)(\Pi q) \int_{\Omega_e} d\Omega \right] \\ &= \sum_e \frac{1}{\eta_e} \left[ \int_{\Omega_e} pq d\Omega - (\Pi p)|\Omega_e|(\Pi q) - (\Pi q)|\Omega_e|(\Pi p) + (\Pi p)(\Pi q)|\Omega_e| \right] \\ &= \sum_e \frac{1}{\eta_e} \left[ \int_{\Omega_e} pq d\Omega - |\Omega_e|(\Pi p)(\Pi q) \right] \end{aligned} \tag{435}$$

where we have used the fact that on each element  $\Pi p^h$  is constant. The left term will obviously yield a  $Q_1$  mass matrix (scaled by the elemental viscosities). Note that this approach is not used in practice as we'll see hereafter.

The pressure inside an element is given by

$$p^h(\vec{x}) = \sum_k N_k^p(\vec{x}) p_k$$

so that

$$\Pi p^h = \frac{1}{|\Omega_e|} \int_{\Omega_e} \sum_k N_k^p p_k d\Omega = \sum_k \left( \underbrace{\frac{1}{|\Omega_e|} \int_{\Omega_e} N_k^p d\Omega}_{\tilde{N}_k^p} \right) p_k \quad (436)$$

and then

$$p^h - \Pi p^h = \sum_k N_k^p(\vec{x}) p_k - \sum_k \tilde{N}_k^p p_k = \sum_k (N_k^p(\vec{x}) - \tilde{N}_k^p) p_k$$

The algorithm is straightforward and as follows: In the loop over elements, a) Compute the average of each shape function  $N_k^p(\vec{x})$  over the element; b) Subtract this average to the shape function; c) Build mass matrix with modified/offset shape functions (taking in account the viscosity).

In the case of rectangular elements of size  $(h_x, h_y)$ ,  $\tilde{N}_k^p$  simplifies even more:

$$\tilde{N}_k^p = \frac{1}{|\Omega_e|} \int_{\Omega_e} N_k^p(\vec{x}) d\Omega = \frac{1}{h_x h_y} \frac{h_x h_y}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_k^p(r, s) dr ds = \frac{1}{4} \int_{-1}^{+1} \int_{-1}^{+1} N_k^p(r, s) dr ds \quad (437)$$

It is easy to show that the average of the  $Q_1$  shape functions over the reference element is 1, so that  $\tilde{N}_k^p = 1/4$ . This explains why in the code we have:

```
Navrg = np.zeros(m, dtype=np.float64)
Navrg[0]=0.25
Navrg[1]=0.25
Navrg[2]=0.25
Navrg[3]=0.25
```

This also means that  $\Pi p^h = (p_1 + p_2 + p_3 + p_4)/4$ , i.e. the projected pressure is the mean of the vertex values. It follows, as shown on p.244 of [182] that the elemental  $\mathbb{C}$  matrix is (omitting the viscosity term)

$$\mathbb{C}_{el} = \mathbb{M}_{el} - \vec{q}^T \vec{q} |\Omega_e| \quad \vec{q} = \left( \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right)$$

The nullspace of  $\mathbb{C}$  consists of constant vectors, i.e.  $\vec{1} \in \text{null}(\mathbb{C})$  which means that the assembled stabilisation operator is consistent.

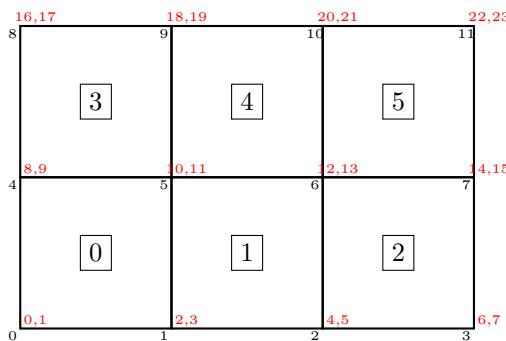
The elemental  $\mathbb{C}_{el}$  matrix is then computed like a mass matrix, although with modified shape function vectors. Inside the loop over quadrature points, we do:

```
Nvect[0:,0:m]=N[0:m]-Navrg[0:m]
C_el+=Nvect.T.dot(Nvect)*jacob*weightq/viscosity(xq,yq,case)
```

It is then assembled inside the big FEM matrix

```
for k1 in range(0,m):
 for k2 in range(0,m):
 C_mat[icon[k1, ie1], icon[k2, ie1]]+=C_el[k1, k2]
```

Non-zero pattern of the  $\mathbb{G}$  matrix: Let us take a simple example: a 3x2 element grid.



The  $\mathbb{K}$  matrix is of size  $NfemV \times NfemV$  with  $NfemV = ndofV \times nnp = 2 \times 12 = 24$ . The  $\mathbb{G}$  matrix is of size  $NfemV \times NfemP$  with  $NfemP = ndofP \times nnp = 1 \times 12 = 12$ . The  $\mathbb{C}$  matrix is of size  $NfemP \times NfemP$ .

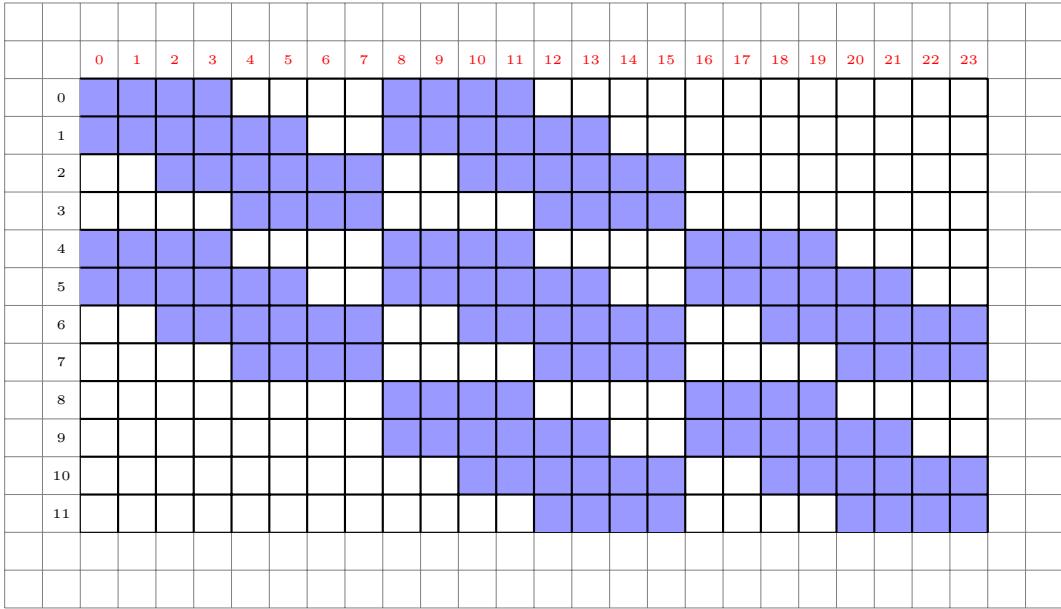
A corner pdof sees 4 vdofs, a side pdof sees 12 vdofs and an inside pdof sees 18 vdofs, so that the total number of nonzeros in  $\mathbb{G}$  can be computed as follows:

$$NZ_{\mathbb{G}} = \underbrace{4}_{\text{corners}} + \underbrace{2(nnx - 2) * 12}_{\text{2hor.sides}} + \underbrace{2(nny - 2) * 12}_{\text{2vert.sides}} + \underbrace{(nnx - 2)(nny - 2) * 18}_{\text{insidenodes}}$$

Concretely,

- pdof #0 sees vdofs 0,1,2,3,8,9,10,11
- pdof #1 sees vdofs 0,1,2,3,4,5,8,9,10,11,12,13
- pdof #5 sees vdofs 0,1,2,3,4,5,8,9,10,11,12,13,16,17,18,19,20,21

so that the  $\mathbb{G}^T$  matrix non-zero structure then is as follows:



Non-zero pattern of the  $\mathbb{C}$  matrix: Let us take a simple example: a 3x2 element grid.

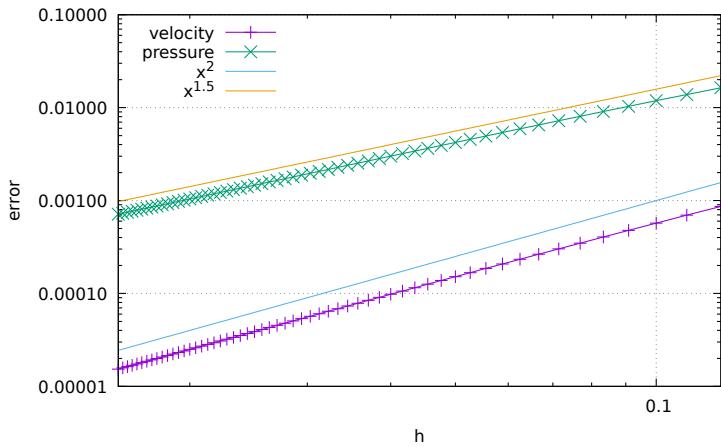
finish structure of C matrix for q1q1

We impose  $\int pdV = 0$  which means that the following constraint is added to the Stokes matrix:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} & 0 \\ \mathbb{G}^T & \mathbb{C} & \mathbb{L} \\ 0 & \mathbb{L}^T & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathcal{V} \\ \mathcal{P} \\ \lambda \end{pmatrix} = \begin{pmatrix} f \\ h \\ 0 \end{pmatrix}$$

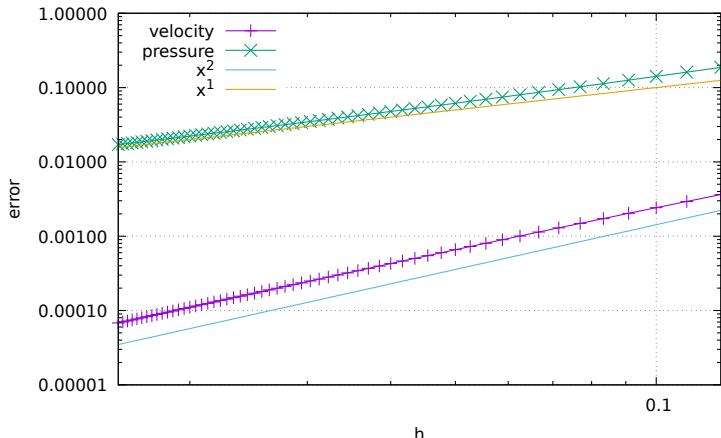
## 29.1 The Donea & Huerta benchmark

As in [165] we solve the benchmark problem presented in section 7.4.1.



## 29.2 The Dohrmann & Bochev benchmark

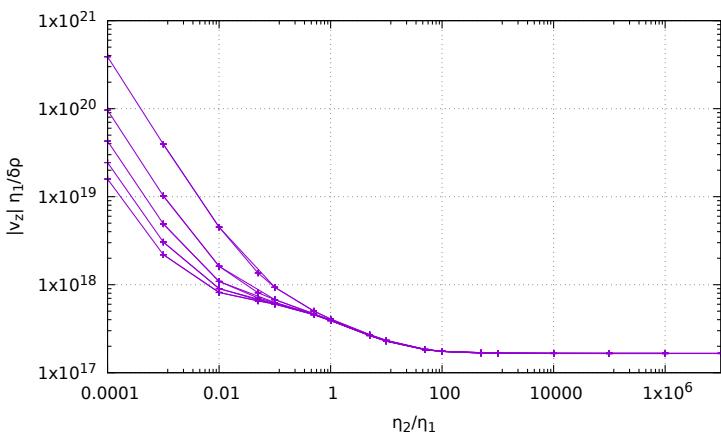
As in [164] we solve the benchmark problem presented in section 7.4.2.



compare my rates with original paper!

## 29.3 The falling block experiment

The setup is described in [545].



## 30 fieldstone\_23: compressible flow (1) - analytical benchmark

This work is part of the MSc thesis of T. Weir (2018).

We first start with an isothermal Stokes flow, so that we disregard the heat transport equation and the equations we wish to solve are simply:

$$-\nabla \cdot \left[ 2\eta \left( \dot{\epsilon}(\mathbf{v}) - \frac{1}{3}(\nabla \cdot \mathbf{v})\mathbf{1} \right) \right] + \nabla p = \rho g \quad \text{in } \Omega, \quad (438)$$

$$\nabla \cdot (\rho \mathbf{v}) = 0 \quad \text{in } \Omega \quad (439)$$

The second equation can be rewritten  $\nabla \cdot (\rho \mathbf{v}) = \rho \nabla \cdot \mathbf{v} + \mathbf{v} \cdot \nabla \rho = 0$  or,

$$\nabla \cdot \mathbf{v} + \frac{1}{\rho} \mathbf{v} \cdot \nabla \rho = 0$$

Note that this presupposes that the density is not zero anywhere in the domain.

We use a mixed formulation and therefore keep both velocity and pressure as unknowns. We end up having to solve the following system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} \\ \mathbb{G}^T + \mathbb{Z} & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathcal{V} \\ \mathcal{P} \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix} \quad \text{or,} \quad \mathbb{A} \cdot X = rhs$$

Where  $\mathbb{K}$  is the stiffness matrix,  $\mathbb{G}$  is the discrete gradient operator,  $\mathbb{G}^T$  is the discrete divergence operator,  $\mathcal{V}$  the velocity vector,  $\mathcal{P}$  the pressure vector. Note that the term  $\mathbb{Z}\mathcal{V}$  derives from term  $\mathbf{v} \cdot \nabla \rho$  in the continuity equation.

Each block  $\mathbb{K}$ ,  $\mathbb{G}$ ,  $\mathbb{Z}$  and vectors  $f$  and  $h$  are built separately in the code and assembled into the matrix  $\mathbb{A}$  and vector  $rhs$  afterwards.  $\mathbb{A}$  and  $rhs$  are then passed to the solver. We will see later that there are alternatives to solve this approach which do not require to build the full Stokes matrix  $\mathbb{A}$ .

*Remark:* the term  $\mathbb{Z}\mathcal{V}$  is often put in the rhs (i.e. added to  $h$ ) so that the matrix  $\mathbb{A}$  retains the same structure as in the incompressible case. This is indeed how it is implemented in ASPECT. This however requires more work since the rhs depends on the solution and some form of iterations is needed.

In the case of a compressible flow the strain rate tensor and the deviatoric strain rate tensor are no more equal (since  $\nabla \cdot \mathbf{v} \neq 0$ ). The deviatoric strainrate tensor is given by<sup>30</sup>

$$\dot{\epsilon}^d(\mathbf{v}) = \dot{\epsilon}(\mathbf{v}) - \frac{1}{3} Tr(\dot{\epsilon})\mathbf{1} = \dot{\epsilon}(\mathbf{v}) - \frac{1}{3}(\nabla \cdot \mathbf{v})\mathbf{1}$$

In that case:

$$\dot{\epsilon}_{xx}^d = \frac{\partial u}{\partial x} - \frac{1}{3} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) = \frac{2}{3} \frac{\partial u}{\partial x} - \frac{1}{3} \frac{\partial v}{\partial y} \quad (440)$$

$$\dot{\epsilon}_{yy}^d = \frac{\partial v}{\partial y} - \frac{1}{3} \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) = -\frac{1}{3} \frac{\partial u}{\partial x} + \frac{2}{3} \frac{\partial v}{\partial y} \quad (441)$$

$$2\dot{\epsilon}_{xy}^d = \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \quad (442)$$

and then

$$\dot{\epsilon}^d(\mathbf{v}) = \begin{pmatrix} \frac{2}{3} \frac{\partial u}{\partial x} - \frac{1}{3} \frac{\partial v}{\partial y} & \frac{1}{2} \frac{\partial u}{\partial y} + \frac{1}{2} \frac{\partial v}{\partial x} \\ \frac{1}{2} \frac{\partial u}{\partial y} + \frac{1}{2} \frac{\partial v}{\partial x} & -\frac{1}{3} \frac{\partial u}{\partial x} + \frac{2}{3} \frac{\partial v}{\partial y} \end{pmatrix}$$

From  $\vec{\tau} = 2\eta \dot{\epsilon}^d$  we arrive at:

$$\begin{pmatrix} \tau_{xx} \\ \tau_{yy} \\ \tau_{xy} \end{pmatrix} = 2\eta \begin{pmatrix} \dot{\epsilon}_{xx}^d \\ \dot{\epsilon}_{yy}^d \\ \dot{\epsilon}_{xy}^d \end{pmatrix} = 2\eta \begin{pmatrix} 2/3 & -1/3 & 0 \\ -1/3 & 2/3 & 0 \\ 0 & 0 & 1/2 \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \end{pmatrix} = \eta \begin{pmatrix} 4/3 & -2/3 & 0 \\ -2/3 & 4/3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} \frac{\partial u}{\partial x} \\ \frac{\partial v}{\partial y} \\ \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \end{pmatrix}$$

or,

$$\vec{\tau} = \mathbf{C}_\eta \mathbf{B} \mathbf{V}$$

<sup>30</sup>See the ASPECT manual for a justification of the 3 value in the denominator in 2D and 3D.

In order to test our implementation we have created a few manufactured solutions:

- benchmark #1 (ibench=1): Starting from a density profile of:

$$\rho(x, y) = xy \quad (443)$$

We derive a velocity given by:

$$v_x(x, y) = \frac{C_x}{x}, v_y(x, y) = \frac{C_y}{y} \quad (444)$$

With  $g_x(x, y) = \frac{1}{x}$  and  $g_y(x, y) = \frac{1}{y}$ , this leads us to a pressure profile:

$$p = -\eta \left( \frac{4C_x}{3x^2} + \frac{4C_y}{3y^2} \right) + xy + C_0 \quad (445)$$

This gives us a strain rate:

$$\dot{\epsilon}_{xx} = \frac{-C_x}{x^2} \quad \dot{\epsilon}_{yy} = \frac{-C_y}{y^2} \quad \dot{\epsilon}_{xy} = 0$$

In what follows, we choose  $\eta = 1$  and  $C_x = C_y = 1$  and for a unit square domain  $[1 : 2] \times [1 : 2]$  we compute  $C_0$  so that the pressure is normalised to zero over the whole domain and obtain  $C_0 = -1$ .

- benchmark #2 (ibench=2): Starting from a density profile of:

$$\rho = \cos(x) \cos(y) \quad (446)$$

We derive a velocity given by:

$$v_x = \frac{C_x}{\cos(x)}, v_y = \frac{C_y}{\cos(y)} \quad (447)$$

With  $g_x = \frac{1}{\cos(y)}$  and  $g_y = \frac{1}{\cos(x)}$ , this leads us to a pressure profile:

$$p = \eta \left( \frac{4C_x \sin(x)}{3 \cos^2(x)} + \frac{4C_y \sin(y)}{3 \cos^2(y)} \right) + (\sin(x) + \sin(y)) + C_0 \quad (448)$$

$$\dot{\epsilon}_{xx} = C_x \frac{\sin(x)}{\cos^2(x)} \quad \dot{\epsilon}_{yy} = C_y \frac{\sin(y)}{\cos^2(y)} \quad \dot{\epsilon}_{xy} = 0$$

We choose  $\eta = 1$  and  $C_x = C_y = 1$ . The domain is the unit square  $[0 : 1] \times [0 : 1]$  and we obtain  $C_0$  as before and obtain

$$C_0 = 2 - 2 \cos(1) + 8/3(\frac{1}{\cos(1)} - 1) \simeq 3.18823730$$

(thank you WolframAlpha)

- benchmark #3 (ibench=3)
- benchmark #4 (ibench=4)
- benchmark #5 (ibench=5)

## features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- Dirichlet boundary conditions (no-slip)
- isothermal
- isoviscous
- analytical solution
- pressure smoothing

ToDo:

- pbs with odd vs even number of elements
- q is 'fine' everywhere except in the corners - revisit pressure smoothing paper?
- redo A v d Berg benchmark (see Tom Weir thesis)

## 31 fieldstone\_24: compressible flow (2) - convection box

This work is part of the MSc thesis of T. Weir (2018).

### 31.1 The physics

Let us start with some thermodynamics. Every material has an equation of state. The equilibrium thermodynamic state of any material can be constrained if any two state variables are specified. Examples of state variables include the pressure  $p$  and specific volume  $\nu = 1/\rho$ , as well as the temperature  $T$ .

After linearisation, the density depends on temperature and pressure as follows:

$$\rho(T, p) = \rho_0 ((1 - \alpha(T - T_0) + \beta_T p)$$

where  $\alpha$  is the coefficient of thermal expansion, also called thermal expansivity:

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

$\alpha$  is the percentage increase in volume of a material per degree of temperature increase; the subscript  $p$  means that the pressure is held fixed.

$\beta_T$  is the isothermal compressibility of the fluid, which is given by

$$\beta_T = \frac{1}{K} = \frac{1}{\rho} \left( \frac{\partial \rho}{\partial P} \right)_T$$

with  $K$  the bulk modulus. Values of  $\beta_T = 10^{-12} - 10^{-11}$  Pa $^{-1}$  are reasonable for Earth's mantle, with values decreasing by about a factor of 5 between the shallow lithosphere and core-mantle boundary. This is the percentage increase in density per unit change in pressure at constant temperature. Both the coefficient of thermal expansion and the isothermal compressibility can be obtained from the equation of state.

The full set of equations we wish to solve is given by

$$-\nabla \cdot [2\eta \dot{\epsilon}^d(\mathbf{v})] + \nabla p = \rho_0 ((1 - \alpha(T - T_0) + \beta_T p) \mathbf{g} \quad \text{in } \Omega \quad (449)$$

$$\nabla \cdot \mathbf{v} + \frac{1}{\rho} \mathbf{v} \cdot \nabla \rho = 0 \quad \text{in } \Omega \quad (450)$$

$$\rho C_p \left( \frac{\partial T}{\partial t} + \mathbf{v} \cdot \nabla T \right) - \nabla \cdot k \nabla T = \rho H + 2\eta \dot{\epsilon}^d : \dot{\epsilon}^d + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \quad \text{in } \Omega, \quad (451)$$

Note that this presupposes that the density is not zero anywhere in the domain.

### 31.2 The numerics

We use a mixed formulation and therefore keep both velocity and pressure as unknowns. We end up having to solve the following system:

$$\begin{pmatrix} \mathbb{K} & \mathbb{G} + \mathbb{W} \\ \mathbb{G}^T + \mathbb{Z} & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathcal{V} \\ \mathcal{P} \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix} \quad \text{or,} \quad \mathbb{A} \cdot X = rhs$$

Where  $\mathbb{K}$  is the stiffness matrix,  $\mathbb{G}$  is the discrete gradient operator,  $\mathbb{G}^T$  is the discrete divergence operator,  $\mathcal{V}$  the velocity vector,  $\mathcal{P}$  the pressure vector. Note that the term  $\mathbb{Z}\mathcal{V}$  derives from term  $\mathbf{v} \cdot \nabla \rho$  in the continuity equation.

As perfectly explained in the step 32 of deal.ii<sup>31</sup>, we need to scale the  $\mathbb{G}$  term since it is many orders of magnitude smaller than  $\mathbb{K}$ , which introduces large inaccuracies in the solving process to the point that the solution is nonsensical. This scaling coefficient is  $\eta/L$ . After building the  $\mathbb{G}$  block, it is then scaled as follows:  $\mathbb{G}' = \frac{\eta}{L} \mathbb{G}$  so that we now solve

<sup>31</sup>[https://www.dealii.org/9.0.0/doxygen/deal.II/step\\_32.html](https://www.dealii.org/9.0.0/doxygen/deal.II/step_32.html)

$$\begin{pmatrix} \mathbb{K} & \mathbb{G}' + \mathbb{W} \\ \mathbb{G}'^T + \mathbb{Z} & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathcal{V} \\ \mathcal{P}' \end{pmatrix} = \begin{pmatrix} f \\ h \end{pmatrix}$$

After the solve phase, we recover the real pressure with  $\mathcal{P} = \frac{\eta}{L}\mathcal{P}'$ .

adapt notes since I should scale  $\mathbb{W}$  and  $\mathbb{Z}$  too.  $h$  should be scaled too !!!!!!!

Each block  $\mathbb{K}$ ,  $\mathbb{G}$ ,  $\mathbb{Z}$  and vectors  $f$  and  $h$  are built separately in the code and assembled into the matrix  $\mathbb{A}$  and vector  $rhs$  afterwards.  $\mathbb{A}$  and  $rhs$  are then passed to the solver. We will see later that there are alternatives to solve this approach which do not require to build the full Stokes matrix  $\mathbb{A}$ .

**Remark 1:** the terms  $\mathbb{Z}\mathcal{V}$  and  $\mathbb{W}\mathcal{P}$  are often put in the rhs (i.e. added to  $h$ ) so that the matrix  $\mathbb{A}$  retains the same structure as in the incompressible case. This is indeed how it is implemented in ASPECT, see also appendix A of [369]. This however requires more work since the rhs depends on the solution and some form of iterations is needed.

**Remark 2:** Very often the adiabatic heating term  $\alpha T(\mathbf{v} \cdot \nabla p)$  is simplified as follows: If you assume the vertical component of the gradient of the dynamic pressure to be small compared to the gradient of the total pressure (in other words, the gradient is dominated by the gradient of the hydrostatic pressure), then  $-\rho\mathbf{g} \simeq \nabla p$  and then  $\alpha T(\mathbf{v} \cdot \nabla p) \simeq -\alpha\rho T\mathbf{v} \cdot \mathbf{g}$ . We will however not be using this approximation in what follows.

We have already established that

$$\vec{\tau} = \mathbf{C}_\eta \mathbf{B} V$$

The following measurements are carried out:

- The root mean square velocity (**vrms**):

$$v_{rms} = \sqrt{\frac{1}{V} \int_V v^2 dV}$$

- The average temperature (**Tavrg**):

$$\langle T \rangle = \frac{1}{V} \int_V T dV$$

- The total mass (**mass**):

$$M = \int_V \rho dV$$

- The Nusselt number (**Nu**):

$$Nu = -\frac{1}{Lx} \frac{1}{\Delta T} \int_0^{L_x} \frac{\partial T(x, y = L_y)}{\partial y} dx$$

- The kinetic energy (**EK**):

$$E_K = \int_V \frac{1}{2} \rho v^2 dV$$

- The work done against gravity

$$\langle W \rangle = - \int_V \rho g_y v_y dV$$

- The total viscous dissipation (**visc\_diss**)

$$\langle \Phi \rangle = \int \Phi dV = \frac{1}{V} \int 2\eta \dot{\epsilon} : \dot{\epsilon} dV$$

- The gravitational potential energy (**EG**)

$$E_G = \int_V \rho g_y (L_y - y) dV$$

- The internal thermal energy (ET)

$$E_T = \int_V \rho_{(0)} C_p T dV$$

**Remark 3:** Measuring the total mass can be misleading: indeed because  $\rho = \rho_0(1 - \alpha T)$ , then measuring the total mass amounts to measuring a constant minus the volume-integrated temperature, and there is no reason why the latter should be zero, so that there is no reason why the total mass should be zero...!

### 31.3 The experimental setup

The setup is as follows: the domain is  $Lx = Ly = 3000\text{km}$ . Free slip boundary conditions are imposed on all four sides. The initial temperature is given by:

$$T(x, y) = \left( \frac{L_y - y}{L_y} - 0.01 \cos\left(\frac{\pi x}{L_x}\right) \sin\left(\frac{\pi y}{L_y}\right) \right) \Delta T + T_{surf}$$

with  $\Delta T = 4000\text{K}$ ,  $T_{surf} = T_0 = 273.15\text{K}$ . The temperature is set to  $\Delta T + T_{surf}$  at the bottom and  $T_{surf}$  at the top. We also set  $k = 3$ ,  $C_p = 1250$ ,  $|g| = 10$ ,  $\rho_0 = 3000$  and we keep the Rayleigh number  $Ra$  and dissipation number  $Di$  as input parameters:

$$Ra = \frac{\alpha g \Delta T L^3 \rho_0^2 C_p}{\eta k} \quad Di = \frac{\alpha g L}{C_p}$$

From the second equation we get  $\alpha = \frac{Di C_p}{gL}$ , which we can insert in the first one:

$$Ra = \frac{Di C_p^2 \Delta T L^2 \rho_0^2}{\eta k} \quad \text{or}, \quad \eta = \frac{Di C_p^2 \Delta T L^2 \rho_0^2}{Ra k}$$

For instance, for  $Ra = 10^4$  and  $Di = 0.75$ , we obtain  $\alpha \simeq 3 \cdot 10^{-5}$  and  $\eta \simeq 10^{25}$  which are quite reasonable values.

### 31.4 Scaling

Following [343], we non-dimensionalize the equations using the reference values for density  $\rho_r$ , thermal expansivity  $\alpha_r$ , temperature contrast  $\Delta T_r$  (`refTemp`), thermal conductivity  $k_r$ , heat capacity  $C_p$ , depth of the fluid layer  $L$  and viscosity  $\eta_r$ . The non-dimensionalization for velocity,  $u_r$ , pressure  $p_r$  and time,  $t_r$  become

$$\begin{aligned} u_r &= \frac{k_r}{\rho_r C_p L} && (\text{refvel}) \\ p_r &= \frac{\eta_r k_r}{\rho_r C_p L^2} && (\text{refpress}) \\ t_r &= \frac{\rho_r C_p L^2}{k_r} && (\text{reftime}) \end{aligned}$$

In the case of the setup described hereabove, and when choosing  $Ra = 10^4$  and  $Di = 0.5$ , we get:

```
alphaT 2.08333e-05
eta 8.437500e+24
reftime 1.125000e+19
refvel 2.666667e-13
refPress 7.500000e+05
```

## 31.5 Conservation of energy 1

### 31.5.1 under BA and EBA approximations

Following [369], we take the dot product of the momentum equation with the velocity  $\mathbf{v}$  and integrate over the whole volume<sup>32</sup>:

$$\int_V [-\nabla \cdot \boldsymbol{\tau} + \nabla p] \cdot \mathbf{v} dV = \int_V \rho \mathbf{g} \cdot \mathbf{v} dV$$

or,

$$-\int_V (\nabla \cdot \boldsymbol{\tau}) \cdot \mathbf{v} dV + \int_V \nabla p \cdot \mathbf{v} dV = \int_V \rho \mathbf{g} \cdot \mathbf{v} dV$$

Let us look at each block separately:

$$-\int_V (\nabla \cdot \boldsymbol{\tau}) \cdot \mathbf{v} dV = -\int_S \underbrace{\boldsymbol{\tau} \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS + \int_V \boldsymbol{\tau} : \nabla \mathbf{v} dV = \int_V \boldsymbol{\tau} : \dot{\boldsymbol{\epsilon}} dV = \int_V \Phi dV$$

which is the volume integral of the shear heating. Then,

$$\int_V \nabla p \cdot \mathbf{v} dV = \int_S \underbrace{p \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS - \int_V \underbrace{\nabla \cdot \mathbf{v}}_{=0 \text{ (incomp.)}} pdV = 0$$

which is then zero in the case of an incompressible flow. And finally

$$\int_V \rho \mathbf{g} \cdot \mathbf{v} dV = W$$

which is the work against gravity.

Conclusion for an *incompressible* fluid: we should have

$$\int_V \Phi dV = \int_V \rho \mathbf{g} \cdot \mathbf{v} dV \quad (452)$$

This formula is hugely problematic: indeed, the term  $\rho$  in the rhs is the full density. We know that to the value of  $\rho_0$  corresponds a lithostatic pressure gradient  $p_L = \rho_0 gy$ . In this case one can write  $\rho = \rho_0 + \rho'$  and  $p = p_L + p'$  so that we also have

$$\int_V [-\nabla \cdot \boldsymbol{\tau} + \nabla p'] \cdot \mathbf{v} dV = \int_V \rho' \mathbf{g} \cdot \mathbf{v} dV$$

which will ultimately yield

$$\int_V \Phi dV = \int_V \rho' \mathbf{g} \cdot \mathbf{v} dV = \int_V (\rho - \rho_0) \mathbf{g} \cdot \mathbf{v} dV \quad (453)$$

Obviously Eqs.(452) and (453) cannot be true at the same time. The problem comes from the nature of the (E)BA approximation:  $\rho = \rho_0$  in the mass conservation equation but it is not constant in the momentum conservation equation, which is of course inconsistent. Since the mass conservation equation is  $\nabla \cdot \mathbf{v} = 0$  under this approximation then the term  $\int_V \nabla p \cdot \mathbf{v} dV$  is always zero for any pressure (full pressure  $p$ , or overpressure  $p - p_L$ ), hence the paradox. This paradox will be lifted when a consistent set of equations will be used (compressible formulation). On a practical note, Eqs.(452) is not verified by the code, while (453) is.

In the end:

$$\int_V \Phi dV = \underbrace{\int_V}_{\text{visc.diss}} (\rho - \rho_0) \underbrace{\mathbf{g} \cdot \mathbf{v} dV}_{\text{work_grav}}$$

(454)

<sup>32</sup>Check: this is akin to looking at the power, force\*velocity, says Arie

### 31.5.2 under no approximation at all

$$\int_V \nabla p \cdot \mathbf{v} dV = \int_S p \underbrace{\mathbf{v} \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS - \int_V \nabla \cdot \mathbf{v} pdV = 0 \quad (455)$$

$$= \int_V \frac{1}{\rho} \mathbf{v} \cdot \nabla \rho pdV = 0 \quad (456)$$

(457)

**ToDo:** see section 3 of [369] where this is carried out with the Adams-Williamson eos.

## 31.6 Conservation of energy 2

Also, following the Reynold's transport theorem [394], p210, we have for a property  $A$  (per unit mass)

$$\frac{d}{dt} \int_V A \rho dV = \int_V \frac{\partial}{\partial t} (A \rho) dV + \int_S A \rho \mathbf{v} \cdot \mathbf{n} dS$$

Let us apply to this to  $A = C_p T$  and compute the time derivative of the internal energy:

$$\frac{d}{dt} \int_V \rho C_p T dV = \int_V \frac{\partial}{\partial t} (\rho C_p T) dV + \int_S \rho C_p \underbrace{\mathbf{v} \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS = \underbrace{\int_V C_p T \frac{\partial \rho}{\partial t} dV}_I + \underbrace{\int_V \rho C_p \frac{\partial T}{\partial t} dV}_{II} \quad (458)$$

In order to expand  $I$ , the mass conservation equation will be used, while the heat transport equation will be used for  $II$ :

$$I = \int_V C_p T \frac{\partial \rho}{\partial t} dV = - \int_V C_p T \nabla \cdot (\rho \mathbf{v}) dV = - \int_V C_p T \rho \underbrace{\mathbf{v} \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS + \int_V \rho C_p \nabla T \cdot \mathbf{v} dV \quad (459)$$

$$II = \int_V \rho C_p \frac{\partial T}{\partial t} dV = \int_V \left[ -\rho C_p \mathbf{v} \cdot \nabla T + \nabla \cdot k \nabla T + \rho H + \Phi + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \right] dV \quad (460)$$

$$= \int_V \left[ -\rho C_p \mathbf{v} \cdot \nabla T + \rho H + \Phi + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \right] dV + \int_V \nabla \cdot k \nabla T dV \quad (461)$$

$$= \int_V \left[ -\rho C_p \mathbf{v} \cdot \nabla T + \rho H + \Phi + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \right] dV + \int_S k \nabla T \cdot \mathbf{n} dS \quad (462)$$

$$= \int_V \left[ -\rho C_p \mathbf{v} \cdot \nabla T + \rho H + \Phi + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \right] dV - \int_S \mathbf{q} \cdot \mathbf{n} dS \quad (463)$$

Finally:

$$I + II = \frac{d}{dt} \underbrace{\int_V \rho C_p T dV}_{\text{ET}} = \int_V \left[ \rho H + \Phi + \alpha T \left( \frac{\partial p}{\partial t} + \mathbf{v} \cdot \nabla p \right) \right] dV - \int_S \mathbf{q} \cdot \mathbf{n} dS \quad (464)$$

$$= \int_V \rho H dV + \underbrace{\int_V \Phi dV}_{\text{visc.diss}} + \underbrace{\int_V \alpha T \frac{\partial p}{\partial t} dV}_{\text{extra}} + \underbrace{\int_V \alpha T \mathbf{v} \cdot \nabla p dV}_{\text{adiab_heating}} - \underbrace{\int_S \mathbf{q} \cdot \mathbf{n} dS}_{\text{heatflux_boundary}} \quad (465)$$

This was of course needlessly complicated as the term  $\partial \rho / \partial t$  is always taken to be zero, so that  $I = 0$  automatically. The mass conservation equation is then simply  $\nabla \cdot (\rho \mathbf{v}) = 0$ . Then it follows that

$$0 = \int_V C_p T \nabla \cdot (\rho \mathbf{v}) dV = - \int_V C_p T \rho \underbrace{\mathbf{v} \cdot \mathbf{n}}_{=0 \text{ (b.c.)}} dS + \int_V \rho C_p \nabla T \cdot \mathbf{v} dV \quad (466)$$

$$= \int_V \rho C_p \nabla T \cdot \mathbf{v} dV \quad (467)$$

so that the same term in Eq.(463) vanishes too, and then Eq.(465) is always valid, although one should be careful when computing  $E_T$  in the BA and EBA cases as it should use  $\rho_0$  and not  $\rho$ .

### 31.7 The problem of the onset of convection

[wiki] In geophysics, the Rayleigh number is of fundamental importance: it indicates the presence and strength of convection within a fluid body such as the Earth's mantle. The mantle is a solid that behaves as a fluid over geological time scales.

The Rayleigh number essentially is an indicator of the type of heat transport mechanism. At low Rayleigh numbers conduction processes dominate over convection ones. At high Rayleigh numbers it is the other way around. There is a so-called critical value of the number with delineates the transition from one regime to the other.

This problem has been studied and approached both theoretically and numerically [554, e.g.] and it was found that the critical Rayleigh number  $Ra_c$  is

$$Ra_c = (27/4)\pi^4 \simeq 657.5$$

in setups similar to ours.

#### VERY BIG PROBLEM

The temperature setup is built as follows:  $T_{surf}$  is prescribed at the top,  $T_{surf} + \Delta T$  is prescribed at the bottom. The initial temperature profile is linear between these two values. In the case of BA, the actual value of  $T_{surf}$  is of no consequence. However, for the EBA the full temperature is present in the adiabatic heating term on the rhs of the hte, and the value of  $T_{surf}$  will therefore influence the solution greatly. This is very problematic as there is no real way to arrive at the surface temperature from the King paper. On top of this, the density uses a reference temperature  $T_0$  which too will influence the solution without being present in the controlling  $Ra$  and  $Di$  numbers!!

In light thereof, it will be very difficult to recover the values of King et al for EBA!

#### features

- $Q_1 \times P_0$  element
- compressible flow
- mixed formulation
- Dirichlet boundary conditions (no-slip)
- isoviscous
- analytical solution
- pressure smoothing

Relevant literature: [52, 321, 532, 369, 343, 370, 384, 299]

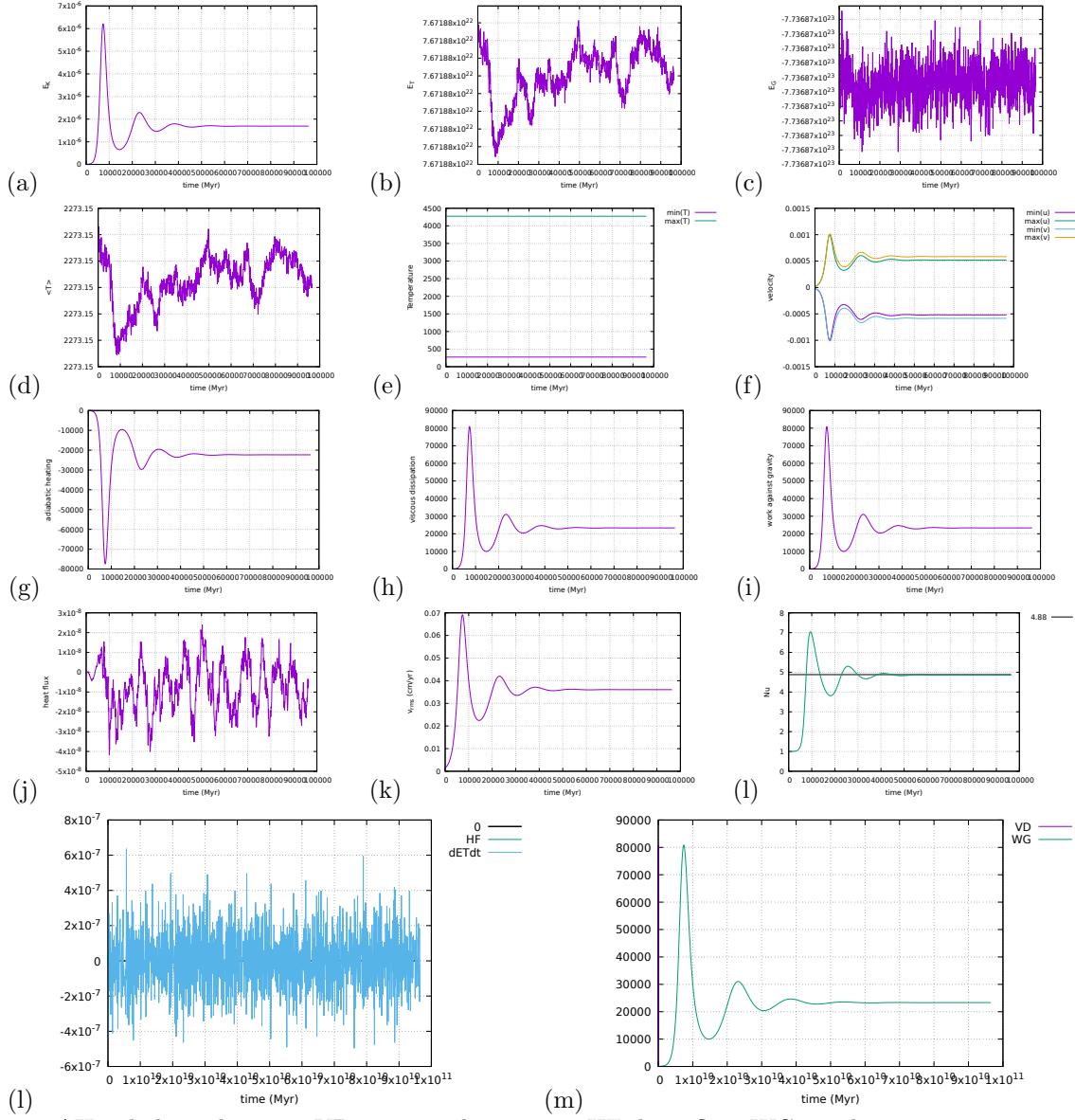
ToDo:

- heat flux is at the moment elemental, so Nusselt and heat flux on boundaries measurements not as accurate as could be.

- implement steady state detection
- do  $Ra = 10^5$  and  $Ra = 10^6$
- velocity average at surface
- non dimensional heat flux at corners [64]
- depth-dependent viscosity (case 2 of [64])

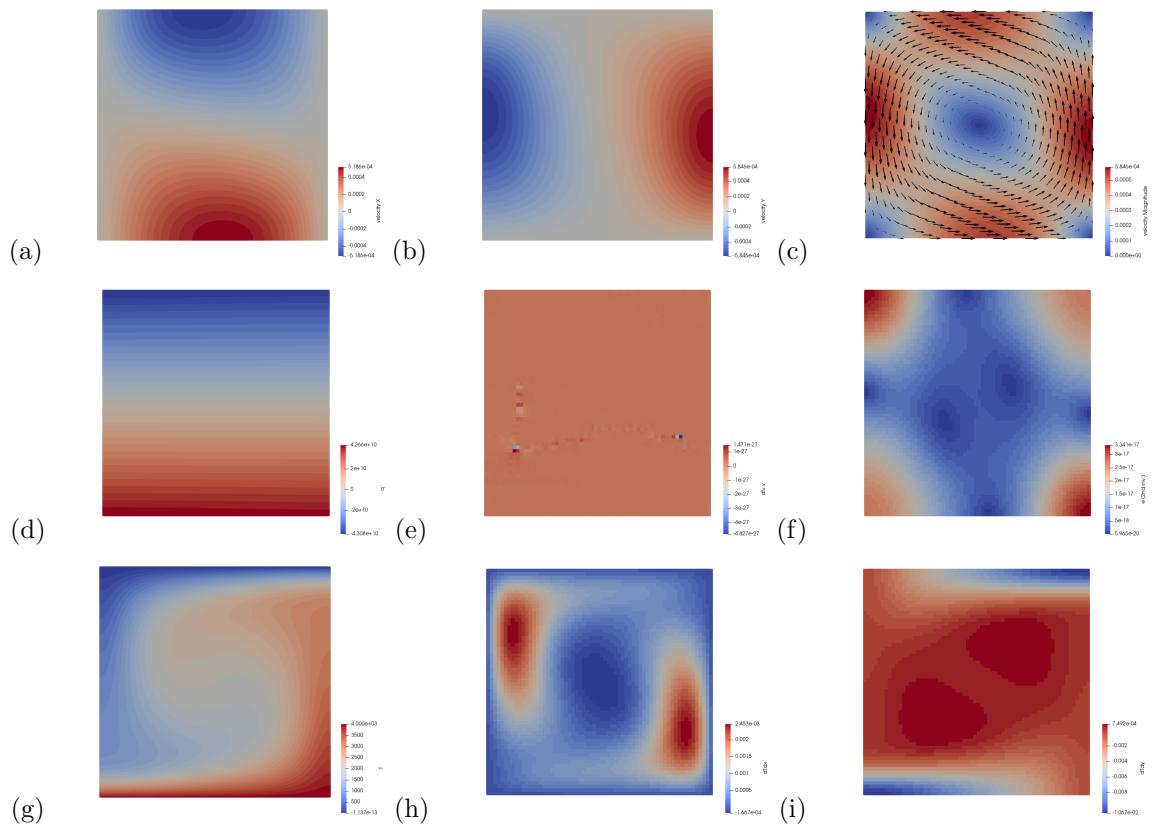
### 31.8 results - BA - $Ra = 10^4$

These results were obtained with a 64x64 resolution, and CFL number of 1. Steady state was reached after about 1250 timesteps.



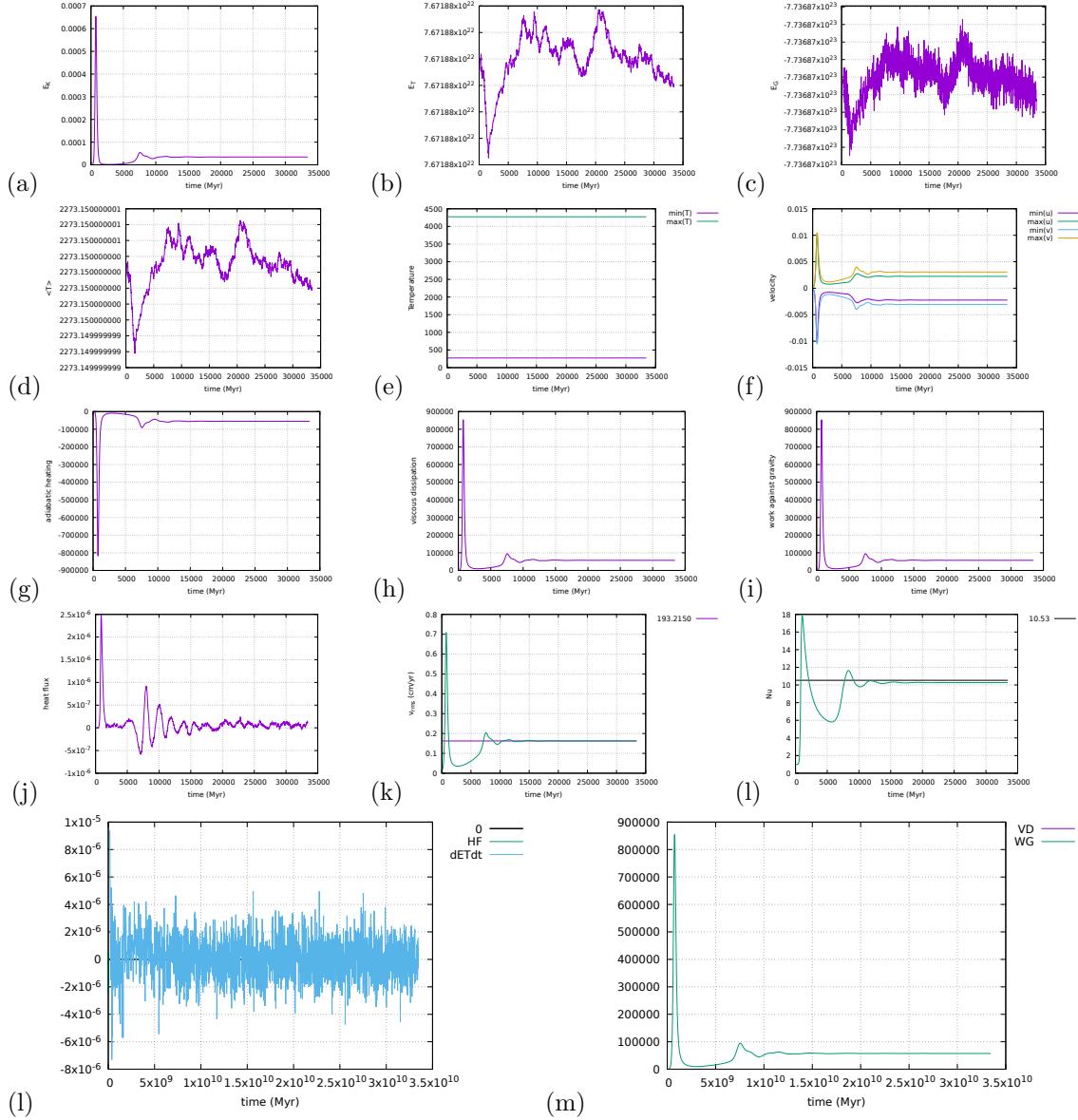
AH: adiabatic heating, VD: viscous dissipation, HF: heat flux, WG: work against gravity

Eq.(465) is verified by (l) and Eq.(454) is verified by (m).



### 31.9 results - BA - $Ra = 10^5$

These results were obtained with a 64x64 resolution, and CFL number of 1. Steady state was reached after about 1250 timesteps.

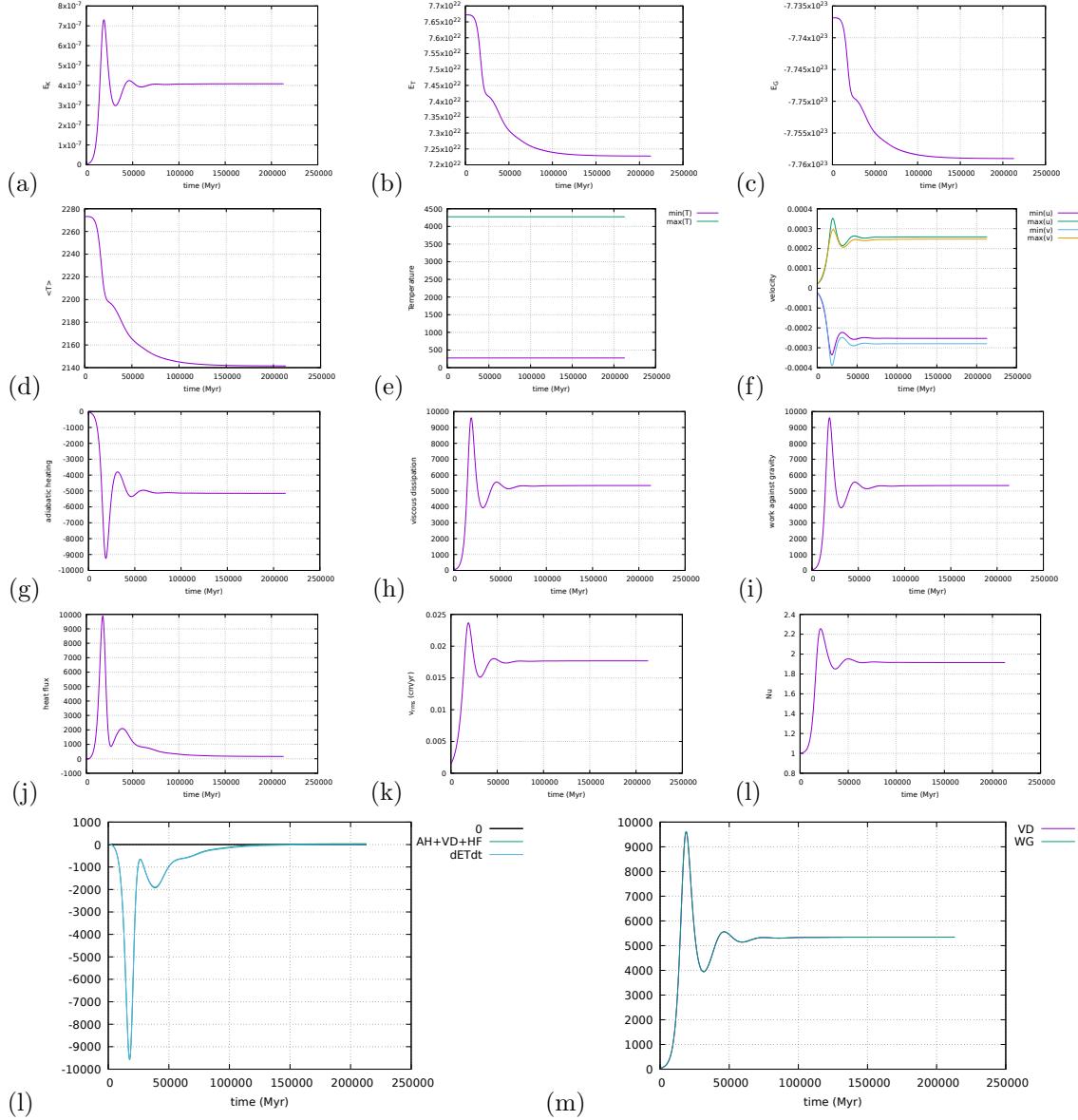


Eq.(465) is verified by (l) and Eq.(454) is verified by (m).

### **31.10 results - BA - $Ra = 10^6$**

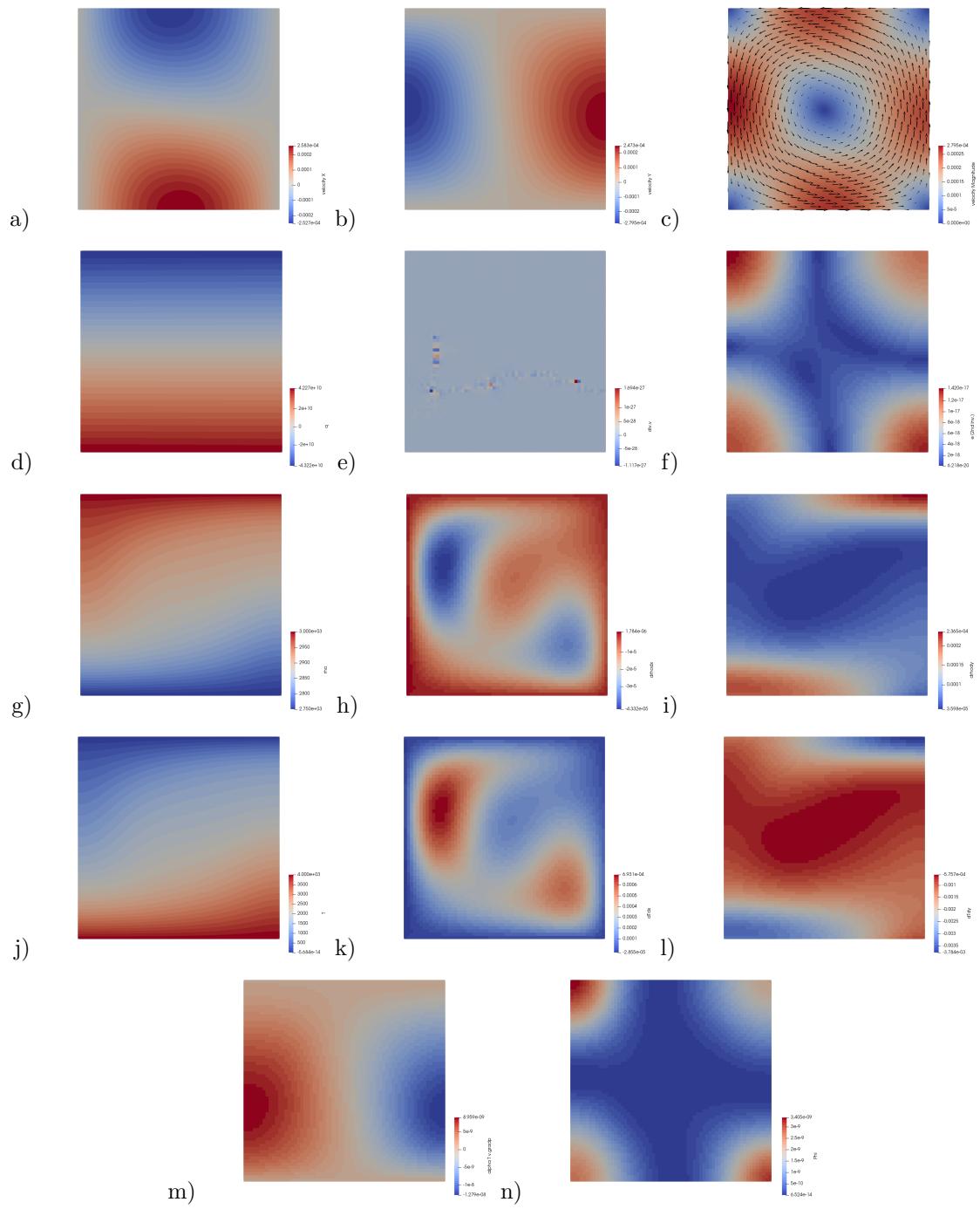
### 31.11 results - EBA - $Ra = 10^4$

These results were obtained with a 64x64 resolution, and CFL number of 1. Steady state was reached after about 2500 timesteps



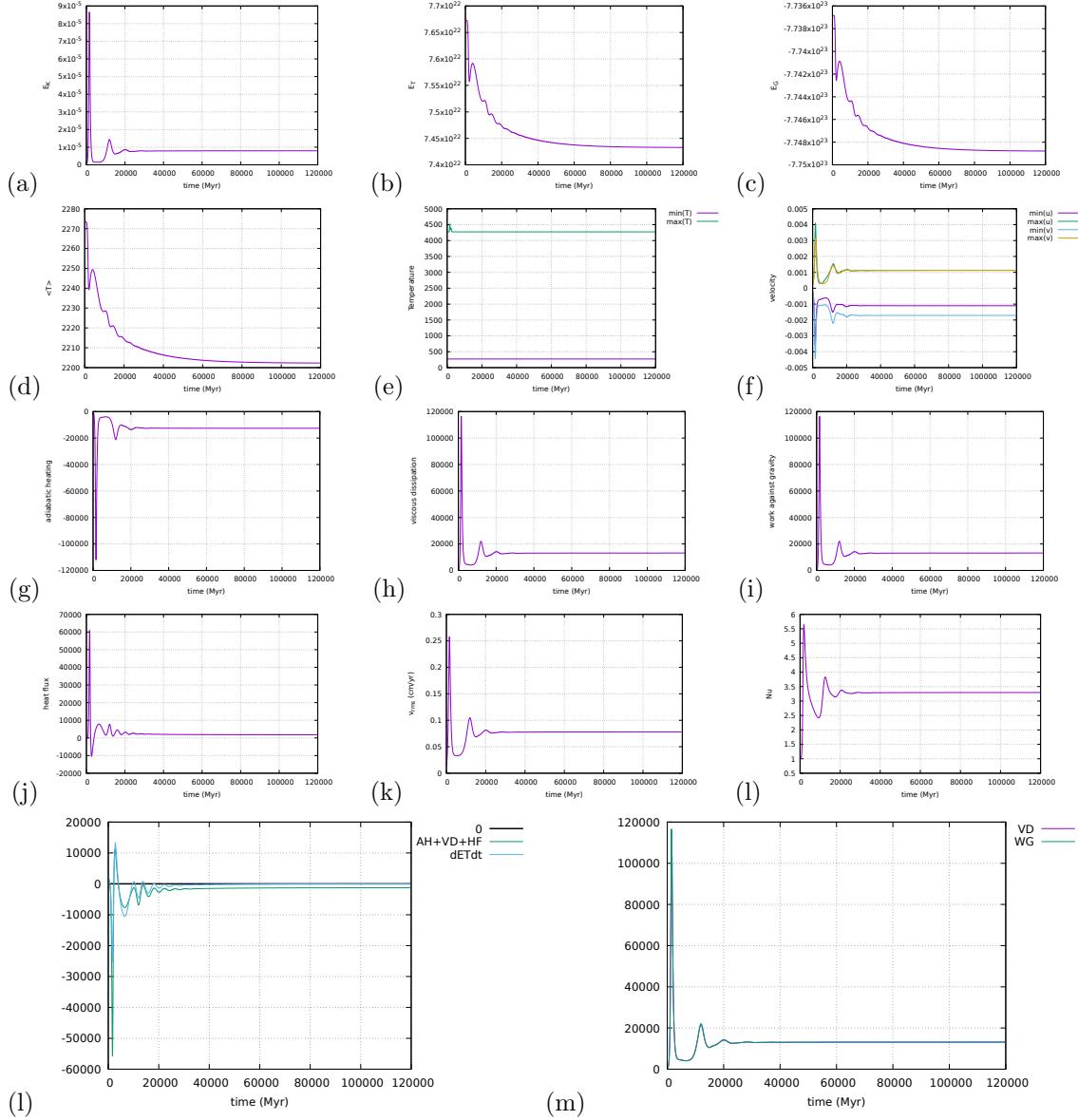
AH: adiabatic heating, VD: viscous dissipation, HF: heat flux, WG: work against gravity

Eq.(465) is verified by (l) and Eq.(454) is verified by (m).



### 31.12 results - EBA - $Ra = 10^5$

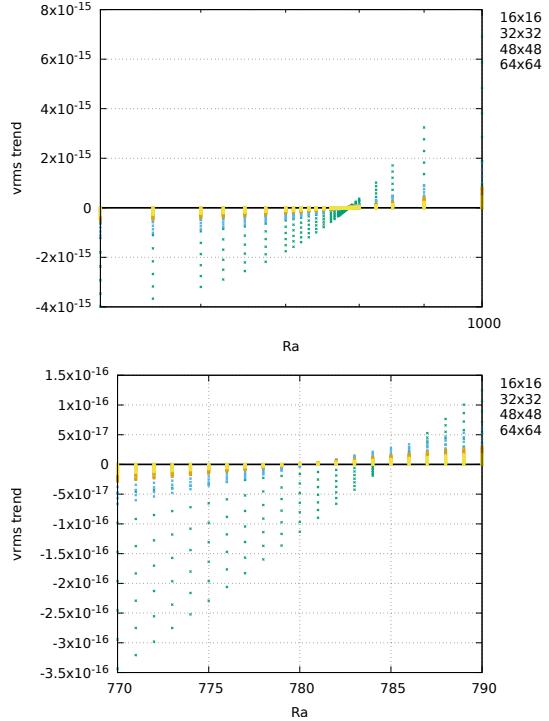
These results were obtained with a 64x64 resolution, and CFL number of 1. Simulation was stopped after about 4300 timesteps.



AH: adiabatic heating, VD: viscous dissipation, HF: heat flux, WG: work against gravity

### 31.13 Onset of convection

The code can be run for values of Ra between 500 and 1000, at various resolutions for the BA formulation. The value  $v_{rms}(t) - v_{rms}(0)$  is plotted as a function of  $Ra$  and for the 10 first timesteps. If the  $v_{rms}$  is found to decrease, then the Rayleigh number is not high enough to allow for convection and the initial temperature perturbation relaxes by diffusion (and then  $v_{rms}(t) - v_{rms}(0) < 0$ ). If the  $v_{rms}$  is found to increase, then  $v_{rms}(t) - v_{rms}(0) > 0$  and the system is going to showcase convection. The zero value of  $v_{rms}(t) - v_{rms}(0)$  gives us the critical Rayleigh number, which is found between 775 and 790.

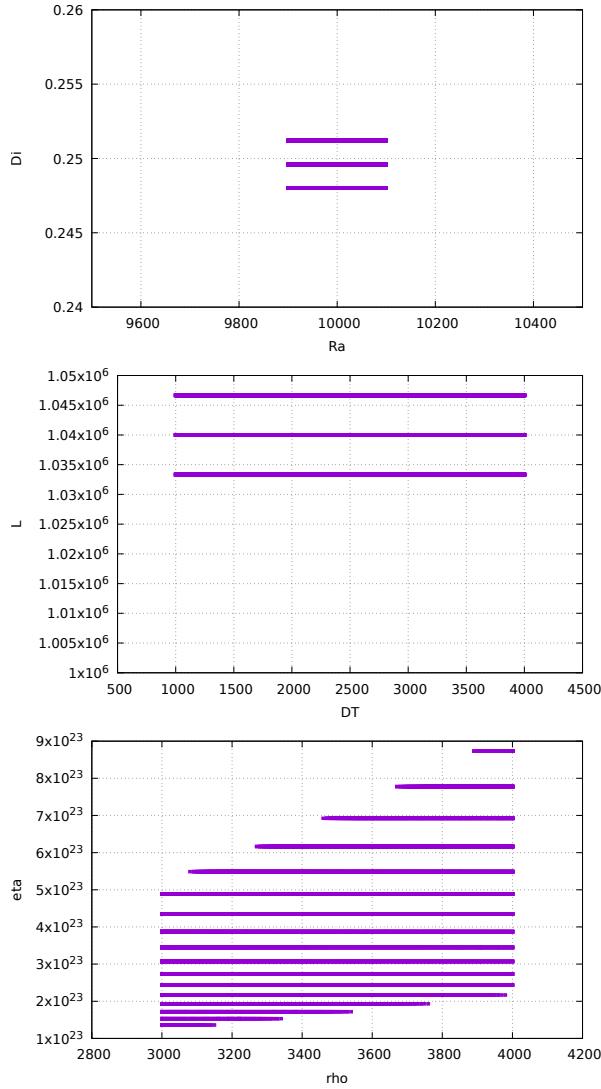


**Appendix:** Looking for the right combination of parameters for the King benchmark.

I run a quadruple do loop over  $L$ ,  $\Delta T$ ,  $\rho_0$  and  $\eta_0$  between plausible values (see code targets.py) and write in a file only the combination which yields the required Rayleigh and Dissipation number values (down to 1% accuracy).

```
alpha=3e-5
g=10
hcapa=1250
hcond=3
DTmin=1000 ; DTmax=4000 ; DTnpts=251
Lmin=1e6 ; Lmax=3e6 ; Lnpts=251
rhomin=3000 ; rhomax=3500 ; rhonpts=41
etamin=19 ; etamax=25 ; etanpts=100
```

On the following plots the 'winning' combinations of these four parameters are shown:



We see that:

- the parameter  $L$  (being to the 3rd power in the  $Ra$  number) cannot vary too much. Although it is varied between 1000 and 3000km there seems to be a 'right' value at about 1040 km. (why?)
- viscosities are within  $10^{23}$  and  $10^{24}$  which are plausible values (although a bit high?).
- densities can be chosen freely between 3000 and 3500
- $\Delta T$  seems to be the most problematic value since it can range from 1000 to 4000K ...

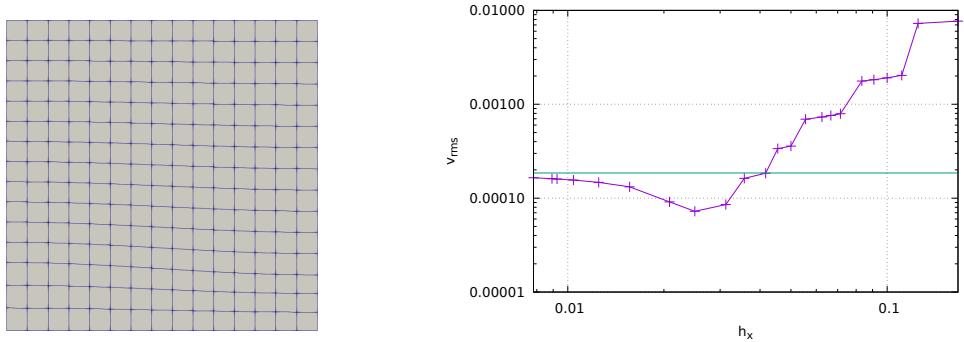
## 32 fieldstone\_25: Rayleigh-Taylor instability (1) - instantaneous

This numerical experiment was first presented in [572]. 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  $y(x) = 0.2 + 0.02 \cos\left(\frac{\pi x}{L_x}\right)$ . The bottom fluid is parametrised by its mass density  $\rho_1$  and its viscosity  $\mu_1$ , while the layer above is parametrised by  $\rho_2$  and  $\mu_2$ .

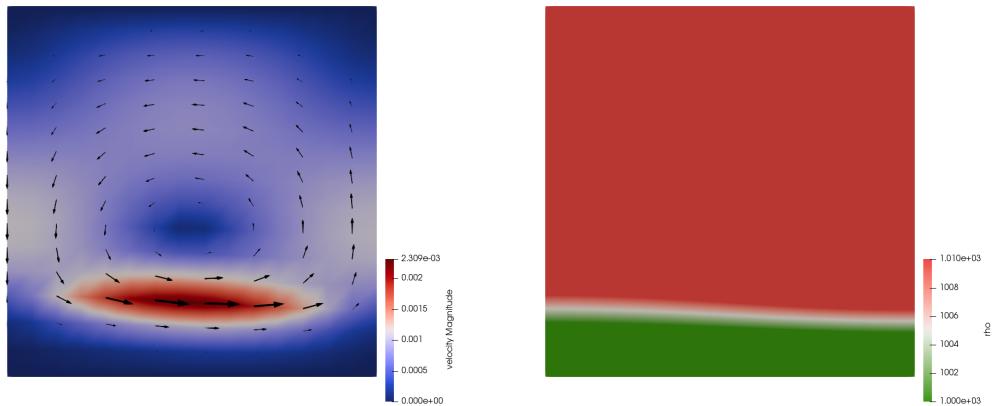
No-slip boundary conditions are applied at the bottom and at the top of the box while free-slip boundary conditions are applied on the sides.

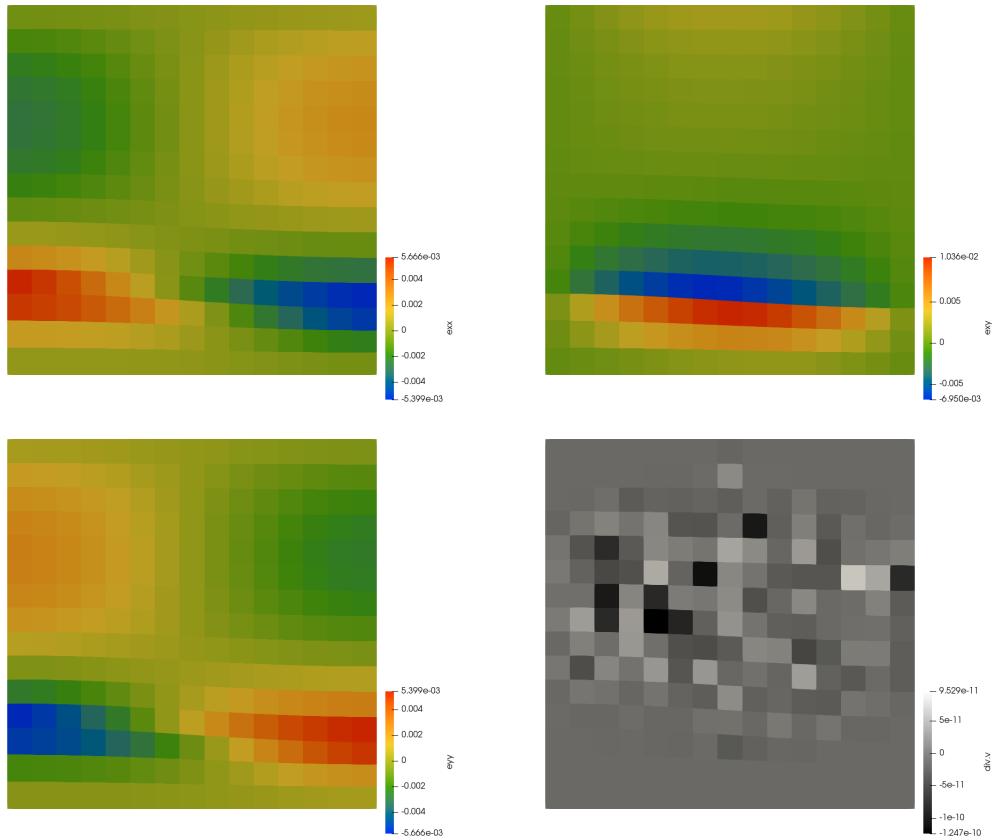
In the original benchmark the system is run over 2000 units of dimensionless time and the timing and position of various upwellings/downwellings is monitored. In this present experiment only the root mean square velocity is measured at  $t = 0$ : the code is indeed not yet foreseen of any algorithm capable of tracking deformation.

Another approach than the ones presented in the extensive literature which showcases results of this benchmark is taken. The mesh is initially fitted to the fluids interface and the resolution is progressively increased. This results in the following figure:



The green line indicates results obtained with my code ELEFANT with grids up to 2000x2000 with the exact same methodology.





### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- isothermal
- numerical benchmark

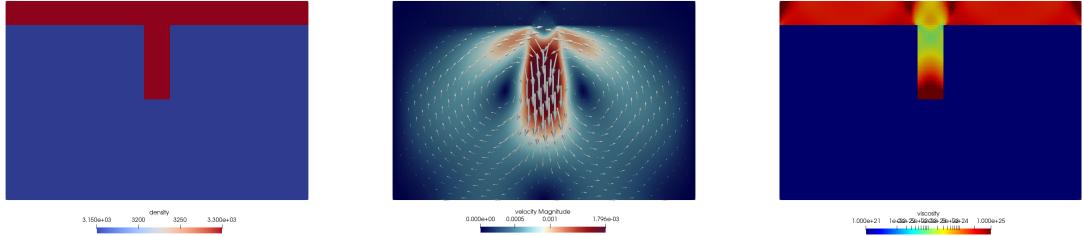
### 33 fieldstone\_26: Slab detachment benchmark (1) - instantaneous

As in [496], the computational domain is  $1000\text{km} \times 660\text{km}$ . No-slip boundary conditions are imposed on the sides of the system while free-slip boundary conditions are imposed at the top and bottom. Two materials are present in the domain: the lithosphere (mat.1) and the mantle (mat.2). The overriding plate (mat.1) is  $80\text{km}$  thick and is placed at the top of the domain. An already subducted slab (mat.1) of  $250\text{km}$  length hangs vertically under this plate. The mantle occupies the rest of the domain.

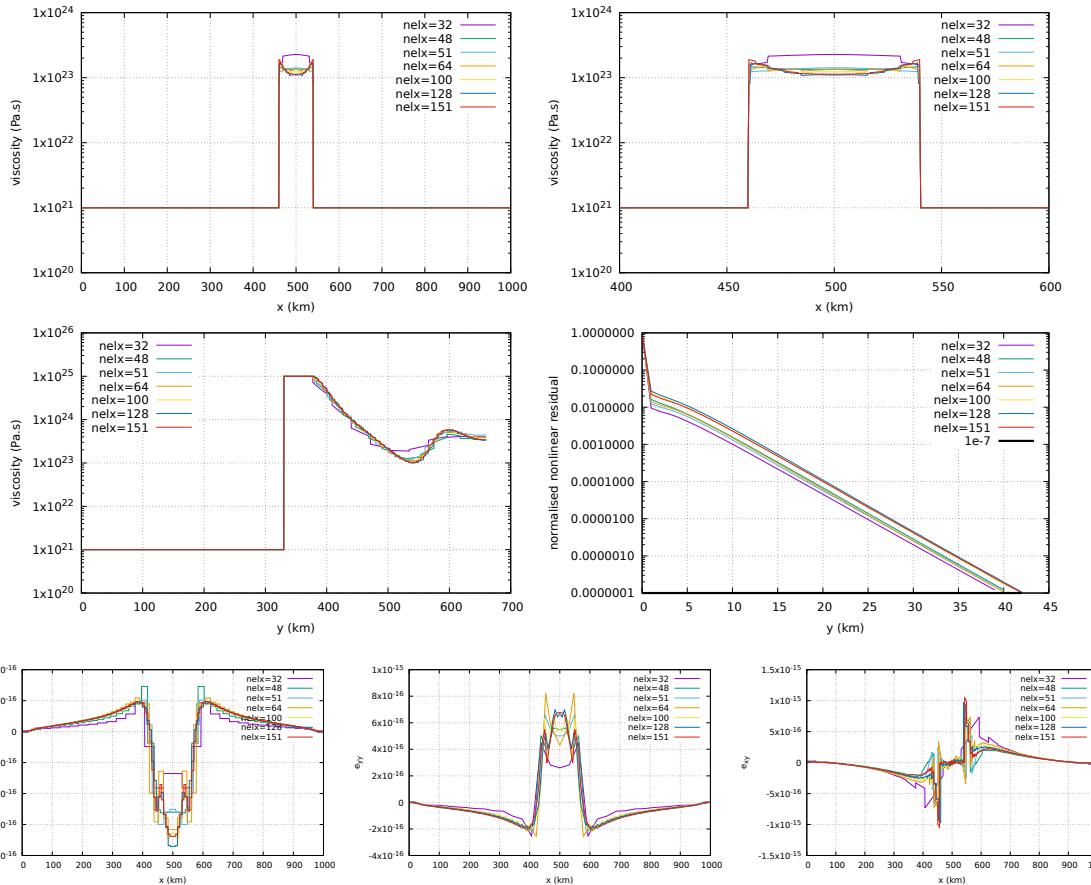
The mantle has a constant viscosity  $\eta_0 = 10^{21}\text{Pa.s}$  and a density  $\rho = 3150\text{kg/m}^3$ . The slab has a density  $\rho = 3300\text{kg/m}^3$  and is characterised by a power-law flow law so that its effective viscosity depends on the second invariant of the strainrate  $I_2$  as follows:

$$\eta_{eff} = \frac{1}{2}A^{-1/n_s}I_2^{1/n_s-1} = \frac{1}{2}[(2 \times 4.75 \times 10^{11})^{-n_s}]^{-1/n_s}I_2^{1/n_s-1} = 4.75 \times 10^{11}I_2^{1/n_s-1} = \eta_0 I_2^{1/n_s-1} \quad (468)$$

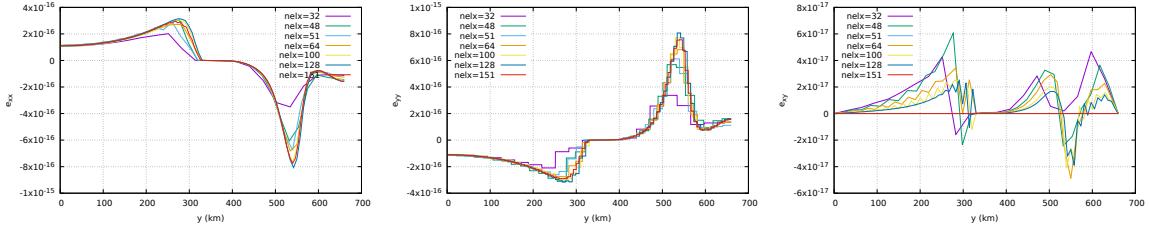
with  $n_s = 4$  and  $A = (2 \times 4.75 \times 10^{11})^{-n_s}$ , or  $\eta_0 = 4.75 \times 10^{11}$ .



Fields at convergence for 151x99 grid.



Along the horizontal line



Along the vertical line

### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- isothermal
- nonlinear rheology
- nonlinear residual

Todo: nonlinear mantle, pressure normalisation

## 34 fieldstone\_27: Consistent Boundary Flux

In what follows we will be re-doing the numerical experiments presented in Zhong et al. [623].

The first benchmark showcases a unit square domain with free slip boundary conditions prescribed on all sides. The resolution is fixed to  $64 \times 64 Q_1 \times P_0$  elements. The flow is isoviscous and the buoyancy force  $\mathbf{f}$  is given by

$$\begin{aligned} f_x &= 0 \\ f_y &= \rho_0 \alpha T(x, y) \end{aligned}$$

with the temperature field given by

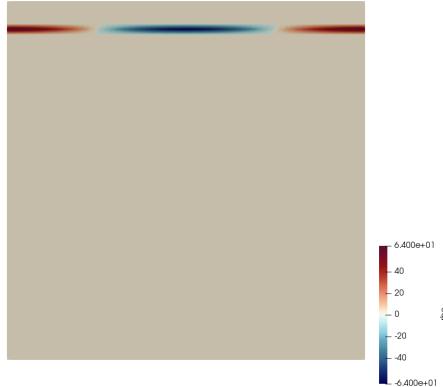
$$T(x, y) = \cos(kx)\delta(y - y_0)$$

where  $k = 2\pi/\lambda$  and  $\lambda$  is a wavelength, and  $y_0$  represents the location of the buoyancy strip. We set  $g_y = -1$  and prescribe  $\rho(x, y) = \rho_0 \alpha \cos(kx)\delta(y - y_0)$  on the nodes of the mesh.

One can prove ([623] and refs. therein) that there is an analytic solution for the surface stress  $\sigma_{zz}$ <sup>33</sup>

$$\frac{\sigma_{yy}}{\rho \alpha g h} = \frac{\cos(kx)}{\sinh^2(k)} [k(1 - y_0) \sinh(k) \cosh(ky_0) - k \sinh(k(1 - y_0)) + \sinh(k) \sinh(ky_0)]$$

We choose  $\rho_0 \alpha = 64$ ,  $\eta = 1$  (note that in this case the normalising coefficient of the stress is exactly 1 (since  $h = L_x/nelx = 1/64$ ) so it is not implemented in the code).  $\lambda = 1$  is set to 1 and we explore  $y_0 = \frac{63}{64}, \frac{62}{64}, \frac{59}{64}$  and  $y_0 = 32/64$ . Under these assumptions the density field for  $y_0 = 59/64$  is:

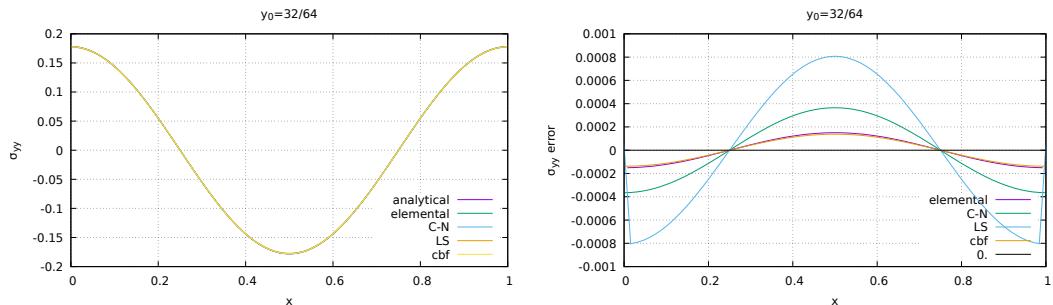


We can recover the stress at the boundary by computing the  $yy$  component of the stress tensor in the top row of elements:

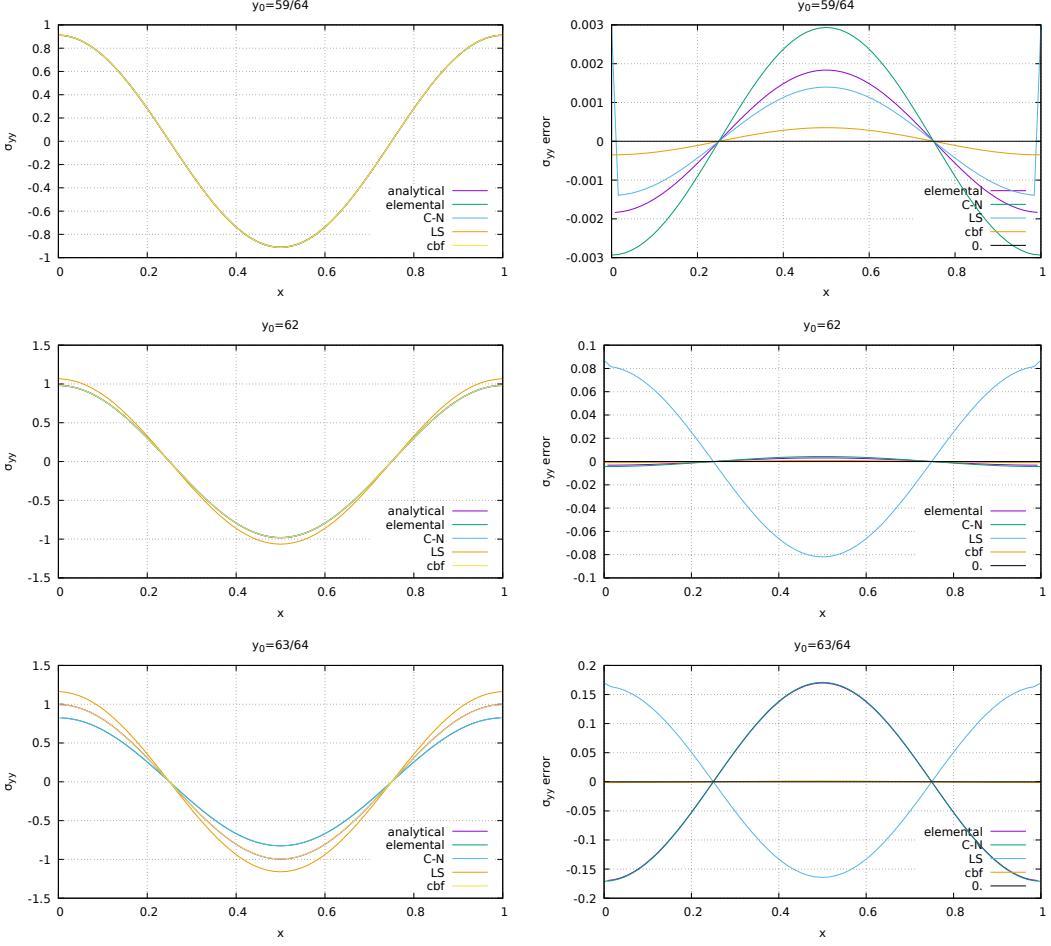
$$\sigma_{yy} = -p + 2\eta\dot{\epsilon}_{yy}$$

Note that pressure is by definition elemental, and that strain rate components are then also computed in the middle of each element.

These elemental quantities can be projected onto the nodes (see section ??) by means of the C→N algorithm or a least square algorithm (LS).



<sup>33</sup>Note that in the paper the authors use  $\rho \alpha g$  which does not have the dimensions of a stress



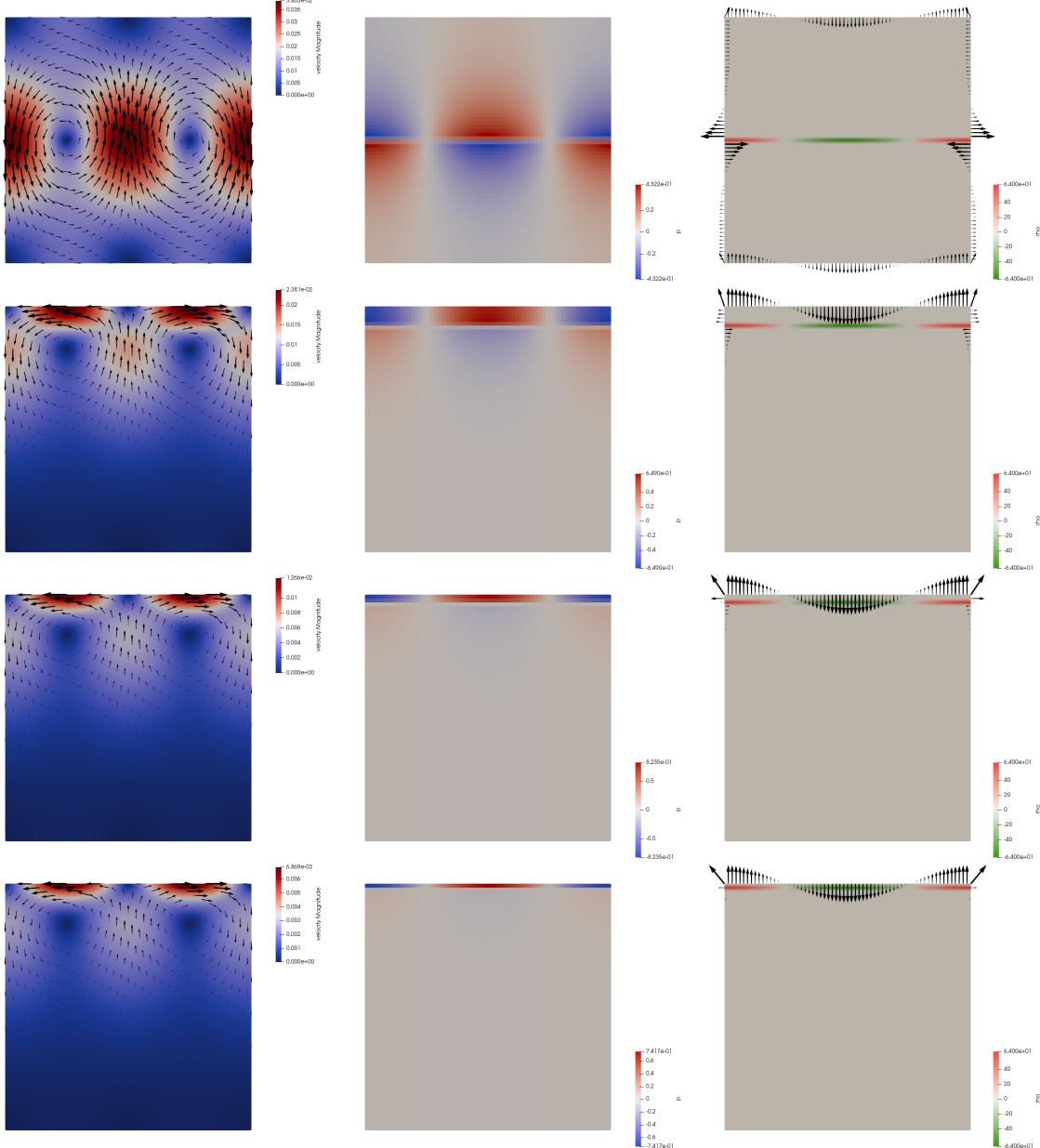
The consistent boundary flux (CBF) method allows us to compute traction vectors  $\mathbf{t} = \boldsymbol{\sigma} \cdot \mathbf{n}$  on the boundary of the domain. On the top boundary,  $\mathbf{n} = (0, 1)$  so that  $\mathbf{t} = (\sigma_{xy}, \sigma_{yy})^T$  and  $t_y$  is the quantity we need to consider and compare to other results.

In the following table are shown the results presented in [623] alongside the results obtained with Fieldstone:

Method	$y_0 = 63/64$	$y_0 = 62/64$	$y_0 = 59/64$ <sup>34</sup>	$y_0 = 32/64$
Analytic solution	0.995476	0.983053	0.912506	0.178136
Pressure smoothing [623]	1.15974	1.06498	0.911109	n.a.
CBF [623]	0.994236	0.982116	0.912157	n.a.
fieldstone: elemental	0.824554 (-17.17 %)	0.978744 (-0.44%)	0.909574 (-0.32 %)	0.177771 (-0.20 %)
fieldstone: nodal (C→N)	0.824554 (-17.17 %)	0.978744 (-0.44%)	0.909574 (-0.32 %)	0.177771 (-0.20 %)
fieldstone: LS	1.165321 ( 17.06 %)	1.070105 ( 8.86%)	0.915496 ( 0.33 %)	0.178182 ( 0.03 %)
fieldstone: CBF	0.994236 ( -0.13 %)	0.982116 (-0.10%)	0.912157 (-0.04 %)	0.177998 (-0.08 %)

We see that we recover the published results with the same exact accuracy, thereby validating our implementation.

On the following figures are shown the velocity, pressure and traction fields for two cases  $y_0 = 32/64$  and  $y_0 = 63/64$ .



Here lies the superiority of our approach over the one presented in the original article: our code computes all traction vectors on all boundaries at once.

[explain how Medge is arrived at!](#)

[compare with ASPECT ???](#)

[gauss-lobatto integration?](#)

[pressure average on surface instead of volume ?](#)

**features**

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- isothermal
- isoviscous
- analytical solution
- pressure smoothing
- consistent boundary flux

## 35 fieldstone\_28: convection 2D box - Tosi et al, 2015

This fieldstone was developed in collaboration with Rens Elbertsen.

The viscosity field  $\mu$  is calculated as the harmonic average between a linear part  $\mu_{lin}$  that depends on temperature only or on temperature and depth  $d$ , and a non-linear, plastic part  $\mu_{plast}$  dependent on the strain rate:

$$\mu(T, z, \dot{\epsilon}) = 2 \left( \frac{1}{\mu_{lin}(T, z)} + \frac{1}{\mu_{plast}(\dot{\epsilon})} \right)^{-1}. \quad (469)$$

The linear part is given by the linearized Arrhenius law (the so-called Frank-Kamenetskii approximation [?]):

$$\mu_{lin}(T, z) = \exp(-\gamma_T T + \gamma_z z), \quad (470)$$

where  $\gamma_T = \ln(\Delta\mu_T)$  and  $\gamma_z = \ln(\Delta\mu_z)$  are parameters controlling the total viscosity contrast due to temperature ( $\Delta\mu_T$ ) and pressure ( $\Delta\mu_z$ ). The non-linear part is given by [?]:

$$\mu_{plast}(\dot{\epsilon}) = \mu^* + \frac{\sigma_Y}{\sqrt{\dot{\epsilon} : \dot{\epsilon}}}, \quad (471)$$

where  $\mu^*$  is a constant representing the effective viscosity at high stresses [522] and  $\sigma_Y$  is the yield stress, also assumed to be constant. In 2-D, the denominator in the second term of equation (471) is given explicitly by

$$\sqrt{\dot{\epsilon} : \dot{\epsilon}} = \sqrt{\dot{\epsilon}_{ij}\dot{\epsilon}_{ij}} = \sqrt{\left(\frac{\partial u_x}{\partial x}\right)^2 + \frac{1}{2} \left(\frac{\partial u_x}{\partial y} + \frac{\partial u_y}{\partial x}\right)^2 + \left(\frac{\partial u_y}{\partial y}\right)^2}. \quad (472)$$

The viscoplastic flow law (equation 469) leads to linear viscous deformation at low stresses (equation (470)) and to plastic deformation for stresses that exceed  $\sigma_Y$  (equation (471)), with the decrease in viscosity limited by the choice of  $\mu^*$  [522].

In all cases that we present, the domain is a two-dimensional square box. The mechanical boundary conditions are for all boundaries free-slip with no flux across, i.e.  $\tau_{xy} = \tau_{yx} = 0$  and  $\mathbf{u} \cdot \mathbf{n} = 0$ , where  $\mathbf{n}$  denotes the outward normal to the boundary. Concerning the energy equation, the bottom and top boundaries are isothermal, with the temperature  $T$  set to 1 and 0, respectively, while side-walls are assumed to be insulating, i.e.  $\partial T / \partial x = 0$ . The initial distribution of the temperature field is prescribed as follows:

$$T(x, y) = (1 - y) + A \cos(\pi x) \sin(\pi y), \quad (473)$$

where  $A = 0.01$  is the amplitude of the initial perturbation.

In the following Table ??, we list the benchmark cases according to the parameters used.

Case	$Ra$	$\Delta\mu_T$	$\Delta\mu_y$	$\mu^*$	$\sigma_Y$	Convective regime
1	$10^2$	$10^5$	1	—	—	Stagnant lid
2	$10^2$	$10^5$	1	$10^{-3}$	1	Mobile lid
3	$10^2$	$10^5$	10	—	—	Stagnant lid
4	$10^2$	$10^5$	10	$10^{-3}$	1	Mobile lid
5a	$10^2$	$10^5$	10	$10^{-3}$	4	Periodic
5b	$10^2$	$10^5$	10	$10^{-3}$	3 – 5	Mobile lid – Periodic – Stagnant lid

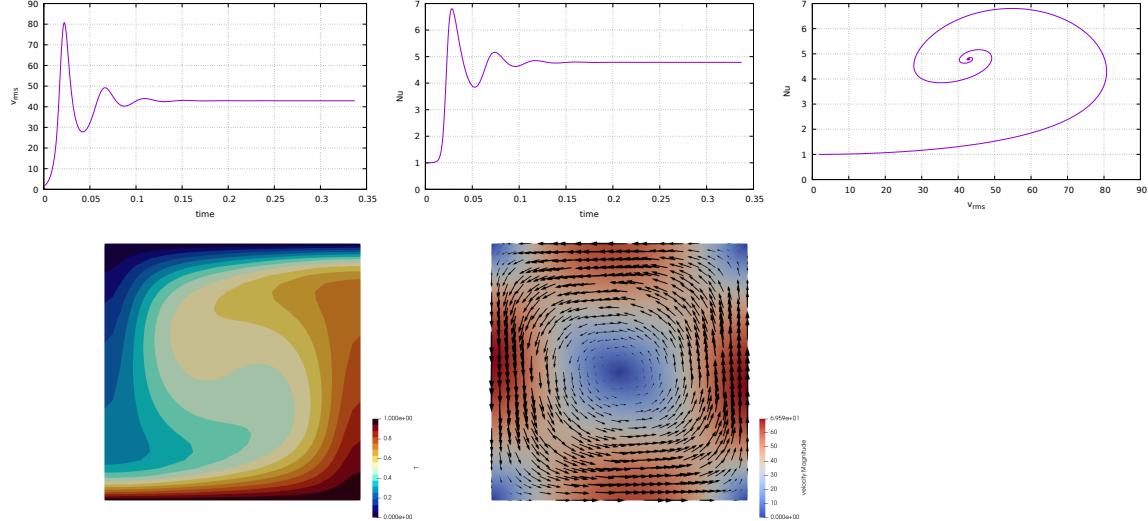
Benchmark cases and corresponding parameters.

In Cases 1 and 3 the viscosity is directly calculated from equation (470), while for Cases 2, 4, 5a, and 5b, we used equation (469). For a given mesh resolution, Case 5b requires running simulations with yield stress varying between 3 and 5

In all tests, the reference Rayleigh number is set at the surface ( $y = 1$ ) to  $10^2$ , and the viscosity contrast due to temperature  $\Delta\mu_T$  is  $10^5$ , implying therefore a maximum effective Rayleigh number of  $10^7$  for  $T = 1$ . Cases 3, 4, 5a, and 5b employ in addition a depth-dependent rheology with viscosity contrast  $\Delta\mu_z = 10$ . Cases 1 and 3 assume a linear viscous rheology that leads to a stagnant lid regime.

Cases 2 and 4 assume a viscoplastic rheology that leads instead to a mobile lid regime. Case 5a also assumes a viscoplastic rheology but a higher yield stress, which ultimately causes the emergence of a strictly periodic regime. The setup of Case 5b is identical to that of Case 5a but the test consists in running several simulations using different yield stresses. Specifically, we varied  $\sigma_Y$  between 3 and 5 in increments of 0.1 in order to identify the values of the yield stress corresponding to the transition from mobile to periodic and from periodic to stagnant lid regime.

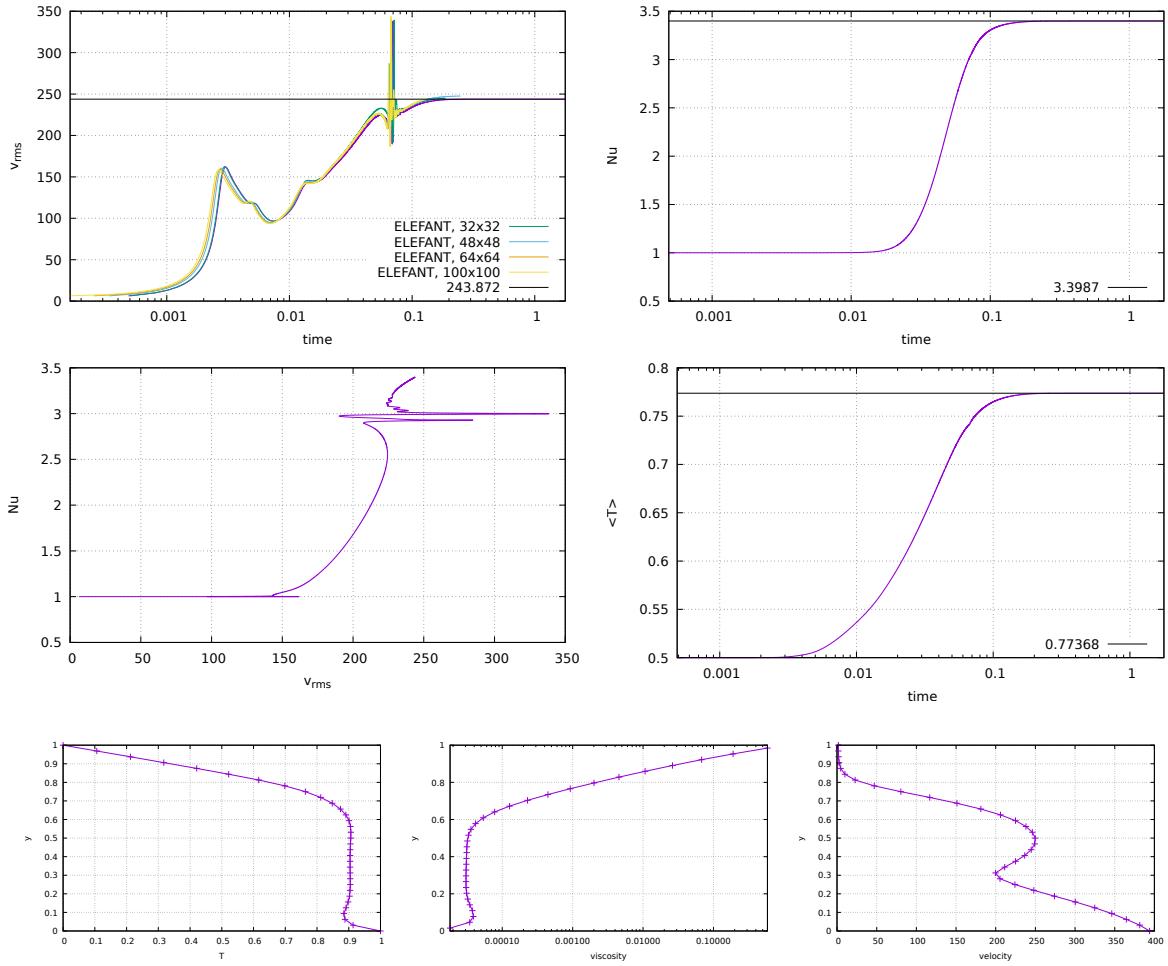
### 35.0.1 Case 0: Newtonian case, a la Blankenbach et al., 1989

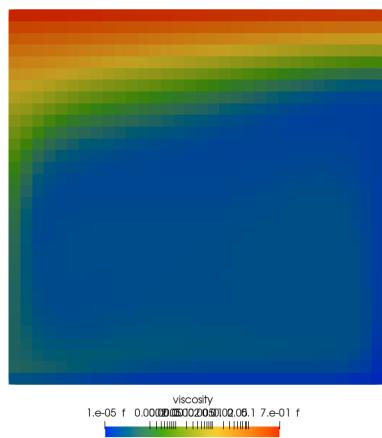
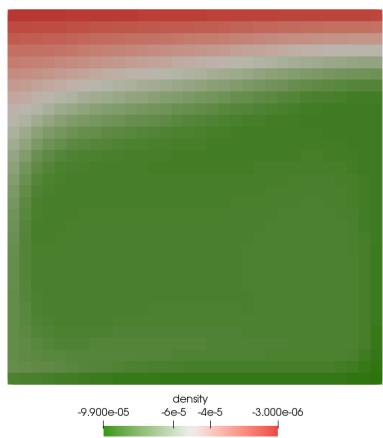
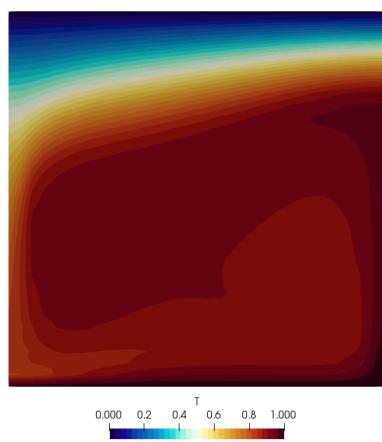
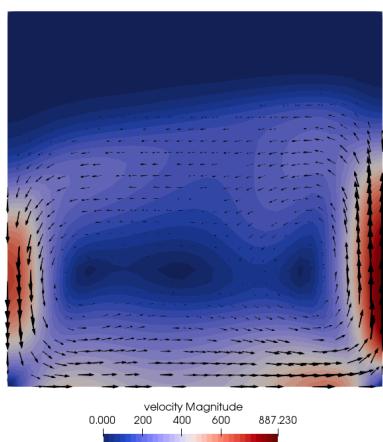
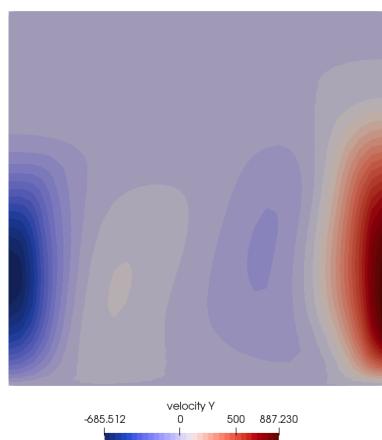
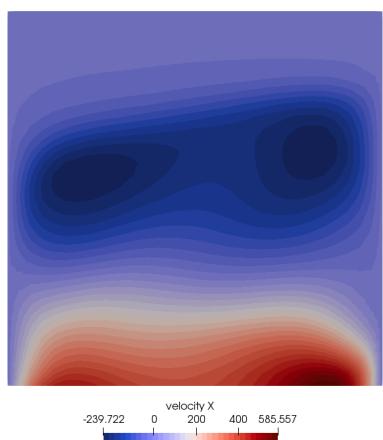


### 35.0.2 Case 1

In this case  $\mu^* = 0$  and  $\sigma_Y = 0$  so that  $\mu_{plast}$  can be discarded. The CFL number is set to 0.5 and the viscosity is given by  $\mu(T, z, \dot{\epsilon}) = \mu_{lin}(T, z)$ . And since  $\Delta\mu_z = 1$  then  $\gamma_z = 0$  so that  $\mu_{lin}(T, z) = \exp(-\gamma_T T)$

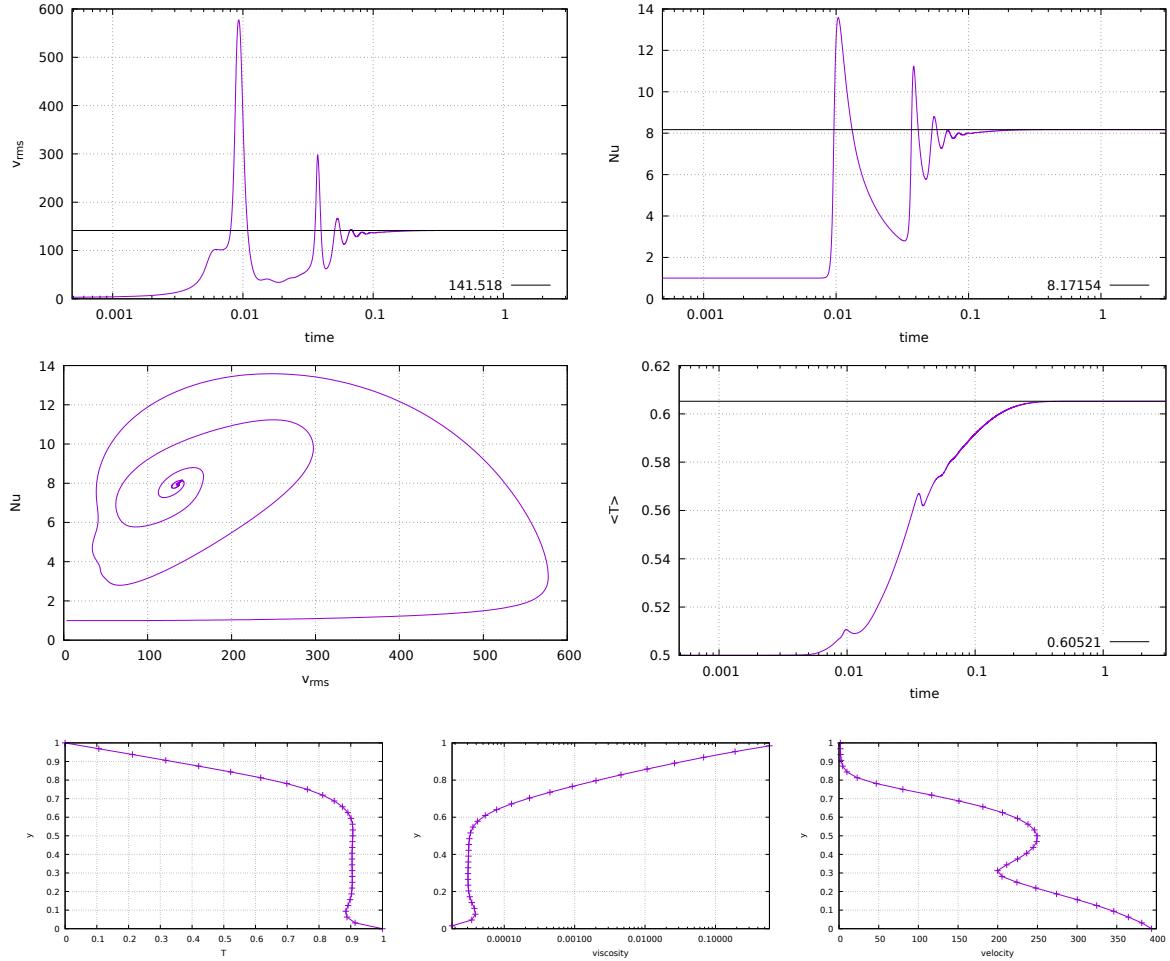
Code Resolution	YACC 100 × 100	Plaatjes 128×128r	CHIC 80 × 80	GAIA 100r×100r	StreamV 80 × 80	StagYY 128×128r	FEniCS 80 × 80	Fluidity 128 × 128	ELEFANT 100 × 100	ASPECT 64 × 64	MC3D 100 × 100
<i>Case 1</i>											
$\langle T \rangle$	0.7767	0.7759	0.7758	0.7759	0.776	0.776	0.7759	0.7758	0.7758	0.7768	0.779
$Nu_{top}$	3.4298	3.4159	3.4260	3.4213	3.4091	3.419	3.5889	3.4253	3.4214	3.4305	3.3129
$Nu_{bot}$	3.3143	3.4159	3.4259	3.4213	3.4091	3.419	3.4231	3.3795	3.313	3.4142	3.3139
$u_{rms}$	251.7997	249.54	249.2985	250.0738	252.0906	249.541	249.5730	248.9252	249.134	251.3069	296.6156
$u_{surf}^{rms}$	1.8298	1.878	1.8999	1.8836	1.8823	1.8723	1.8698	1.8474	1.8642	1.8695	1.1114
$u_{surf}^{max}$	2.5516	2.618	2.6477	2.6254	2.64	2.6104	2.6066	2.5761	2.6119	2.6064	1.5329
$\langle W \rangle$	2.4583	2.369	2.431	2.4121	2.4071	2.4189	2.4246	2.4148	2.4316	2.4282	2.5548
$\langle \Phi \rangle / Ra$	2.4333	2.4119	2.4189	2.4165	2.392	2.4182	2.4246	2.4148	2.4276	2.4281	2.31916
$\delta$	1.02%	1.78%	0.50%	0.18%	0.63%	0.03%	< 0.01%	< 0.01%	0.16%	< 0.01%	9.22%

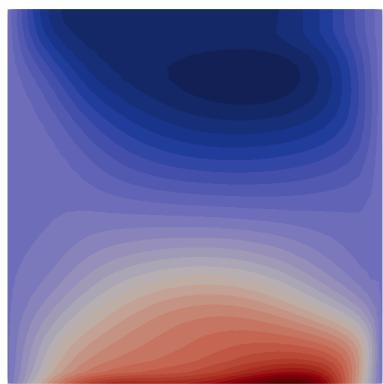




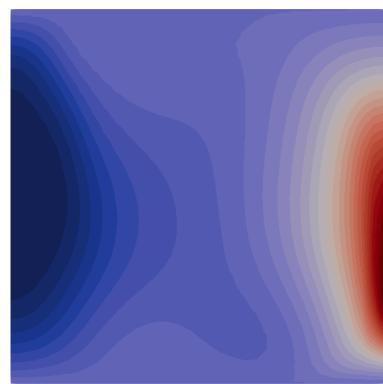
### 35.0.3 Case 2

$\langle T \rangle$	0.5289	0.5276	0.5276	0.5289	0.5283	0.527304	0.527521	0.5274	0.5277	0.5278	0.5364
$Nu_{top}$	6.5572	6.6156	6.6074	6.5913	6.6356	6.61082	6.68224	6.6401	6.5912	6.6249	7.4376
$Nu_{bot}$	6.5243	6.6158	6.6073	6.5913	6.6356	6.61082	6.66899	6.6326	6.5834	6.6267	7.4376
$U_{RMS}$	79.6202	79.1358	79.0181	78.6652	79.4334	78.9903	79.0684	79.0318	79.1105	79.1996	98.8912
$U_{surf}$	75.4814	75.1727	75.0434	74.1719	74.8587	75.0606	75.0975	75.0827	74.7596	75.1903	93.5057
$U_{max}$	89.2940	88.9715	88.8130	87.6118	89.1444	88.823	88.8753	88.85	88.9146	88.9848	123.046
$\eta_{min}$	$1.9174 \times 10^{-4}$	$1.9220 \times 10^{-4}$	$1.9448 \times 10^{-4}$	$1.9204 \times 10^{-4}$	$1.9900 \times 10^{-4}$	$1.9574 \times 10^{-4}$	$1.9167 \times 10^{-4}$	$1.9200 \times 10^{-4}$	$1.9860 \times 10^{-4}$	$1.9178 \times 10^{-4}$	$1.93 \times 10^{-4}$
$\eta_{max}$	1.6773	1.9834	1.6508	1.9670	1.1800	1.6665	1.7446	1.8891	1.5200	1.8831	$5.67 \times 10^{-3}$
$\langle W \rangle$	5.6512	5.6251	5.6076	5.5903	5.629	5.61012	5.61425	5.6136	5.6216	5.6235	6.751
$\langle \Phi \rangle / Ra$	5.6463	5.6174	5.6024	5.5434	5.6325	5.60152	5.61425	5.6136	5.6182	5.6235	7.786
$\delta$	0.09%	0.14%	0.09%	0.84%	0.06%	0.15%	<0.01%	<0.01%	0.06%	<0.01%	13.29%

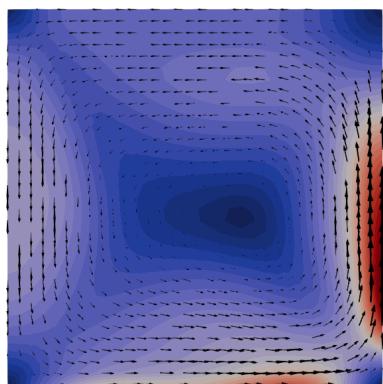




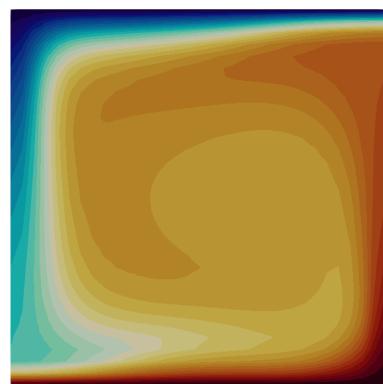
velocity X  
-1.437e+02 0 100 200 3.247e+02



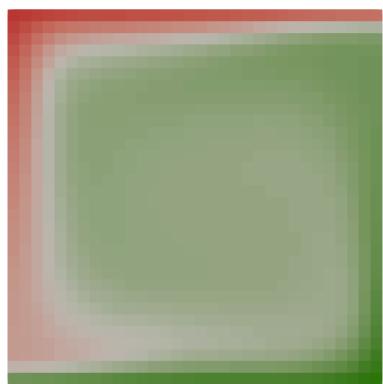
velocity Y  
-1.764e+02 0 200 4.549e+02



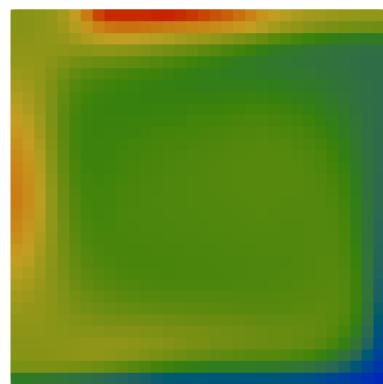
velocity Magnitude  
0.000e+00 100 200 300 4.549e+02



T  
0.000 0.2 0.4 0.6 0.8 1.000



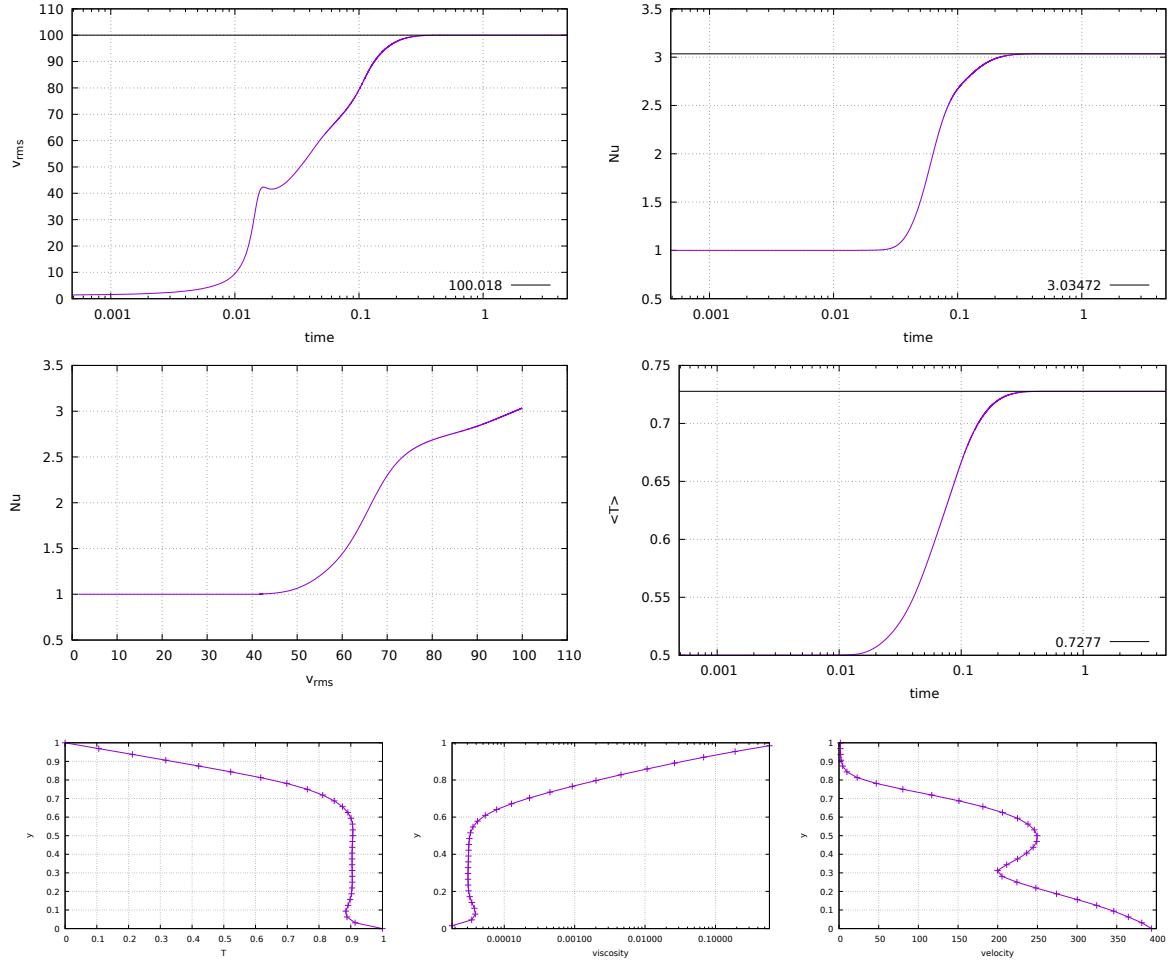
density  
-9.900e-05 -6e-5 -4e-5 -1e-5 -2.000e-06

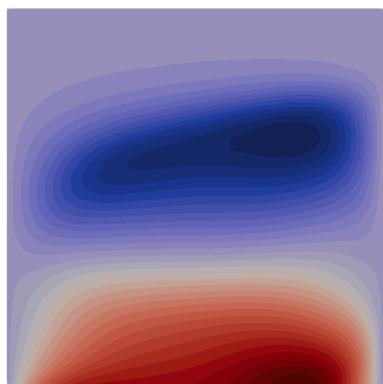


viscosity  
2.300e-05 0.00002 0.00003 0.00004 0.00005 6.326e-01

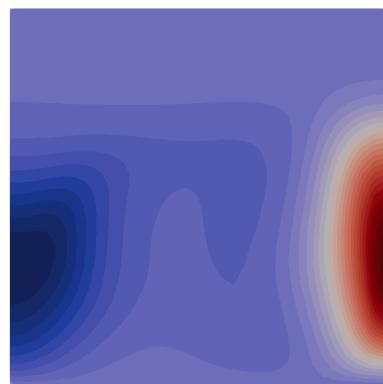
### 35.0.4 Case 3

$\langle T \rangle$	0.7286	0.7275	0.7271	0.7272	0.7241	0.7274	0.727464	0.7275	0.7275	0.7278	0.7305
$Nu_{top}$	3.0374	3.0298	3.0324	3.0314	3.0253	3.03025	3.0918	3.0399	3.0347	3.0371	2.9311
$Nu_{bot}$	2.9628	3.0298	3.0323	3.0314	3.0253	3.03025	3.03487	3.0376	2.9908	3.0410	2.9311
$u_{RMS}$	100.9467	100.024	99.8701	99.9917	100.197	100.018	100.127	100.0396	100.1208	100.3368	111.6121
$u_{surf}^{surf}$	2.0374	2.0785	2.0916	2.0835	2.0789	2.07299	2.07301	2.0569	2.0652	2.0727	1.356
$u_{max}$	2.8458	2.9029	2.9201	2.9094	2.9207	2.89495	2.89501	2.873	2.9019	2.8946	1.8806
$\eta_{min}$	$4.7907 \times 10^{-5}$	$4.8140 \times 10^{-5}$	$4.8014 \times 10^{-5}$	$4.8047 \times 10^{-5}$	$4.8400 \times 10^{-5}$	$4.7951 \times 10^{-5}$	$4.8081 \times 10^{-5}$	$4.8000 \times 10^{-5}$	$4.8080 \times 10^{-5}$	$4.7972 \times 10^{-5}$	$10^{-4}$
$\eta_{max}$	1	0.9987	0.9857	0.9988	0.9010	0.9637	0.9999	1	0.9023	1	1
$\langle W \rangle$	2.0400	2.0028	2.0340	2.0269	2.0235	2.03002	2.03482	2.0298	2.0384	2.0362	2.056
$\langle \Phi \rangle / Ra$	2.0335	2.0277	2.0304	2.0286	2.0164	2.0302	2.03482	2.0298	2.037	2.0362	1.865
$\delta$	0.32%	1.23%	0.18%	0.08%	0.35%	0.01%	<0.01%	<0.01%	0.07%	<0.01%	9.29%

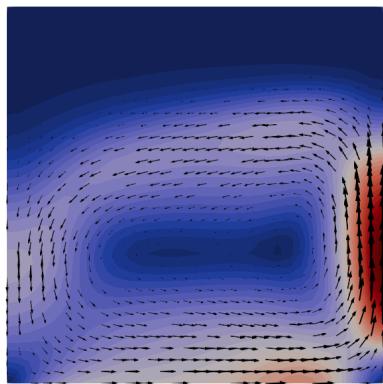




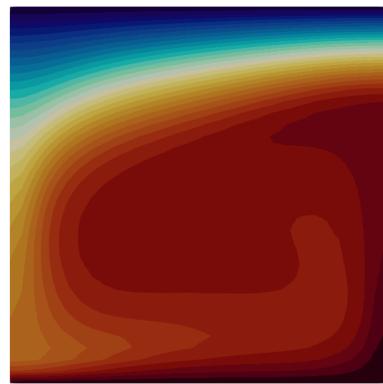
-133.934      0      100      204.210  
velocity X



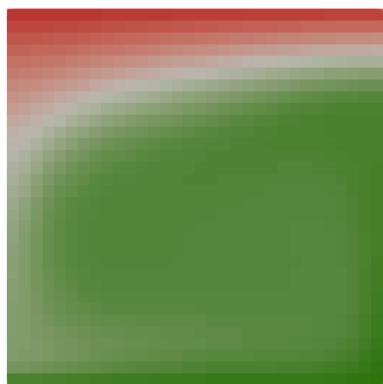
-137.969      0      100      200      330.512  
velocity Y



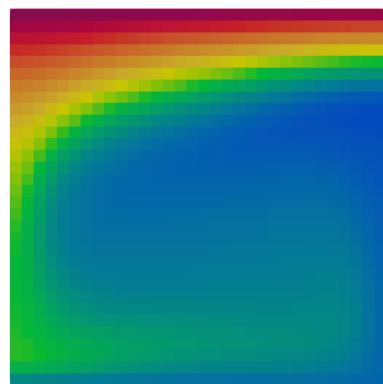
0.000      100      200      330.512  
velocity Magnitude



0.000      0.2      0.4      0.6      0.8      1.000  
T



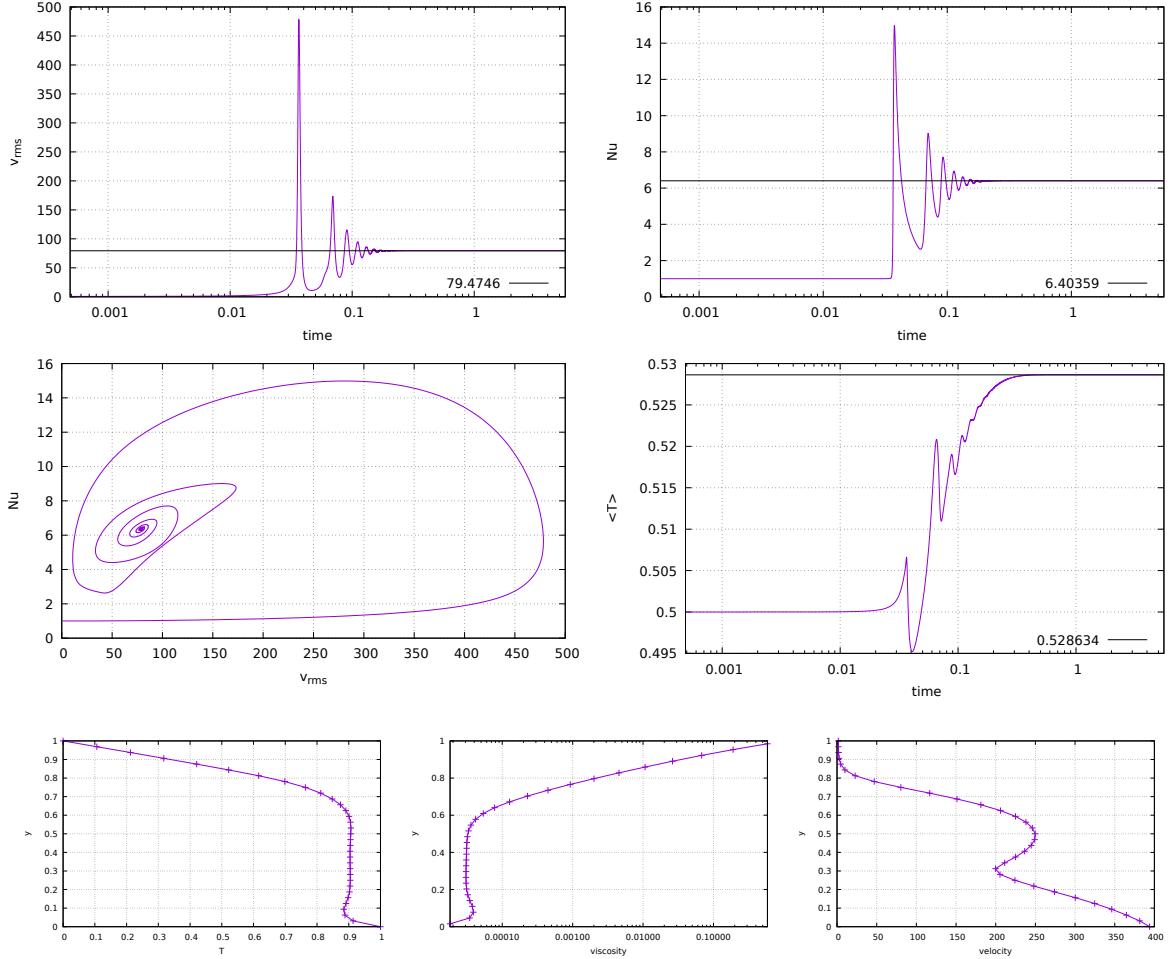
-1.000e-04      -6e-5      -4e-5      -3.000e-06  
density

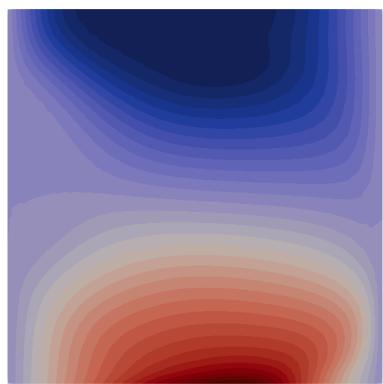


4.800e-05      0.000020000020.001      7.802e-01  
viscosity

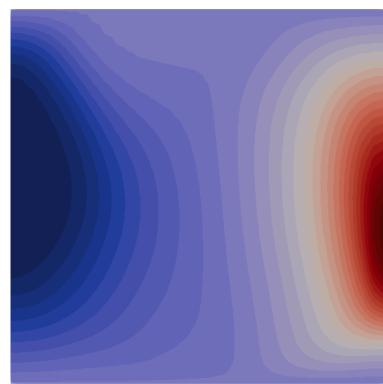
### 35.0.5 Case 4

$\langle T \rangle$	0.5289	0.5276	0.5276	0.5289	0.5283	0.527304	0.527521	0.5274	0.5277	0.5278	0.5364
$Nu_{top}$	6.5572	6.6156	6.6074	6.5913	6.6356	6.61082	6.68224	6.6401	6.5912	6.6249	7.4376
$Nu_{bot}$	6.5243	6.6158	6.6073	6.5913	6.6356	6.61082	6.66899	6.6326	6.5834	6.6267	7.4376
$U_{RMS}$	79.6202	79.1358	79.0181	78.6652	79.4334	78.9903	79.0684	79.0318	79.1105	79.1996	98.8912
$U_{RMS}^{surf}$	75.4814	75.1727	75.0434	74.1719	74.8587	75.0606	75.0975	75.0827	74.7596	75.1903	93.5057
$U_{max}^{surf}$	89.2940	88.9715	88.8130	87.6118	89.1444	88.823	88.8753	88.85	88.9146	88.9848	123.046
$\eta_{min}$	$1.9174 \times 10^{-4}$	$1.9220 \times 10^{-4}$	$1.9448 \times 10^{-4}$	$1.9204 \times 10^{-4}$	$1.9900 \times 10^{-4}$	$1.9574 \times 10^{-4}$	$1.9167 \times 10^{-4}$	$1.9200 \times 10^{-4}$	$1.9860 \times 10^{-4}$	$1.9178 \times 10^{-4}$	$1.93 \times 10^{-4}$
$\eta_{max}$	1.6773	1.9834	1.6508	1.9670	1.1800	1.6665	1.7446	1.8891	1.5200	1.8831	$5.67 \times 10^{-3}$
$\langle W \rangle$	5.6512	5.6251	5.6076	5.5903	5.629	5.61012	5.61425	5.6136	5.6216	5.6235	6.751
$\langle \Phi \rangle / Ra$	5.6463	5.6174	5.6024	5.5434	5.6325	5.60152	5.61425	5.6136	5.6182	5.6235	7.786
$\delta$	0.09%	0.14%	0.09%	0.84%	0.06%	0.15%	<0.01%	<0.01%	0.06%	<0.01%	13.29%

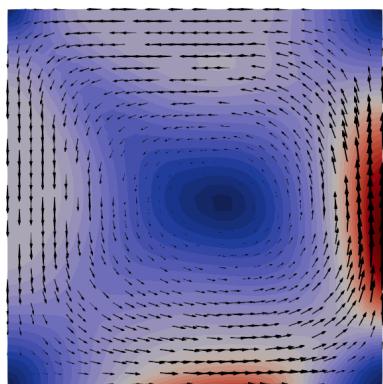




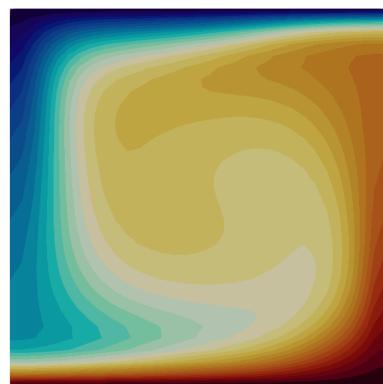
velocity X  
-92.305 -50 0 50 100 153.873



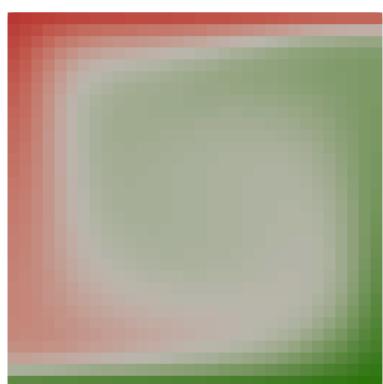
velocity Y  
-97.677 0 100 209.053



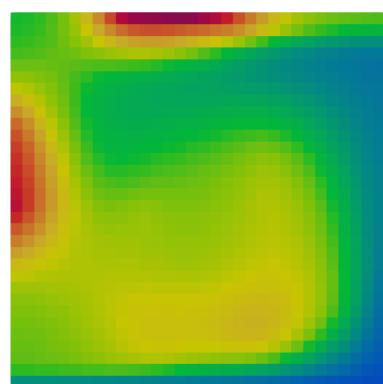
velocity Magnitude  
0.000 50 100 150 209.053



T  
0.000 0.2 0.4 0.6 0.8 1.000

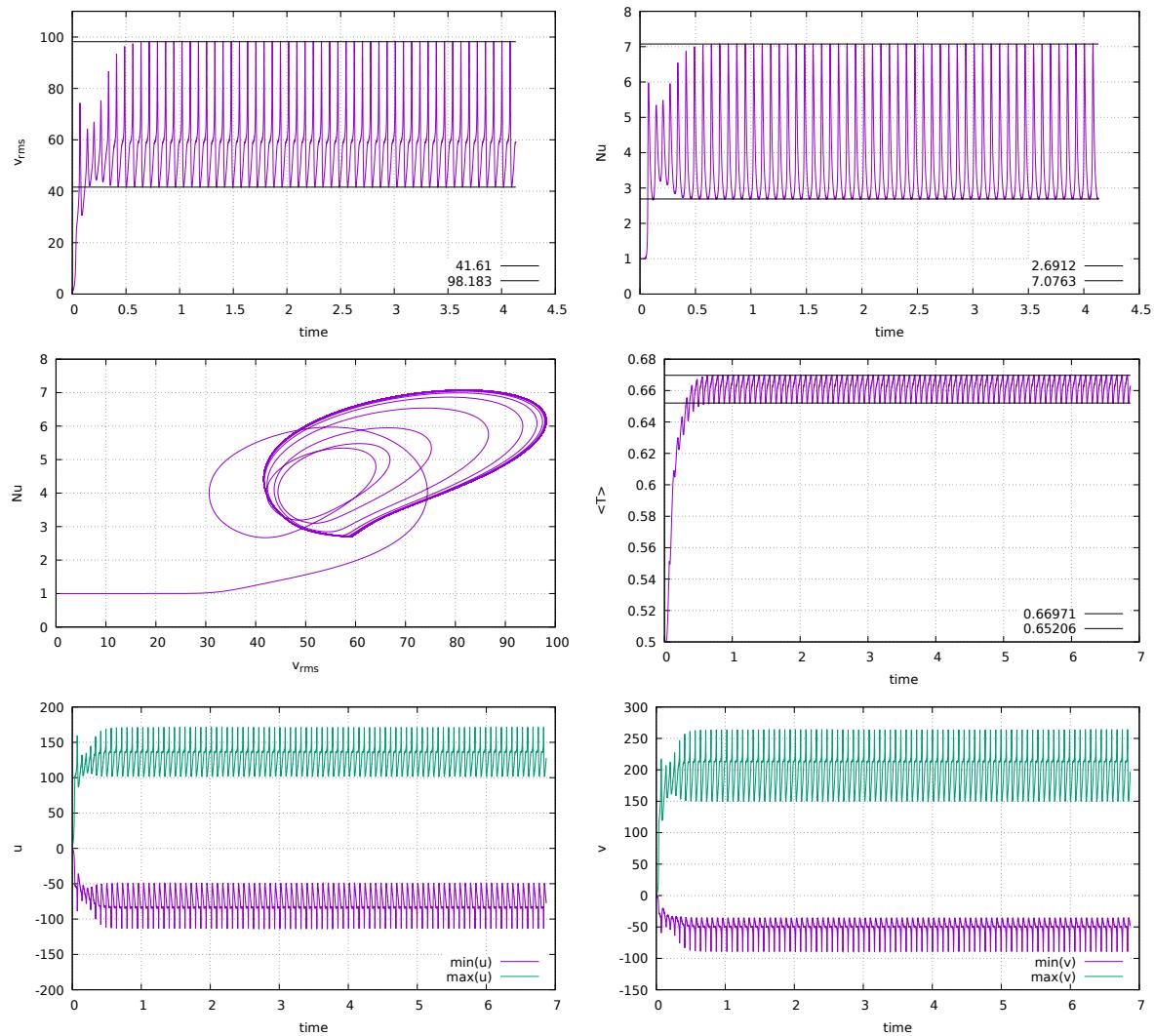


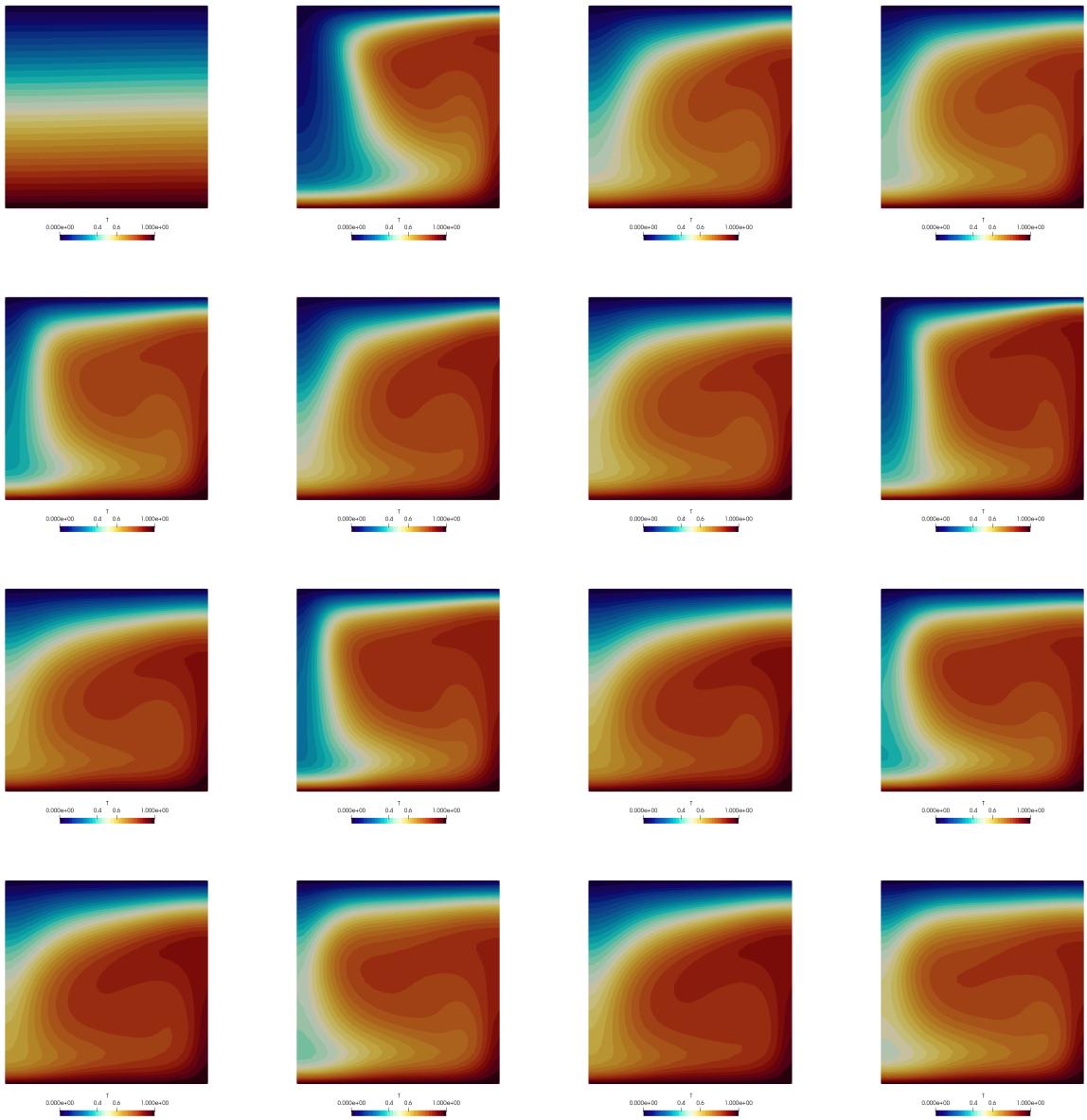
density  
-9.900e-05 -5e-05 -4e-05 -1.000e-06



viscosity  
0.00021 0.00203 0.00205 0.1 0.75499

### 35.0.6 Case 5

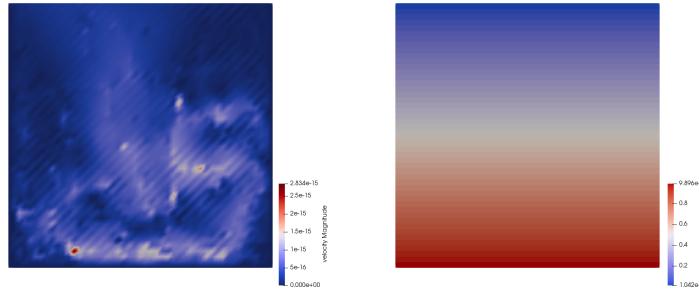




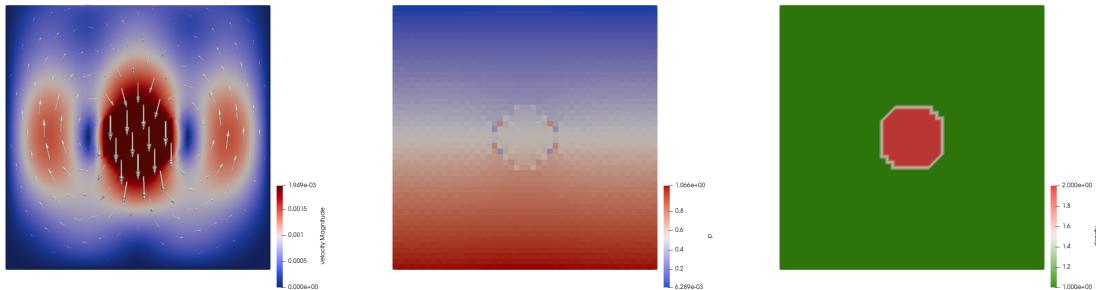
## 36 fieldstone\_29: open boundary conditions

In what follows we will investigate the use of the so-called open boundary conditions in the very simple context of a 2D Stokes sphere experiment.

We start with a domain without the sphere. Essentially, it is what people would call an aquarium. Free slip boundary conditions are prescribed on the sides and no-slip conditions at the bottom. The top surface is left free. The fluid has a density  $\rho_0 = 1$  and a viscosity  $\eta_0 = 1$ . In the absence of any density difference in the domain there is no active buoyancy force so that we expect a zero velocity field and a lithostatic pressure field. This is indeed what we recover:



If we now implement a sphere parametrised by its density  $\rho_s = \rho_0 + 1$ , its viscosity  $\eta_s = 10^3$  and its radius  $R_s = 0.123$  in the middle of the domain, we see clear velocity field which logically shows the sphere falling downward and a symmetric return flow of the fluid on each side:

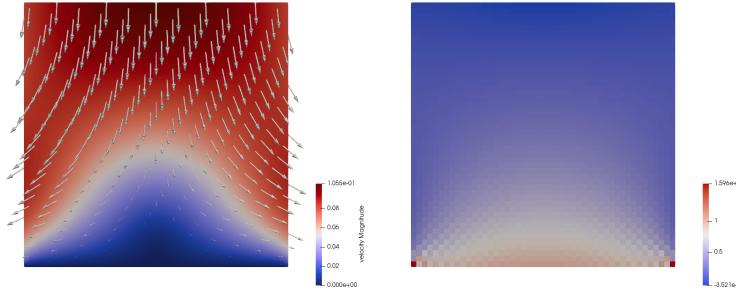


Unfortunately it has been widely documented that the presence of free-slip boundary conditions affects the evolution of subduction [127], even when these are placed rather far from the subduction zone. A proposed solution to this problem is the use of 'open boundary conditions' which are in fact stress boundary conditions. The main idea is to prescribe a stress on the lateral boundaries (instead of free slip) so that it balances out exactly the existing lithostatic pressure inside the domain along the side walls. Only pressure deviations with respect to the lithostatic are responsible for flow and such boundary conditions allow flow across the boundaries.

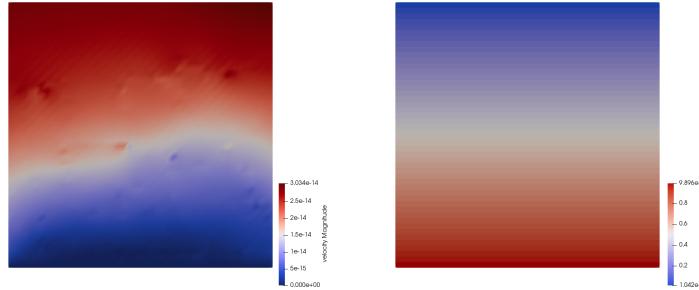
We need the lithostatic pressure and compute it before hand (which is trivial in our case but can prove to be a bit more tedious in real life situations when for instance density varies in the domain as a function of temperature and/or pressure).

```
plith = np.zeros(nnp, dtype=np.float64)
for i in range(0,nnp):
 plith[i]=(Ly-y[i])*rho0*abs(gy)
```

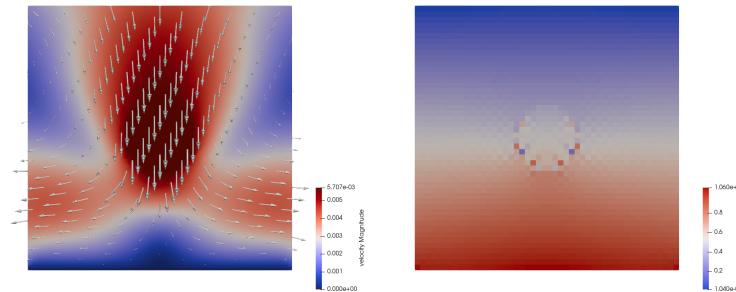
Let us start with a somewhat pathological case: even in the absence of the sphere, what happens when no boundary conditions are prescribed on the sides? The answer is simple: think about an aquarium without side walls, or a broken dam. The velocity field indeed shows a complete collapse of the fluid left and right of the bottom.



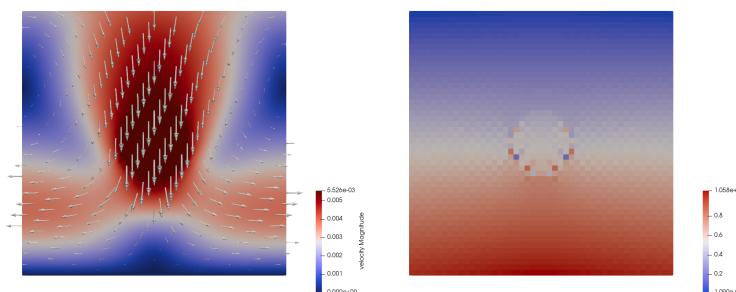
Let us then continue (still with no sphere) but let us now switch on the open boundary conditions. Since the side boundary conditions match the lithostatic pressure we expect no flow at all in the absence of any density perturbation in the system. This is indeed what is recovered:



Finally, let us reintroduce the sphere. This time flow is allowed through the left and right side boundaries:



Finally, although horizontal velocity Dirichlet boundary conditions and open boundary conditions are not compatible, the same is not true for the vertical component of the velocity: the open b.c. implementation acts on the horizontal velocity dofs only, so that one can fix the vertical component to zero, as is shown hereunder:



We indeed see that the in/outflow on the sides is perpendicular to the boundaries.

Turning now to the actual implementation, we see that it is quite trivial, since all element edges are vertical, and all have the same vertical dimension  $h_x$ . Since we use a  $Q_0$  approximation for the pressure

we need to prescribe a single pressure value in the middle of the element. Finally because of the sign of the normal vector projection onto the  $x$ -axis, we obtain:

```
if open_bc_left and x[icon[0, iel]]<eps: # left side
 pmid=0.5*(plith[icon[0, iel]]+plith[icon[3, iel]])
 f_el[0]+=0.5*hy*pmid
 f_el[6]+=0.5*hy*pmid
if open_bc_right and x[icon[1, iel]]>Lx-eps: # right side
 pmid=0.5*(plith[icon[1, iel]]+plith[icon[2, iel]])
 f_el[2]=-0.5*hy*pmid
 f_el[4]=-0.5*hy*pmid
```

These few lines of code are added after the elemental matrices and rhs are built, and before the application of other Dirichlet boundary conditions, and assembly.

#### features

- $Q_1 \times P_0$  element
- incompressible flow
- mixed formulation
- open boundary conditions
- isoviscous

## 37 fieldstone\_30: conservative velocity interpolation

In this the Stokes equations are not solved. It is a 2D implementation of the cvi algorithm as introduced in [587] which deals with the advection of markers.  $Q_1$  and  $Q_2$  basis functions are used and in both cases the cvi algorithm can be toggled on/off. Markers can be distributed regularly or randomly at startup.

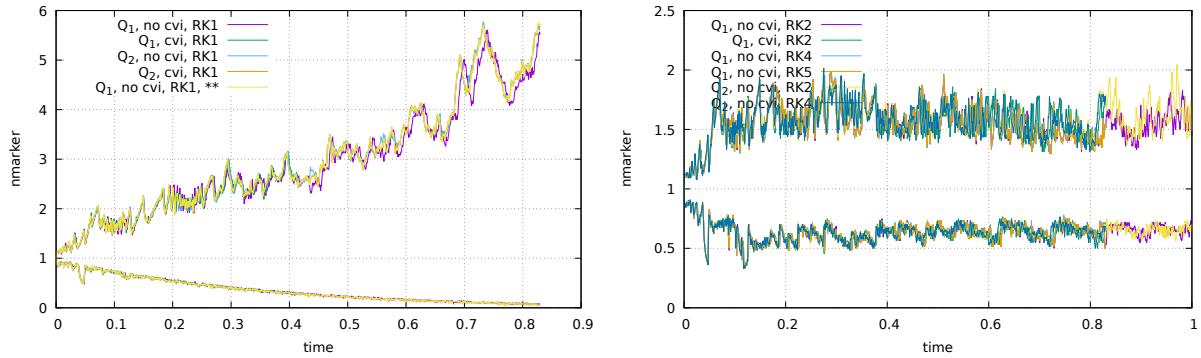
Three velocity fields are prescribed on the mesh:

- the so-called Couette flow of [587]
- the SolCx solution
- a flow created by means of a stream line function (see fieldstone 32)

### 37.1 Couette flow

### 37.2 SolCx

### 37.3 Streamline flow



In this case RK order seems to be more important than cvi.

Explore why ?!

#### features

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions (free-slip)
- direct solver
- isothermal
- non-isoviscous
- analytical solution

## 38 fieldstone\_31: conservative velocity interpolation 3D

## 39 fieldstone\_32: 2D analytical sol. from stream function

### 39.1 Background theory

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  $\Psi = \text{const}$  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

#### REFERENCES.

In two dimensions the velocity is obtained as follows:

$$\mathbf{v} = \left( \frac{\partial \Psi}{\partial y}, -\frac{\partial \Psi}{\partial x} \right) \quad (474)$$

Provided the function  $\Psi$  is a smooth enough function, this automatically insures that the flow is incompressible:

$$\nabla \cdot \mathbf{v} = \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = \frac{\partial^2 \Psi}{\partial x \partial y} - \frac{\partial^2 \Psi}{\partial y \partial x} = 0 \quad (475)$$

Assuming constant viscosity, the Stokes equation writes:

$$-\nabla p + \mu \Delta \mathbf{v} = \rho \mathbf{g} \quad (476)$$

Let us introduce the vector  $\mathbf{W}$  for convenience such that in each dimension:

$$W_x = -\frac{\partial p}{\partial x} + \mu \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial x \partial y} \right) = \rho g_x \quad (477)$$

$$W_y = -\frac{\partial p}{\partial y} + \mu \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial x \partial y} \right) = \rho g_y \quad (478)$$

Taking the curl of the vector  $\mathbf{W}$  and only considering the component perpendicular to the  $xy$ -plane:

$$\frac{\partial V_y}{\partial x} - \frac{\partial V_x}{\partial y} = \frac{\partial \rho g_y}{\partial x} - \frac{\partial \rho g_x}{\partial y} \quad (479)$$

The advantage of this approach is that the pressure terms cancel out (the curl of a gradient is always zero), so that:

$$\frac{\partial}{\partial x} \mu \left( \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial x \partial y} \right) - \frac{\partial}{\partial y} \mu \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial x \partial y} \right) \right) = \frac{\partial \rho g_y}{\partial x} - \frac{\partial \rho g_x}{\partial y} \quad (480)$$

and then replacing  $u, v$  by their stream function derivatives yields (for a constant viscosity):

$$\mu \left( \frac{\partial^4 \Psi}{\partial x^4} + \frac{\partial^4 \Psi}{\partial y^4} + 2 \frac{\partial^4 \Psi}{\partial x^2 \partial y^2} \right) = \frac{\partial \rho g_y}{\partial x} - \frac{\partial \rho g_x}{\partial y} \quad (481)$$

or,

$$\nabla^4 \Psi = \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \Psi = \frac{\partial \rho g_y}{\partial x} - \frac{\partial \rho g_x}{\partial y} \quad (482)$$

These equations are also to be found in the geodynamics literature, eee Eq. 1.43 of Tackley book, p 70-71 of Gerya book.

### 39.2 A simple application

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  $\mu = 1$  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:

$$\Psi(x, y) = \sin(m\pi x) \sin(n\pi y) \quad (483)$$

We have the velocity being defined as:

$$\mathbf{v} = (u, v) = \left( \frac{\partial \Psi}{\partial y}, -\frac{\partial \Psi}{\partial x} \right) = (n\pi \sin(m\pi x) \cos(n\pi y), -m\pi \cos(m\pi x) \sin(n\pi y)) \quad (484)$$

The strain rate components are then:

$$\dot{\varepsilon}_{xx} = \frac{\partial u}{\partial x} = mn\pi^2 \cos(m\pi x) \cos(n\pi y) \quad (485)$$

$$\dot{\varepsilon}_{yy} = \frac{\partial v}{\partial y} = -mn\pi^2 \cos(m\pi x) \cos(n\pi y) \quad (486)$$

$$2\dot{\varepsilon}_{xy} = \frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} \quad (487)$$

$$= \frac{\partial^2 \Psi}{\partial y^2} - \frac{\partial^2 \Psi}{\partial x^2} \quad (488)$$

$$= -n^2\pi^2\Psi + m^2\pi^2\Psi \quad (489)$$

$$= (m^2 - n^2)\pi^2 \sin(m\pi x) \sin(n\pi y) \quad (490)$$

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  $\dot{\varepsilon}_{xy}$  is null on all four sides, thereby guaranteeing free slip.

Our choice of stream function yields:

$$\nabla^4 \Psi = \frac{\partial^4 \Psi}{\partial x^4} + \frac{\partial^4 \Psi}{\partial y^4} + 2 \frac{\partial^2 \Psi}{\partial x^2 \partial y^2} = \pi^4 (m^4 \Psi + n^4 \Psi + 2m^2 n^2 \Psi) = (m^4 + n^4 + 2m^2 n^2) \pi^4 \Psi$$

We assume  $g_x = 0$  and  $g_y = -1$  so that we simply have

$$(m^4 + n^4 + 2m^2 n^2) \pi^4 \Psi = -\frac{\partial \rho}{\partial x} \quad (491)$$

so that (assuming the integration constant to be zero):

$$\rho(x, y) = \frac{m^4 + n^4 + 2m^2 n^2}{m} \pi^3 \cos(m\pi x) \sin(n\pi y)$$

The  $x$ -component of the momentum equation is

$$-\frac{\partial p}{\partial x} + \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = -\frac{\partial p}{\partial x} - m^2 n \pi^3 \sin(m\pi x) \cos(n\pi y) - n^3 \pi^3 \sin(m\pi x) \cos(n\pi y) = 0$$

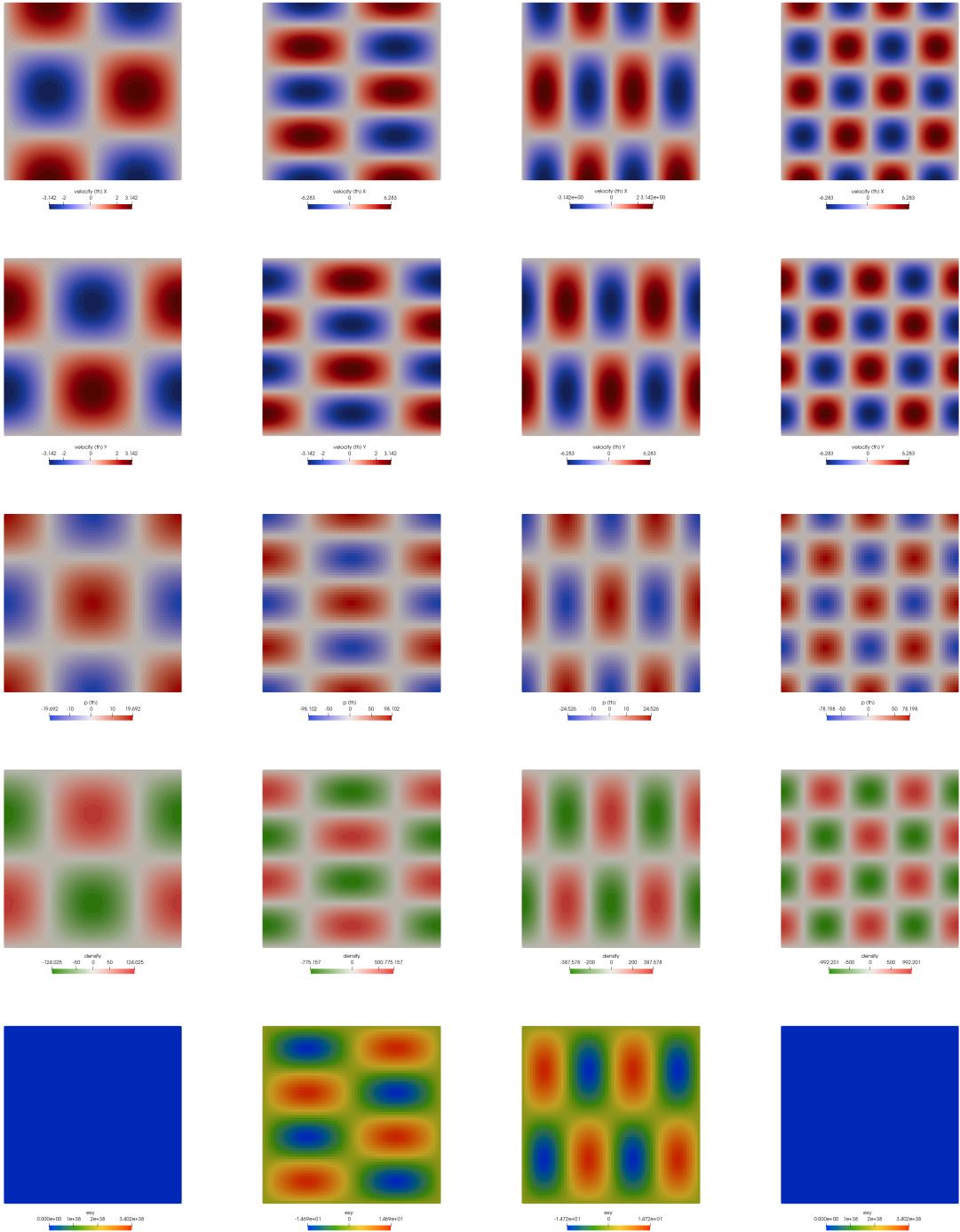
so

$$\frac{\partial p}{\partial x} = -(m^2 n + n^3) \pi^3 \sin(m\pi x) \cos(n\pi y)$$

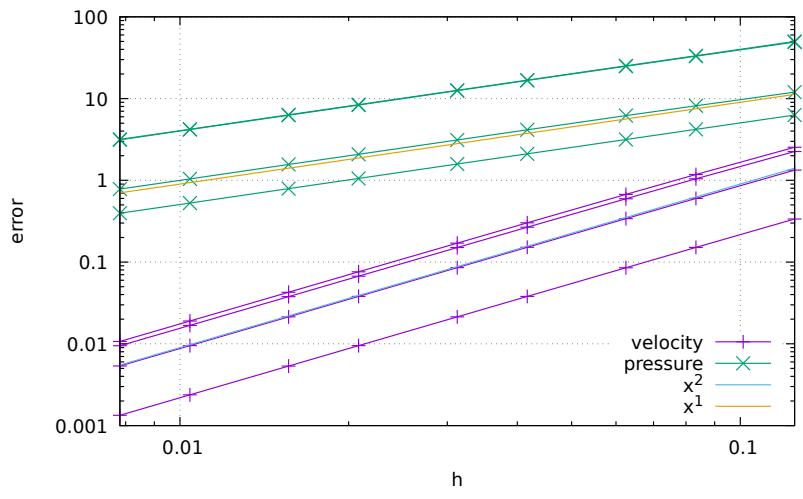
and the pressure field is then (once again neglecting the integration constant):

$$p(x, y) = \frac{m^2 n + n^3}{m} \pi^2 \cos(\pi x) \cos(\pi y)$$

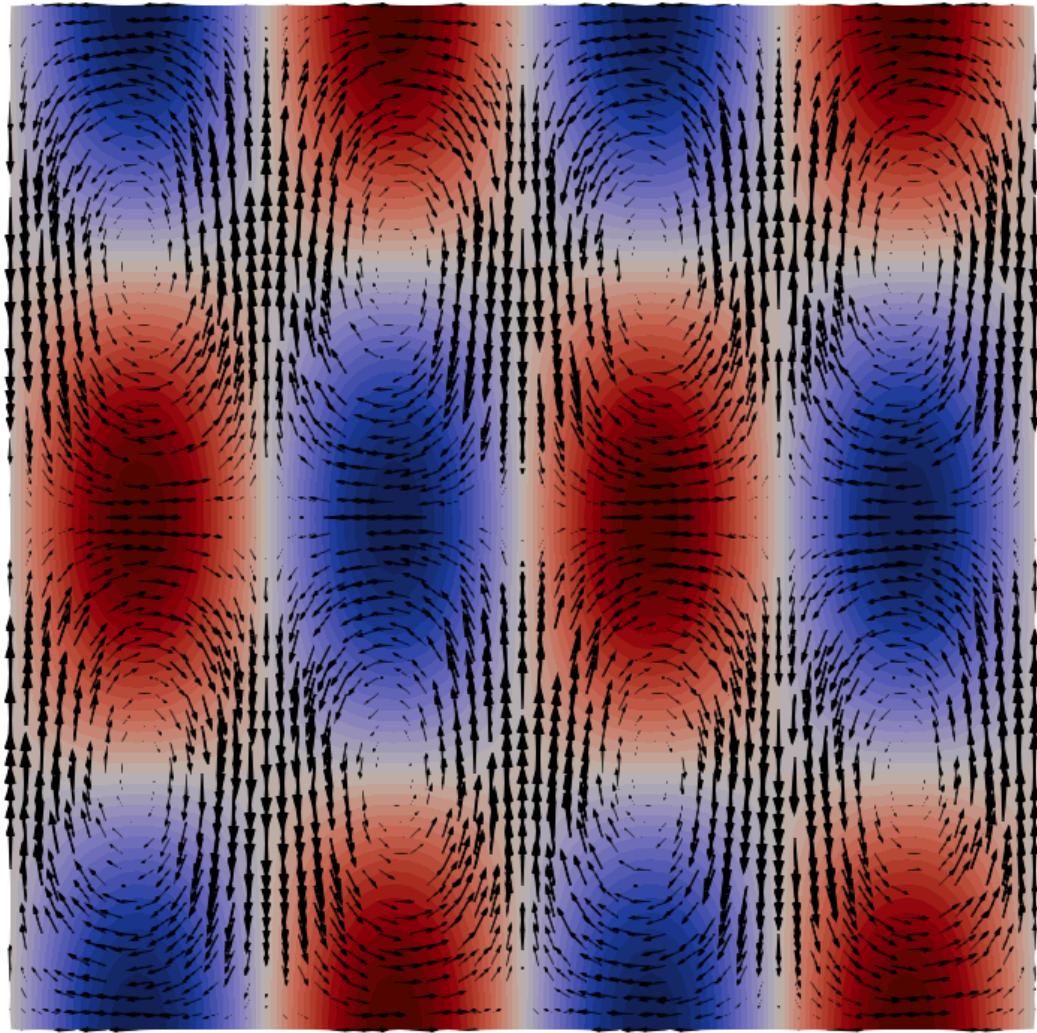
Note that in this case  $\int p dV = 0$  so that volume normalisation of the pressure is turned on (when free slip boundary conditions are prescribed on all sides the pressure is known up to a constant and this undeterminacy can be lifted by adding an additional constraint to the pressure field).



Top to bottom: Velocity components  $u$  and  $v$ , pressure  $p$ , density  $\rho$  and strain rate  $\dot{\epsilon}_{xy}$ . From left to right:  
 $(m, n) = (1, 1)$ ,  $(m, n) = (1, 2)$ ,  $(m, n) = (2, 1)$ ,  $(m, n) = (2, 2)$



Errors for velocity and pressure for  $(m, n) = (1, 1), (1, 2), (2, 1), (2, 2)$



Velocity arrows for  $(m, n) = (2, 1)$

## 40 fieldstone\_33: Convection in an annulus



*This fieldstone was developed in collaboration with Rens Elbertsen.*

This is based on the community benchmark for viscoplastic thermal convection in a 2D square box [550] as already carried out in ??.

In this experiment the geometry is an annulus of inner radius  $R_1 = 1.22$  and outer radius  $R_2 = 2.22$ . The rheology and buoyancy forces are identical to those of the box experiment. The initial temperature is now given by:

$$T(r, \theta) = T_c(r) + A s(1 - s) \cos(N_0 \theta) \quad s = \frac{R_2 - r}{R_2 - R_1} \in [0, 1]$$

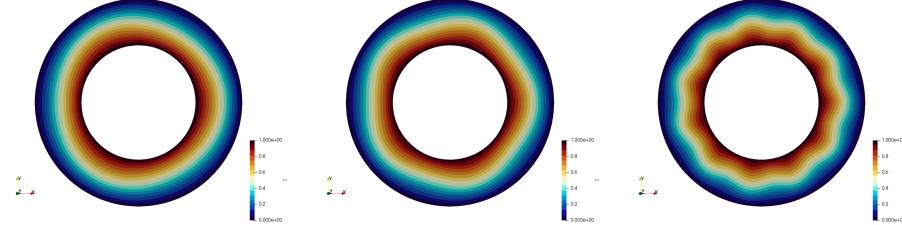
where  $s$  is the normalised depth,  $A$  is the amplitude of the perturbation and  $N_0$  the number of lobes. In this equation  $T_c(r)$  stands for the steady state purely conductive temperature solution which is obtained by solving the Laplace's equation in polar coordinates (all terms in  $\theta$  are dropped because of radial symmetry) supplemented with two boundary conditions:

$$\Delta T_c = \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial T}{\partial r} \right) = 0 \quad T(r = R_1) = T_1 = 1 \quad T(r = R_2) = T_2 = 0$$

We obtain

$$T_c(r) = \frac{\log(r/R_2)}{\log(R_1/R_2)}$$

Note that this profile differs from the straight line that is used in [550] and in section 35.



Examples of initial temperature fields for  $N_0 = 3, 5, 11$

Boundary conditions can be either no-slip or free-slip on both inner and outer boundary. However, when free-slip is used on both a velocity null space exists and must be filtered out. In other words, the solver may be able to come up with a solution to the Stokes operator, but that solution plus an arbitrary rotation is also an equally valid solution. This additional velocity field can be problematic since it is used for advecting temperature (and/or compositions) and it also essentially determines the time step value for a chosen mesh size (CFL condition).

For these reasons the nullspace must be removed from the obtained solution after every timestep. There are two types of nullspace removal: removing net angular momentum, and removing net rotations.

We calculate the following output parameters:

- the average temperature  $\langle T \rangle$

$$\langle T \rangle = \frac{\int_{\Omega} T d\Omega}{\int_{\Omega} d\Omega} = \frac{1}{V_{\Omega}} \int_{\Omega} T d\Omega \quad (492)$$

- the root mean square velocity  $v_{rms}$  as given by equation (27).
- the root mean square of the radial and tangential velocity components as given by equations (29) and (30).

- the heat transfer through both boundaries  $Q$ :

$$Q_{inner,outer} = \int_{\Gamma_{i,o}} \mathbf{q} \cdot \mathbf{n} \, d\Gamma \quad (493)$$

- the Nusselt number at both boundaries  $Nu$  as given by equations (32) and (33).
- the power spectrum of the temperature field:

$$PS_n(T) = \left| \int_{\Omega} T(r, \theta) e^{in\theta} d\Omega \right|^2. \quad (494)$$

**features**

- $Q_1 \times P_0$  element
- incompressible flow
- penalty formulation
- Dirichlet boundary conditions
- non-isothermal
- non-isoviscous
- annulus geometry

## 41 fieldstone\_34: the Cartesian geometry elastic aquarium

This fieldstone was developed in collaboration with Lukas van de Wiel.

The setup is as follows: a 2D square of elastic material of size  $L$  is subjected to the following boundary conditions: free slip on the sides, no slip at the bottom and free at the top. It has a density  $\rho$  and is placed in a gravity field  $\mathbf{g} = -ge_y$ . For an isotropic elastic medium the stress tensor is given by:

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

where  $\lambda$  is the Lamé parameter and  $\mu$  is the shear modulus. The displacement field is  $\mathbf{u} = (0, u_y(y))$  because of symmetry reasons (we do not expect any of the dynamic quantities to depend on the  $x$  coordinate and also expect the horizontal displacement to be exactly zero). The velocity divergence is then  $\nabla \cdot \mathbf{u} = \partial u_y / \partial y$  and the strain tensor:

$$\boldsymbol{\varepsilon} = \begin{pmatrix} 0 & 0 \\ 0 & \frac{\partial u_y}{\partial y} \end{pmatrix}$$

so that the stress tensor is:

$$\boldsymbol{\sigma} = \begin{pmatrix} \lambda \frac{\partial u_y}{\partial y} & 0 \\ 0 & (\lambda + 2\mu) \frac{\partial u_y}{\partial y} \end{pmatrix}$$

$$\nabla \cdot \boldsymbol{\sigma} = (\partial_x \quad \partial_y) \cdot \begin{pmatrix} \lambda \frac{\partial u_y}{\partial y} & 0 \\ 0 & (\lambda + 2\mu) \frac{\partial u_y}{\partial y} \end{pmatrix} = \begin{pmatrix} 0 \\ (\lambda + 2\mu) \frac{\partial^2 u_y}{\partial y^2} \end{pmatrix} = \begin{pmatrix} 0 \\ \rho g \end{pmatrix}$$

so that the vertical displacement is then given by:

$$u_y(y) = \frac{1}{2} \frac{\rho g}{\lambda + 2\mu} y^2 + \alpha y + \beta$$

where  $\alpha$  and  $\beta$  are two integration constants. We need now to use the two boundary conditions: the first one states that the displacement is zero at the bottom, i.e.  $u_y(y=0) = 0$  which immediately implies  $\beta = 0$ . The second states that the stress at the top is zero (free surface), which implies that  $\partial u_y / \partial y(y=L) = 0$  which allows us to compute  $\alpha$ . Finally:

$$u_y(y) = \frac{\rho g}{\lambda + 2\mu} \left( \frac{y^2}{2} - Ly \right)$$

The pressure is given by

$$p = -(\lambda + \frac{2}{3}\mu)\nabla \cdot \mathbf{u} = (\lambda + \frac{2}{3}\mu) \frac{\rho g}{\lambda + 2\mu} (L - y) = \frac{\lambda + \frac{2}{3}\mu}{\lambda + 2\mu} \rho g (L - y) = \frac{1 + \frac{2\mu}{3\lambda}}{1 + 2\mu/\lambda} \rho g (L - y)$$

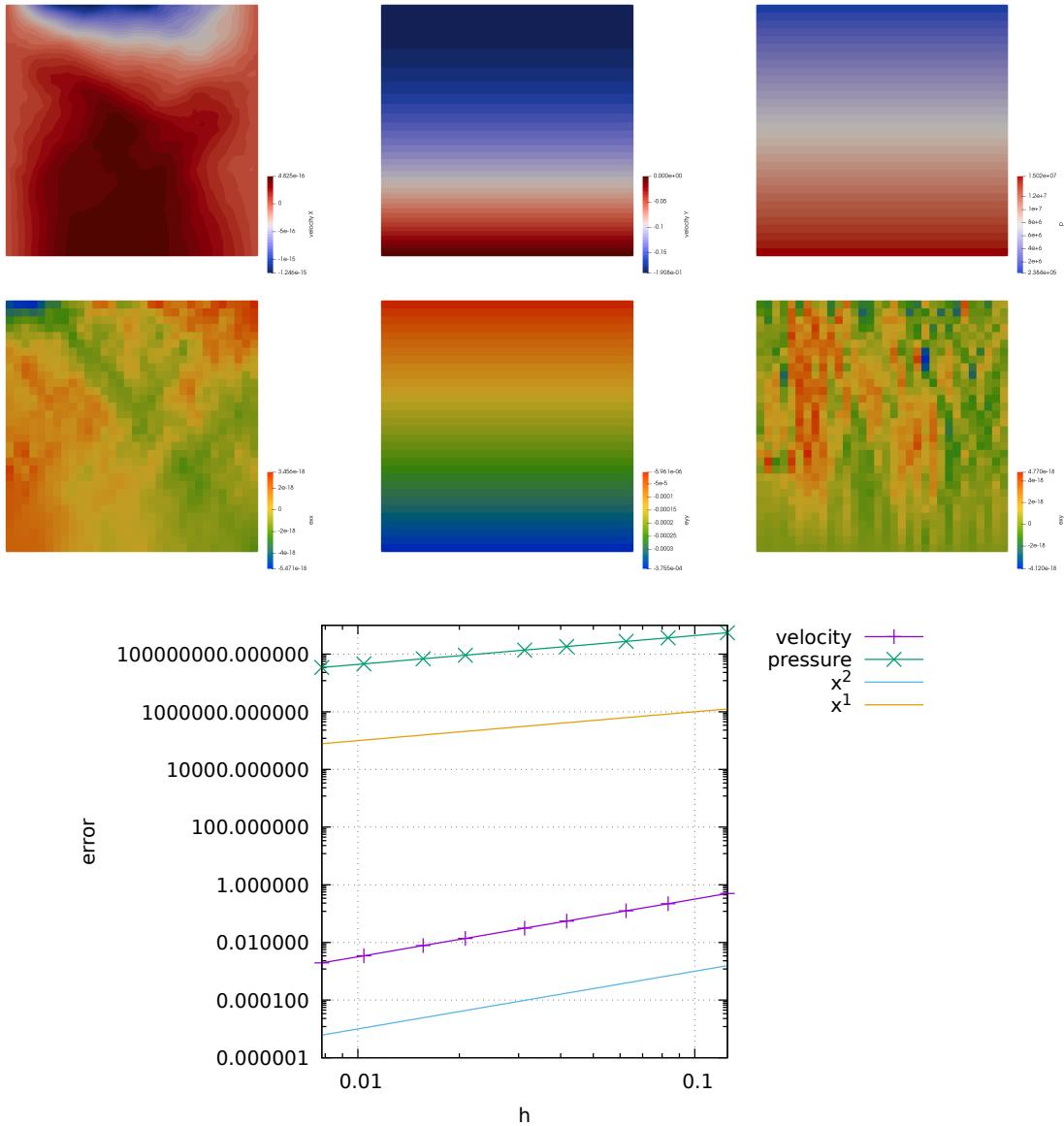
In the incompressible limit, the poisson ratio is  $\nu \sim 0.5$ . Materials are characterised by a finite Young's modulus  $E$ , which is related to  $\nu$  and  $\lambda$ :

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

It is then clear that for incompressible parameters  $\lambda$  becomes infinite while  $\mu$  remains finite. In that case the pressure then logically converges to the well known formula:

$$p = \rho g(L - y)$$

In what follows we set  $L = 1000\text{m}$ ,  $\rho = 2800$ ,  $\nu = 0.25$ ,  $E = 6 \cdot 10^{10}$ ,  $g = 9.81$ .



## 42 fieldstone\_35: 2D analytical sol. in annulus from stream function



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  $\mathbf{e}_r$  and  $\mathbf{e}_\theta$ , the inner radius of the annulus by  $R_1$  and the outer radius by  $R_2$ . Further, we assume that the viscosity  $\mu$  is constant, which we set to  $\mu = 1$  we set the gravity vector to  $\mathbf{g} = -g_r \mathbf{e}_r$  with  $g_r = 1$ .

Given these assumptions, the incompressible Stokes equations in the annulus are [500]

$$A_r = \frac{\partial^2 v_r}{\partial r^2} + \frac{1}{r} \frac{\partial v_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 v_r}{\partial \theta^2} - \frac{v_r}{r^2} - \frac{2}{r^2} \frac{\partial u_\theta}{\partial \theta} - \frac{\partial p}{\partial r} = \rho g_r \quad (495)$$

$$A_\theta = \frac{\partial^2 v_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial v_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 v_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r^2} - \frac{1}{r} \frac{\partial p}{\partial \theta} = 0 \quad (496)$$

$$\frac{1}{r} \frac{\partial(rv_r)}{\partial r} + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} = 0 \quad (497)$$

Equations (495) and (496) are the momentum equations in polar coordinates while Equation (497) is the incompressibility constraint. The components of the velocity are obtained from the stream function as follows:

$$v_r = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \quad v_\theta = -\frac{\partial \Psi}{\partial r}$$

where  $v_r$  is the radial component and  $v_\theta$  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  $\mathbf{u} \cdot \nabla \Psi = 0$  using the formula for  $\mathbf{u}$  in terms of  $\Psi$  which proves that level curves of  $\Psi$  are streamlines:

$$\mathbf{u} \cdot \nabla \Psi = v_r \frac{\partial \Psi}{\partial r} + v_\theta \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \frac{\partial \Psi}{\partial r} - \frac{\partial \Psi}{\partial r} \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = 0$$

In polar coordinates the curl of a vector  $\mathbf{A}$  is:

$$\nabla \times \mathbf{A} = \frac{1}{r} \left( \frac{\partial(rA_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} \right)$$

Taking the curl of vector  $\mathbf{A}$  yields:

$$\frac{1}{r} \left( \frac{\partial(rA_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} \right) = \frac{1}{r} \left( -\frac{\partial(\rho g_r)}{\partial \theta} \right)$$

Multiplying on each side by  $r$

$$\frac{\partial(rA_\theta)}{\partial r} - \frac{\partial A_r}{\partial \theta} = -\frac{\partial \rho g_r}{\partial \theta}$$

If we now replace  $A_r$  and  $A_\theta$  by their expressions as a function of velocity and pressure, we see that the pressure terms cancel out and assuming the viscosity and the gravity vector to be constant we get: Let us assume the following separation of variables  $\boxed{\Psi(r, \theta) = \phi(r)\xi(\theta)}$ . Then

$$v_r = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{\phi \xi'}{r} \quad v_\theta = -\frac{\partial \Psi}{\partial r} = -\phi' \xi$$

Let us first express  $A_r$  and  $A_\theta$  as functions of  $\Psi$  and

$$A_r = \frac{\partial^2 v_r}{\partial r^2} + \frac{1}{r} \frac{\partial v_r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 v_r}{\partial \theta^2} - \frac{v_r}{r^2} - \frac{2}{r^2} \frac{\partial u_\theta}{\partial \theta} \quad (498)$$

$$= \frac{\partial^2}{\partial r^2} \left( \frac{\phi \xi'}{r} \right) + \frac{1}{r} \frac{\partial}{\partial r} \left( \frac{\phi \xi'}{r} \right) + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} \left( \frac{\phi \xi'}{r} \right) - \frac{1}{r^2} \left( \frac{\phi \xi'}{r} \right) - \frac{2}{r^2} \frac{\partial}{\partial \theta} (-\phi' \xi) \quad (499)$$

$$= \left( \frac{\phi''}{r} - 2 \frac{\phi'}{r^2} + 2 \frac{\phi}{r^3} \right) \xi' + \left( \frac{\phi'}{r^2} - \frac{\phi}{r^3} \right) \xi' + \frac{\phi}{r^3} \xi''' - \frac{\phi \xi'}{r^3} + \frac{2}{r^2} \phi' \xi' \quad (500)$$

$$= \frac{\phi'' \xi'}{r} + \frac{\phi' \xi'}{r^2} + \frac{\phi \xi'''}{r^3} \quad (501)$$

$$\frac{\partial A_r}{\partial \theta} = \frac{\phi'' \xi''}{r} + \frac{\phi' \xi''}{r^2} + \frac{\phi \xi''''}{r^3} \quad (502)$$

(503)

$$A_\theta = \frac{\partial^2 v_\theta}{\partial r^2} + \frac{1}{r} \frac{\partial v_\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 v_\theta}{\partial \theta^2} + \frac{2}{r^2} \frac{\partial v_r}{\partial \theta} - \frac{v_\theta}{r^2} \quad (504)$$

$$= \frac{\partial^2}{\partial r^2} (-\phi' \xi) + \frac{1}{r} \frac{\partial}{\partial r} (-\phi' \xi) + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2} (-\phi' \xi) + \frac{2}{r^2} \frac{\partial}{\partial \theta} \left( \frac{\phi \xi'}{r} \right) - \frac{1}{r^2} (-\phi' \xi) \quad (505)$$

$$= -\phi''' \xi - \frac{\phi'' \xi}{r} - \frac{\phi' \xi''}{r^2} + \frac{2\phi \xi''}{r^2} + \frac{\phi' \xi}{r^2} \quad (506)$$

$$\frac{\partial(rA_\theta)}{\partial r} = \quad (507)$$

WRONG:

$$\frac{\partial(r\Delta v)}{\partial r} = \frac{\partial}{\partial r} \left( \frac{\partial}{\partial r} \left( r \frac{\partial v}{\partial r} \right) + \frac{1}{r} \frac{\partial^2 v}{\partial \theta^2} \right) \quad (508)$$

$$= \frac{\partial^2}{\partial r^2} \left( r \frac{\partial v}{\partial r} \right) + \frac{\partial}{\partial r} \left( \frac{1}{r} \frac{\partial^2 v}{\partial \theta^2} \right) \quad (509)$$

$$= \frac{\partial^2}{\partial r^2} \left( r \frac{\partial}{\partial r} \left( -\frac{\partial \Psi}{\partial r} \right) \right) + \frac{\partial}{\partial r} \left( \frac{1}{r} \frac{\partial^2}{\partial \theta^2} \left( -\frac{\partial \Psi}{\partial r} \right) \right) \quad (510)$$

$$= -\frac{\partial^2}{\partial r^2} \left( r \frac{\partial^2 \Psi}{\partial r^2} \right) - \frac{\partial}{\partial r} \left( \frac{1}{r} \frac{\partial^3 \Psi}{\partial \theta^2 \partial r} \right) \quad (511)$$

$$= -2 \frac{\partial^3 \Psi}{\partial r^3} - r \frac{\partial^4 \Psi}{\partial r^4} + \frac{1}{r^2} \frac{\partial^3 \Psi}{\partial \theta^2 \partial r} - \frac{1}{r} \frac{\partial^4 \Psi}{\partial \theta^2 \partial r^2} \quad (512)$$

$$\frac{\partial \Delta u}{\partial \theta} = \frac{\partial}{\partial \theta} \left( \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial u}{\partial r} \right) + \frac{1}{r^2} \frac{\partial^2 u}{\partial \theta^2} \right) \quad (513)$$

$$= \frac{\partial}{\partial \theta} \left( \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial u}{\partial r} \right) \right) + \frac{1}{r^2} \frac{\partial^3 u}{\partial \theta^3} \quad (514)$$

$$= \frac{\partial}{\partial \theta} \left( \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial}{\partial r} \left( \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \right) \right) \right) + \frac{1}{r^2} \frac{\partial^3}{\partial \theta^3} \left( \frac{1}{r} \frac{\partial \Psi}{\partial \theta} \right) \quad (515)$$

$$= \frac{1}{r^3} \frac{\partial^2 \Psi}{\partial \theta^2} - \frac{1}{r^2} \frac{\partial^3 \Psi}{\partial r \partial \theta^2} + \frac{1}{r} \frac{\partial^4 \Psi}{\partial r^2 \partial \theta^2} + \frac{1}{r^3} \frac{\partial^4 \Psi}{\partial \theta^4} \quad (516)$$

Assuming the following separation of variables  $\boxed{\Psi(r, \theta) = \phi(r)\xi(\theta)}$ :

$$\frac{\partial(r\Delta v)}{\partial r} = -2\phi'''\xi - r\phi''''\xi + \frac{1}{r^2}\phi'\xi'' - \frac{1}{r}\phi''\xi'' \quad (517)$$

$$\frac{\partial \Delta u}{\partial \theta} = \frac{1}{r^3}\phi\xi'' - \frac{1}{r^2}\phi'\xi'' + \frac{1}{r}\phi''\xi'' + \frac{1}{r^3}\phi\xi''' \quad (518)$$

so that

$$\frac{\partial(r\Delta v)}{\partial r} - \frac{\partial \Delta u}{\partial \theta} = -2\phi'''\xi - r\phi''''\xi + \frac{1}{r^2}\phi'\xi'' - \frac{1}{r}\phi''\xi'' - \frac{1}{r^3}\phi\xi'' + \frac{1}{r^2}\phi'\xi'' - \frac{1}{r}\phi''\xi'' - \frac{1}{r^3}\phi\xi'''$$

Further assuming  $\boxed{\xi(\theta) = \cos(k\theta)}$ , then  $\xi'' = -k^2\xi$  and  $\xi''' = k^4\xi$  then

$$\frac{\partial(r\Delta v)}{\partial r} - \frac{\partial \Delta u}{\partial \theta} = -2\phi'''\xi - r\phi''''\xi - k^2 \frac{1}{r^2} \phi' \xi + k^2 \frac{1}{r} \phi'' \xi + k^2 \frac{1}{r^3} \phi \xi - k^2 \frac{1}{r^2} \phi' \xi + k^2 \frac{1}{r} \phi'' \xi - k^4 \frac{1}{r^3} \phi \xi$$

By choosing  $\rho$  such that  $\rho = \lambda(r)\Upsilon(\theta)$  and such that  $\partial_\theta \Upsilon = \xi = \cos(k\theta)$  then we have

$$-2\phi'''\xi - r\phi''''\xi - k^2 \frac{1}{r^2} \phi' \xi + k^2 \frac{1}{r} \phi'' \xi + k^2 \frac{1}{r^3} \phi \xi - k^2 \frac{1}{r^2} \phi' \xi + k^2 \frac{1}{r} \phi'' \xi - k^4 \frac{1}{r^3} \phi \xi = -\frac{1}{\eta} \lambda \xi g_r$$

and then dividing by  $\xi$ : (IS THIS OK ?)

$$-2\phi'''' - r\phi''''' - k^2 \frac{1}{r^2} \phi' + k^2 \frac{1}{r} \phi'' + k^2 \frac{1}{r^3} \phi - k^2 \frac{1}{r^2} \phi' + k^2 \frac{1}{r} \phi'' - k^4 \frac{1}{r^3} \phi = -\frac{1}{\eta} \lambda g_r$$

$$-2\phi'''' - r\phi''''' - 2k^2 \frac{1}{r^2} \phi' + 2k^2 \frac{1}{r} \phi'' + (k^2 - k^4) \frac{1}{r^3} \phi = -\frac{1}{\eta} \lambda g_r$$

so

$$\lambda(r) = \frac{\eta}{g_r} \left( 2\phi'''' + r\phi''''' + 2k^2 \frac{1}{r^2} \phi' - 2k^2 \frac{1}{r} \phi'' - (k^2 - k^4) \frac{1}{r^3} \phi \right)$$

Also not forget  $\Upsilon = \frac{1}{k} \sin(k\theta)$

## 42.1 Linking with our paper

We have

$$\phi(r) = -rg(r) \quad (519)$$

$$\phi'(r) = -g(r) - rg'(r) = -f(r) \quad (520)$$

$$\phi''(r) = -f'(r) \quad (521)$$

$$\phi'''(r) = -f''(r) \quad (522)$$

$$\phi''''(r) = -f'''(r) \quad (523)$$

$$f(r) = \frac{\eta_0}{g_0} \left( 2f''(r) + rf'''(r) + 2k^2 \frac{1}{r^2} f(r) - 2k^2 \frac{1}{r} f'(r) + (k^2 - k^4) \frac{1}{r^2} g(r) \right)$$

## 42.2 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.

$$\mathbf{v}(r = R_1) = \mathbf{v}(r = R_2) = \mathbf{0}$$

Due to the separation of variables, and since  $\xi(\theta) = \cos(k\theta)$  we have

$$u(r, \theta) = \frac{1}{r} \frac{\partial \Psi}{\partial \theta} = \frac{1}{r} \phi \xi' = -\frac{1}{r} \phi(r) k \sin(k\theta) \quad v(r, \theta) = -\frac{\partial \Psi}{\partial r} = -\phi'(r) \xi = -\phi'(r) \cos(k\theta)$$

It is obvious that the only way to insure no-slip boundary conditions is to have

$$\phi(R_1) = \phi(R_2) = \phi'(R_1) = \phi'(R_2) = 0$$

We could then choose

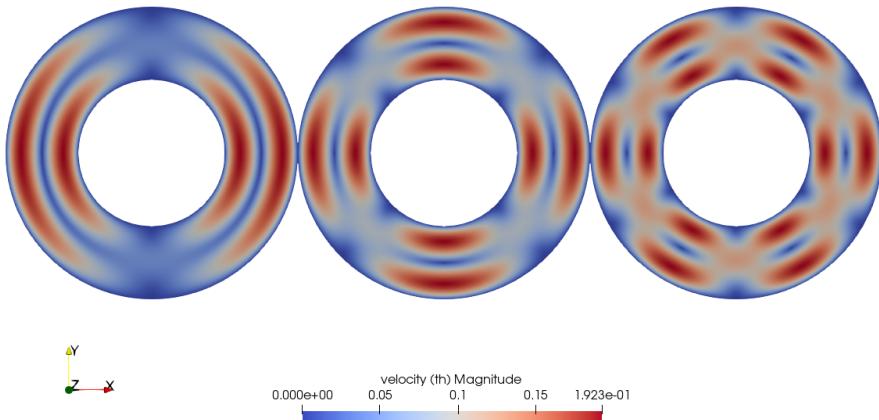
$$\phi(r) = (r - R_1)^2 (r - R_2)^2 \mathcal{F}(r) \quad (524)$$

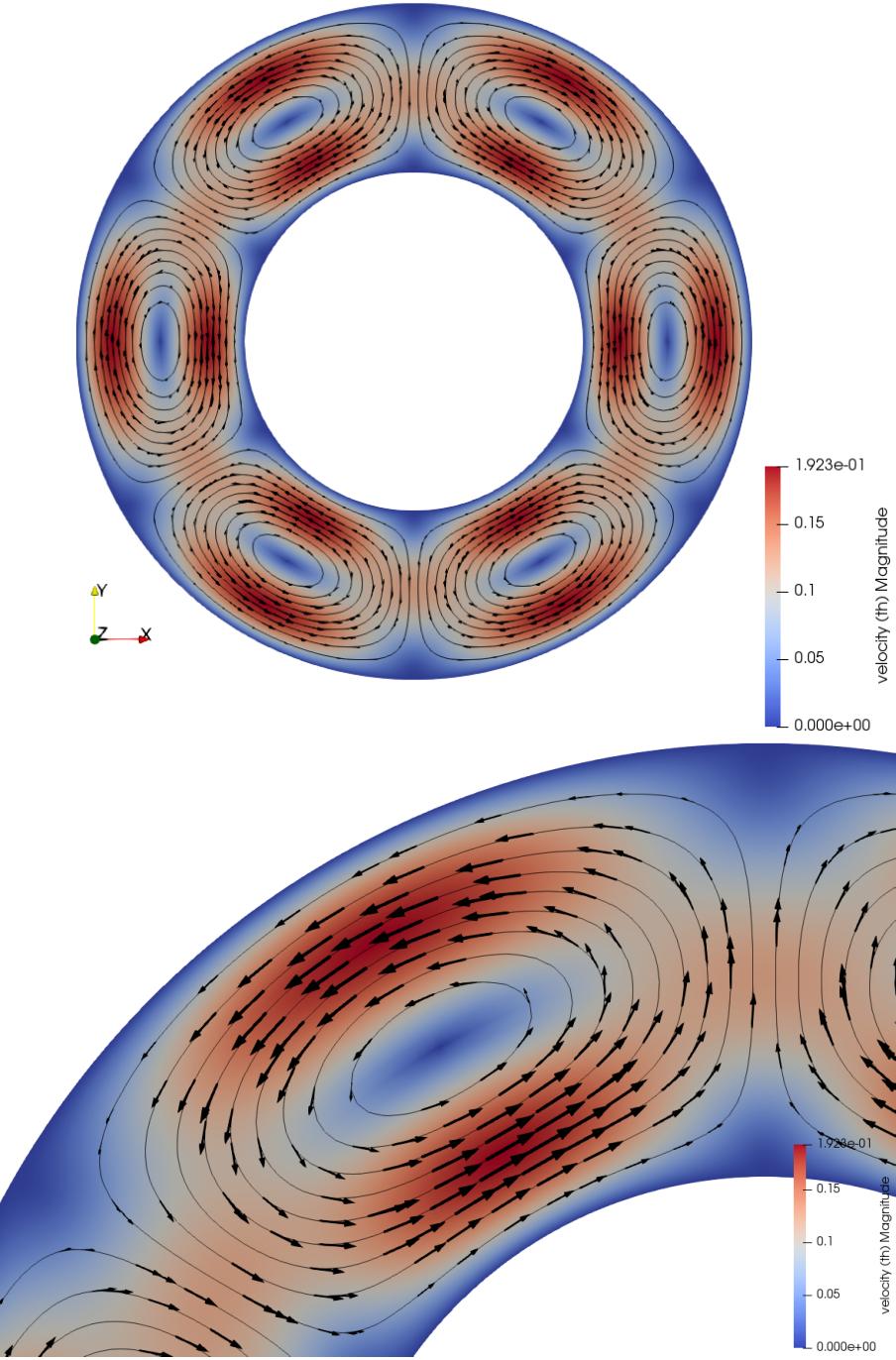
$$\phi'(r) = 2(r - R_1)(r - R_2)^2 \mathcal{F}(r) + 2(r - R_1)^2(r - R_2) \mathcal{F}(r) + (r - R_1)^2(r - R_2)^2 \mathcal{F}'(r) \quad (525)$$

which are indeed identically zero on both boundaries. Here  $\mathcal{F}(r)$  is any (smooth enough) function of  $r$ . We would then have

$$\boxed{\Psi(r, \theta) = (r - R_1)^2 (r - R_2)^2 \mathcal{F}(r) \cos(k\theta)}$$

In what follows we will take  $\mathcal{F}(r) = 1$  for simplicity.





COMPUTE  $f$  from  $\phi$  and then the pressure.

### 42.3 Free slip boundary conditions

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  $\Psi(r, \theta) = \phi(r)\xi(\theta)$ :

$$\phi(r = R_1) = \phi(r = R_2) = 0$$

- $(\sigma \cdot n) \times n = \mathbf{0}$  (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!):

$$-\frac{\partial^2 \Psi}{\partial r^2} + \frac{1}{r} \frac{\partial \Psi}{\partial r} + \frac{1}{r^2} \frac{\partial^2 \Psi}{\partial \theta^2} = 0$$

$$-\phi''\xi + \frac{1}{r}\phi'\xi + \frac{1}{r^2}\phi\xi'' = 0$$

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>35</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>35</sup><http://mathworld.wolfram.com/EulerDifferentialEquation.html>

## 43 fieldstone\_36: the annulus geometry elastic aquarium

This fieldstone was developed in collaboration with Lukas van de Wiel.

The domain is an annulus with inner radius  $R_1$  and outer radius  $R_2$ . It is filled with a single elastic material characterised by a Young's modulus  $E$  and a Poisson ratio  $\nu$ , a density  $\rho_0$ . The gravity  $\mathbf{g} = -g_0 \mathbf{e}_r$  is pointing towards the center of the domain.

The problem at hand is axisymmetric so that the tangential component of the displacement vector  $v_\theta$  is assumed to be zero as well as all terms containing  $\partial_\theta$ . The components of the strain tensor are

$$\varepsilon_{rr} = \frac{\partial v_r}{\partial r} \quad (526)$$

$$\varepsilon_{\theta\theta} = \frac{v_r}{r} + \frac{1}{r} \frac{\partial v_\theta}{\partial \theta} = \frac{v_r}{r} \quad (527)$$

$$\varepsilon_{r\theta} = \frac{1}{2} \left( \frac{\partial v_\theta}{\partial r} - \frac{v_\theta}{r} + \frac{1}{r} \frac{\partial v_r}{\partial \theta} \right) = 0 \quad (528)$$

so that the tensor simply is

$$\boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_{rr} & \varepsilon_{r\theta} \\ \varepsilon_{r\theta} & \varepsilon_{\theta\theta} \end{pmatrix} = \begin{pmatrix} \frac{\partial v_r}{\partial r} & 0 \\ 0 & \frac{v_r}{r} \end{pmatrix} \quad (529)$$

The pressure is

$$p = -\lambda \nabla \cdot \mathbf{v} = -\lambda \left( \frac{1}{r} \frac{\partial(rv_r)}{\partial r} \right) \quad (530)$$

and finally the stress tensor:

$$\boldsymbol{\sigma} = -p \mathbf{1} + 2\mu \boldsymbol{\varepsilon} = \begin{pmatrix} \lambda \frac{1}{r} \frac{\partial(rv_r)}{\partial r} + 2\mu \frac{\partial v_r}{\partial r} & 0 \\ 0 & \lambda \frac{1}{r} \frac{\partial(rv_r)}{\partial r} + 2\mu \frac{v_r}{r} \end{pmatrix} \quad (531)$$

The divergence of the stress tensor is given by [500]:

$$\nabla \cdot \boldsymbol{\sigma} = \begin{pmatrix} \frac{\partial \sigma_{rr}}{\partial r} + \frac{\sigma_{rr} - \sigma_{\theta\theta}}{r} + \frac{1}{r} \frac{\partial \sigma_{\theta r}}{\partial \theta} \\ \frac{\partial \sigma_{r\theta}}{\partial r} + \frac{1}{r} \frac{\sigma_{\theta\theta}}{\partial \theta} + \frac{\sigma_{r\theta} + \sigma_{\theta r}}{r} \end{pmatrix} \quad (532)$$

Given the diagonal nature of the stress tensor this simplifies to (also remember that  $\partial_\theta = 0$ ):

$$\nabla \cdot \boldsymbol{\sigma} = \begin{pmatrix} \frac{\partial \sigma_{rr}}{\partial r} + \frac{\sigma_{rr} - \sigma_{\theta\theta}}{r} \\ 0 \end{pmatrix} \quad (533)$$

Focusing on the  $r$ -component of the stress divergence:

$$(\nabla \cdot \boldsymbol{\sigma})_r = \frac{\partial \sigma_{rr}}{\partial r} + \frac{\sigma_{rr} - \sigma_{\theta\theta}}{r} \quad (534)$$

$$= \frac{\partial}{\partial r} \left[ \lambda \frac{1}{r} \frac{\partial(rv_r)}{\partial r} + 2\mu \frac{\partial v_r}{\partial r} \right] + \frac{1}{r} \left[ \lambda \frac{1}{r} \frac{\partial(rv_r)}{\partial r} + 2\mu \frac{\partial v_r}{\partial r} - \lambda \frac{1}{r} \frac{\partial(rv_r)}{\partial r} - 2\mu \frac{v_r}{r} \right] \quad (535)$$

$$= \lambda \frac{\partial}{\partial r} \frac{1}{r} \frac{\partial(rv_r)}{\partial r} + 2\mu \frac{\partial^2 v_r}{\partial r^2} + \lambda \frac{1}{r^2} \frac{\partial(rv_r)}{\partial r} + \frac{2\mu}{r} \frac{\partial v_r}{\partial r} - \lambda \frac{1}{r^2} \frac{\partial(rv_r)}{\partial r} - \frac{2\mu v_r}{r^2} \quad (536)$$

$$= \lambda \left( -\frac{v_r}{r^2} + \frac{1}{r} \frac{\partial v_r}{\partial r} + \frac{\partial^2 v_r}{\partial r^2} \right) + 2\mu \frac{\partial^2 v_r}{\partial r^2} + \frac{2\mu}{r} \frac{\partial v_r}{\partial r} - \frac{2\mu v_r}{r^2} \quad (537)$$

$$= -(2\mu + \lambda) \frac{v_r}{r^2} + (2\mu + \lambda) \frac{1}{r} \frac{\partial v_r}{\partial r} + (2\mu + \lambda) \frac{\partial^2 v_r}{\partial r^2} \quad (538)$$

So the momentum conservation in the  $r$  direction is

$$(\nabla \cdot \boldsymbol{\sigma} + \rho_0 \mathbf{g})_r = -(2\mu + \lambda) \frac{v_r}{r^2} + (2\mu + \lambda) \frac{1}{r} \frac{\partial v_r}{\partial r} + (2\mu + \lambda) \frac{\partial^2 v_r}{\partial r^2} - \rho_0 g_0 = 0 \quad (539)$$

or,

$$\frac{\partial^2 v_r}{\partial r^2} + \frac{1}{r} \frac{\partial v_r}{\partial r} - \frac{v_r}{r^2} = \frac{\rho_0 g_0}{\lambda + 2\mu} \quad (540)$$

We now look at the boundary conditions. On the inner boundary we prescribe  $v_r(r = R_1) = 0$  while free surface boundary conditions are prescribed on the outer boundary, i.e.  $\boldsymbol{\sigma} \cdot \mathbf{n} = 0$  (i.e. there is no force applied on the surface).

The general form of the solution can then be obtained:

$$v_r(r) = C_1 r^2 + C_2 r + \frac{C_3}{r} \quad (541)$$

with

$$C_1 = \frac{\rho_0 g_0}{3(\lambda + 2\mu)} \quad C_2 = -C_1 R_1 - \frac{C_3}{R_1^2} \quad C_3 = \frac{k_1 + k_2}{(R_1^2 + R_2^2)(2\mu + \lambda) + (R_2^2 - R_1^2)\lambda} \quad (542)$$

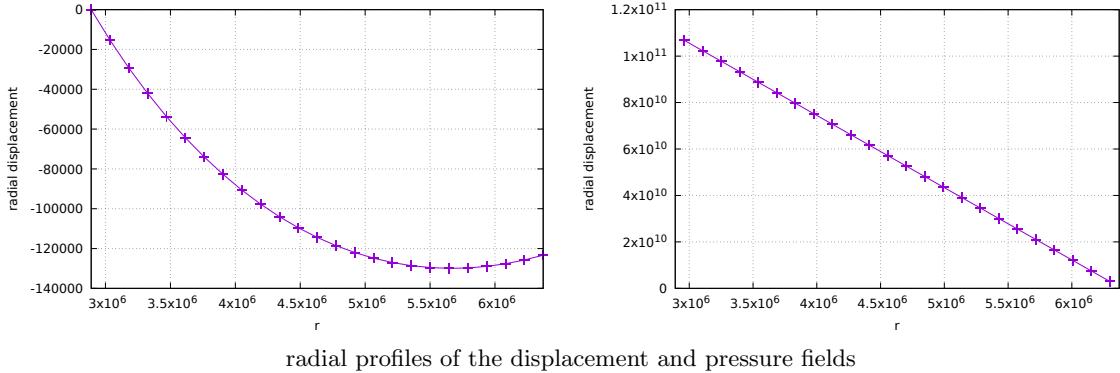
and

$$k_1 = (2\mu + \lambda)C_1(2R_1^2 R_2^3 - R_1^3 R_2^2) \quad k_2 = \lambda C_1(R_1^2 R_2^3 - R_1^3 R_2^2) \quad (543)$$

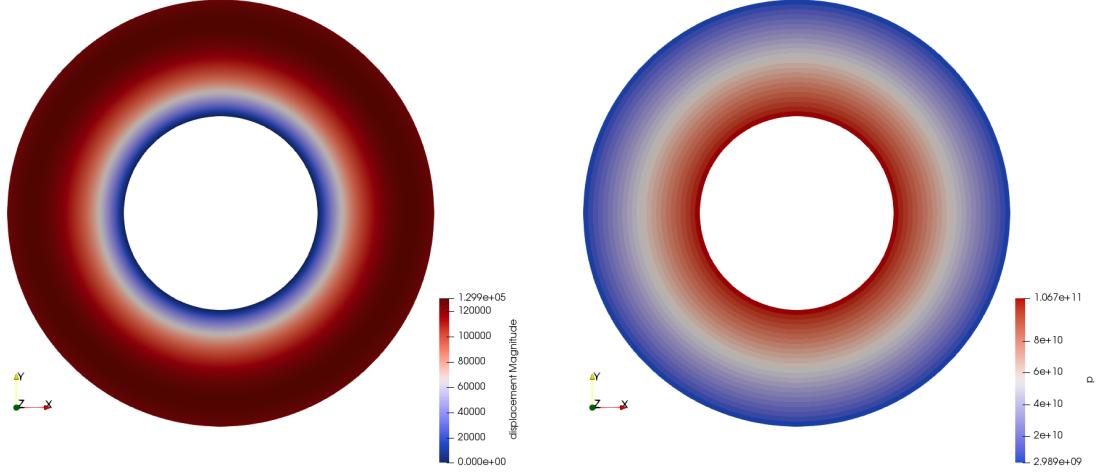
Pressure can then be computed as follows:

$$p = -\lambda \nabla \cdot \mathbf{v} = -\lambda \left( \frac{1}{r} \frac{\partial(r v_r)}{\partial r} \right) = -\lambda \left( \frac{1}{r} (3C_1 r^2 + 2C_2 r) \right) = -\lambda (3C_1 r + 2C_2) \quad (544)$$

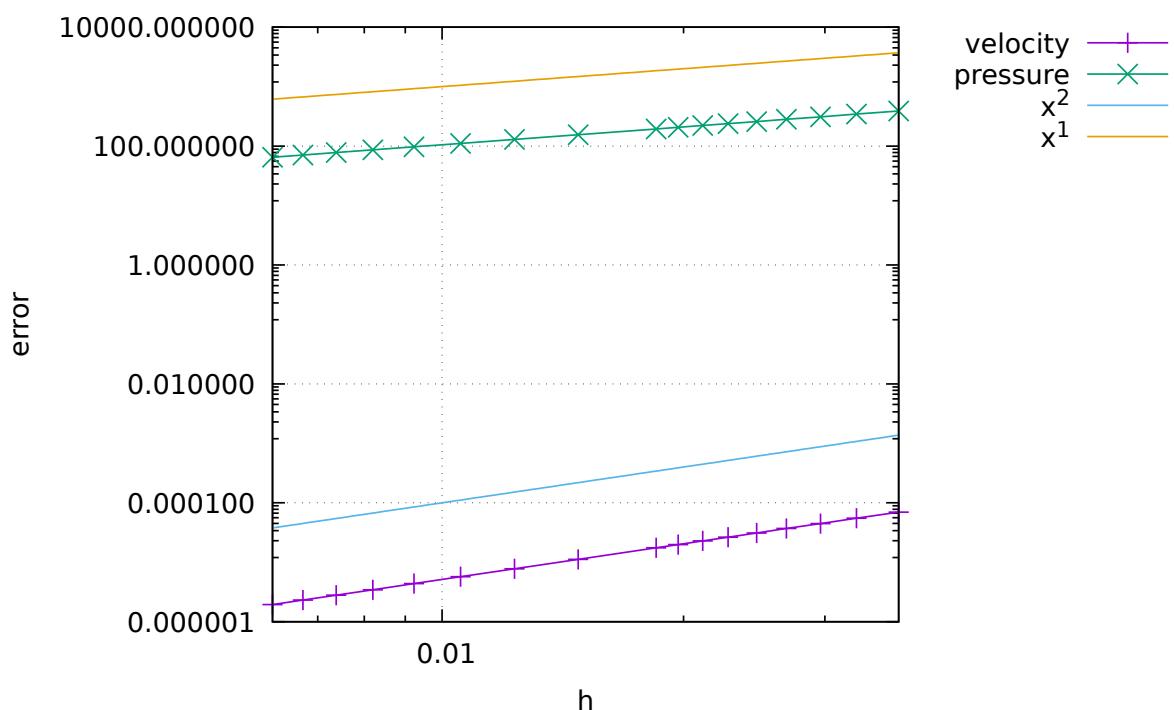
We choose  $R_1 = 2890\text{km}$ ,  $R_2 = 6371\text{km}$ ,  $g_0 = 9.81\text{ms}^{-2}$ ,  $\rho_0 = 3300$ ,  $E = 6 \cdot 10^{10}$ ,  $\nu = 0.49$ .



radial profiles of the displacement and pressure fields



displacement and pressure fields in the domain



## 44 fieldstone\_37: marker advection and population control



The domain is a unit square. The Stokes equations are not solved, the velocity is prescribed everywhere in the domain as follows:

$$u = -(z - 0.5) \quad (545)$$

$$v = 0 \quad (546)$$

$$w = (x - 0.5) \quad (547)$$

At the moment, velocity is computed on the marker itself (rk0 algorithm). When markers are advected outside, they are arbitrarily placed at location (-0.0123,-0.0123).

in construction.

## 45 fieldstone\_38: Critical Rayleigh number

This fieldstone was developed in collaboration with Arie van den Berg.

The system is a layer of fluid between  $y = 0$  and  $y = 1$ , with boundary conditions  $T(x, y = 0) = 1$  and  $T(x, y = 1) = 0$ , characterized by  $\rho$ ,  $c_p$ ,  $k$ ,  $\eta_0$ . The Rayleigh number of the system is

$$\text{Ra} = \frac{\rho_0 g_0 \alpha \Delta T h^3}{\eta_0 \kappa}$$

We have  $\Delta T = 1$ ,  $h = 1$  and choose  $\kappa = 1$  so that the Rayleigh number simplifies to  $\text{Ra} = \rho_0 g_0 \alpha / \eta_0$ .

The Stokes equation is  $\vec{\nabla} \cdot \boldsymbol{\sigma} + \vec{b} = \vec{0}$  with  $\vec{b} = \rho \vec{g}$ . Then the components of this equation on the  $x$ - and  $y$ -axis are:

$$(\vec{\nabla} \cdot \boldsymbol{\sigma})_x = -\rho \vec{g} \cdot \vec{e}_x = 0 \quad (548)$$

$$(\vec{\nabla} \cdot \boldsymbol{\sigma})_y = -\rho \vec{g} \cdot \vec{e}_y = \rho g_0 \quad (549)$$

since  $\vec{g}$  and  $\vec{e}_y$  are in opposite directions ( $\vec{g} = -g_0 \vec{e}_y$ , with  $g_0 > 0$ ). The stream function formulation of the incompressible isoviscous Stokes equation is then

$$\nabla^4 \Psi = \frac{g_0}{\eta_0} \frac{\partial \rho}{\partial x}$$

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

$$\boxed{\nabla^4 \Psi = -\frac{\rho_0 g_0 \alpha}{\eta_0} g \frac{\partial T}{\partial x} = -\text{Ra} \frac{\partial T}{\partial x}} \quad (550)$$

For small perturbations of the conductive state  $T_0(y) = 1 - y$  we define the temperature perturbation  $T_1(x, y)$  such that

$$T(x, y) = T_0(y) + T_1(x, y)$$

The temperature perturbation  $T_1$  satisfies the homogeneous boundary conditions  $T_1(x, y = 0) = 0$  and  $T_1(x, y = 1) = 0$ . The temperature equation is

$$\rho c_p \frac{DT}{Dt} = \rho c_p \left( \frac{\partial T}{\partial t} + \vec{v} \cdot \vec{\nabla} T \right) = \rho c_p \left( \frac{\partial T_0 + T_1}{\partial t} + \vec{v} \cdot \vec{\nabla}(T_0 + T_1) \right) = k \Delta(T_0 + T_1)$$

and can be simplified as follows:

$$\rho c_p \left( \frac{\partial T_1}{\partial t} + \vec{v} \cdot \vec{\nabla} T_0 \right) = k \Delta T_1$$

since  $T_0$  does not depend on time,  $\Delta T_0 = 0$  and we assume the nonlinear term  $\vec{v} \cdot \vec{\nabla} T_1$  to be second order (temperature perturbations and coupled velocity changes are assumed to be small). Using the relationship between velocity and stream function  $v_y = -\partial_x \Psi$  we have  $\vec{v} \cdot \vec{\nabla} T_0 = -v_y = \partial_x \Psi$  and since  $\kappa = k/\rho c_p = 1$  we get

$$\boxed{\frac{\partial T_1}{\partial t} - \kappa \Delta T_1 = -\frac{\partial \Psi}{\partial x}} \quad (551)$$

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  $T_1$  can be written:

$$\Psi(x, y, t) = A_\Psi \exp(pt) \exp(\pm i k_x x) \exp(\pm i k_y y) = A_\Psi E_\psi(x, y, t) \quad (552)$$

$$T_1(x, y, t) = A_T \exp(pt) \exp(\pm i k_x x) \exp(\pm i k_y y) = A_T E_T(x, y, t) \quad (553)$$

where  $k_x$  and  $k_y$  are the horizontal and vertical wave number respectively. Note that we then have

$$\nabla^2 \Psi = -(k_x^2 + k_y^2) \Psi \quad \nabla^2 T_1 = -(k_x^2 + k_y^2) T_1$$

The boundary conditions on  $T_1$ , coupled with a choice of a real function for the  $x$  dependence yields:

$$E_T(x, y, t) = \exp(pt) \cos(k_x x) \sin(n\pi y).$$

**from here onwards check for minus signs!**

The velocity boundary conditions are  $v_y(x, y=0) = 0$  and  $v_y(x, y=1) = 0$  which imposes conditions on  $\partial\Psi/\partial x$  and we find that we can use the same  $y$  dependence as for  $T_1$ . Choosing again for a real function for the  $x$  dependence yields:

$$E_\Psi(x, y, t) = \exp(pt) \sin(k_x x) \sin(n\pi z)$$

We then have

$$\Psi(x, y, t) = A_\Psi \exp(pt) \sin(k_x x) \sin(n\pi z) = A_\Psi E_\Psi(x, y, t) \quad (554)$$

$$T_1(x, y, t) = A_T \exp(pt) \cos(k_x x) \sin(n\pi z) = A_T E_T(x, y, t) \quad (555)$$

In what follows we simplify notations:  $k = k_x$ . Then the two PDEs become:

$$pT_1 + \kappa(k^2 + n^2\pi^2) - kA_\Psi \exp(pt) \cos(k_x x) \sin(n\pi z) = kA_\Psi E_\theta \quad (556)$$

$$-RaA_T \cos(kx) \sin(n\pi z) + \kappa(k^2 + n^2\pi^2)^2 A_\Psi = -RaA_T E_\Psi + \kappa(k^2 + n^2\pi^2)^2 A_\Psi = 0 \quad (557)$$

These equations must then be verified for all ... which leads to write:

$$\begin{pmatrix} p + (k^2 + n^2\pi^2) & -k \\ -Ra k & (k^2 + n^2\pi^2)^2 \end{pmatrix} \begin{pmatrix} A_\theta \\ A_\Psi \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

The determinant of such system must be nul otherwise there is only a trivial solution to the problem, i.e.  $A_\theta = 0$  and  $A_\Psi = 0$  which is not helpful. CHECK/REPHRASE

$$D = [p + (k^2 + n^2\pi^2)](k^2 + n^2\pi^2)^2 - Ra k^2 = 0$$

or,

$$p = \frac{Ra k^2 - (k^2 + n^2\pi^2)^3}{(k^2 + n^2\pi^2)^2}$$

The coefficient  $p$  determines the stability of the system: if it is negative, the system is stable and both  $\Psi$  and  $T_1$  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. The threshold is then  $p = 0$  and the solution of the above system is

## 46 fieldstone\_39: chpe15

The Drucker-Prager yield function is given by the function  $f$ :

$$f = p \sin \phi + c \cos \phi - \tau$$

where  $\tau$  is the square root of the second invariant of the deviatoric stress. We have

$$p = \frac{1}{2}(\sigma_1 + \sigma_3)$$

and

$$\tau = \frac{1}{2}(\sigma_1 - \sigma_3)$$

Inserting these into  $f$  yields:

$$f = \frac{1}{2}(\sigma_1 + \sigma_3) \sin \phi + c \cos \phi - \frac{1}{2}(\sigma_1 - \sigma_3)$$

The yield condition  $f = 0$  can be reworked as follows:

$$\sigma_1 - \frac{1 + \sin \phi}{1 - \sin \phi} \sigma_3 - 2 \frac{\cos \phi}{1 - \sin \phi} c = 0$$

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  $N_\phi$  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 [130].

This paper offers a solution to the problem of the angle of shear bands in geodynamic models. The underlying idea is based on simple modifications brought to existing incompressible flow codes. Note that the codes featured in that paper also implemented elastic behaviour but this can be easily switched off by setting  $Z = 1$  in their equations.

Their plasticity implementation starts with a modification of the continuity equation:

$$\vec{\nabla} \cdot \vec{v} = R = 2 \sin \psi \dot{\varepsilon}_p$$

where  $R$  is the dilation rate,  $\Psi$  is the dilation angle and  $\dot{\varepsilon}_p$  is the square root of the second invariant of the plastic strain rate.

Under this assumption, the deviatoric strain rate tensor is given by

$$\dot{\varepsilon}^d(\vec{v}) = \dot{\varepsilon}(\vec{v}) - \frac{1}{3} \text{Tr}[\dot{\varepsilon}(\vec{v})] \mathbf{1} = \dot{\varepsilon}(\vec{v}) - \frac{1}{3} \vec{\nabla} \cdot \vec{v} \mathbf{1} = \dot{\varepsilon}(\vec{v}) - \frac{1}{3} R \mathbf{1}$$

Turning now to the momentum conservation equation:

$$\begin{aligned} -\vec{\nabla} p + \vec{\nabla} \cdot \boldsymbol{\tau} &= -\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\varepsilon}^d(\vec{v})) \\ &= -\vec{\nabla} p + \vec{\nabla} \cdot \left[ 2\eta \left( \dot{\varepsilon}(\vec{v}) - \frac{1}{3} R \mathbf{1} \right) \right] \\ &= -\vec{\nabla} p + \vec{\nabla} \cdot (2\eta \dot{\varepsilon}(\vec{v})) - \frac{2}{3} \vec{\nabla} (\eta R) \end{aligned} \tag{558}$$

The last term is then an addition to the right hand side of the momentum equation and its weak form is as follows:

$$\vec{f}' = \int_{\Omega} N_v \frac{2}{3} \vec{\nabla}(\eta R) dV = \frac{4}{3} \sin \Psi \int_{\Omega} N_v \vec{\nabla}(\eta \dot{\varepsilon}_p) dV$$

This formulation proves to be problematic since in order to compute the gradient, we would need the viscosity and the plastic strain rate on the mesh nodes and both these quantities are effectively computed on the quadrature points. One option could be to project those quadrature values onto the nodes, which may introduce interpolation errors/artefacts and/or smoothing. Another option is to resort to integration by parts:

$$\int_{\Omega} N_v \vec{\nabla}(\eta \dot{\varepsilon}_p) dV = [N_v \eta \dot{\varepsilon}_p]_{\Gamma} - \int_{\Omega} \vec{\nabla} N_v (\eta \dot{\varepsilon}_p) dV$$

The last term is now trivial to compute since the shape function derivatives, the viscosity and the plastic strain rate are known at the quadrature points. Remains the surface term. We will neglect it for now to simplify our implementation and note that a) it will not directly affect what happens inside the domain, b) it could be somewhat important when shear bands intersect with the free surface.

$$\vec{f}' = -\frac{4}{3} \sin \psi \int_{\Omega} \vec{\nabla} N_v (\eta \dot{\varepsilon}_p) dV = -\frac{2}{3} \int_{\Omega} \vec{\nabla} N_v (\eta R) dV$$

Although the authors do indicate that they add a term in each rhs, it is not very clear how they deal with the implementation issue above. We then propose an alternative: instead of explicitly removing the deviatoric part of the strain rate as in Eq. 558 and replace the trace of the tensor by  $R$ , one could leave the term inside the matrix, thereby using a compressible form of the viscous block of the Stokes matrix. We will recover the same converged solution as before, but the path to convergence will be different than the first approach. In what follows, we denote the original approach by Choi & Petersen 'method 1' and the latter 'method 2'.

Finally, we need to define what the plastic strain rate tensor is. When using a rigid plastic rheology, the only deformation mechanism *is* plasticity so that the plastic strain rate *is* the strain rate. When using a visco-plastic rheology, the plastic strain rate is the strain rate of the zones above/at yield (the shear bands, where the vrm is active).

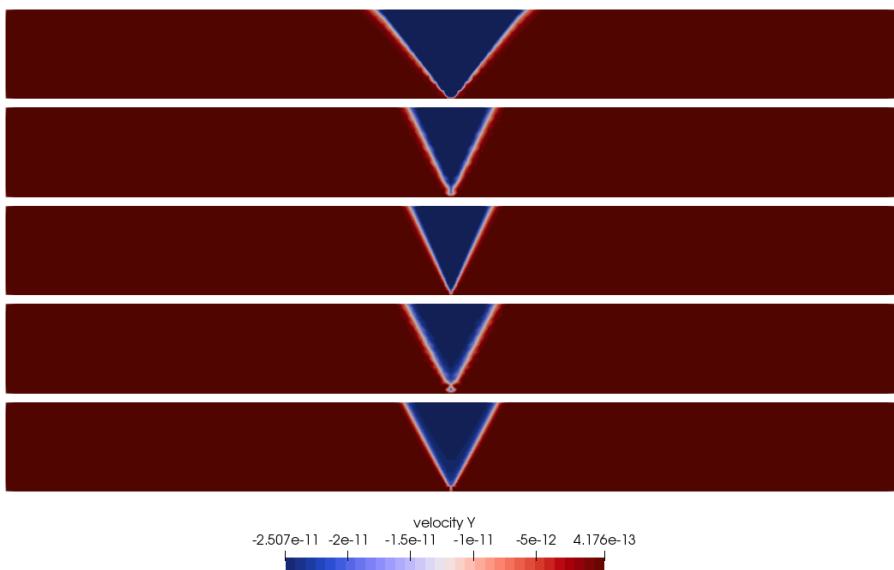
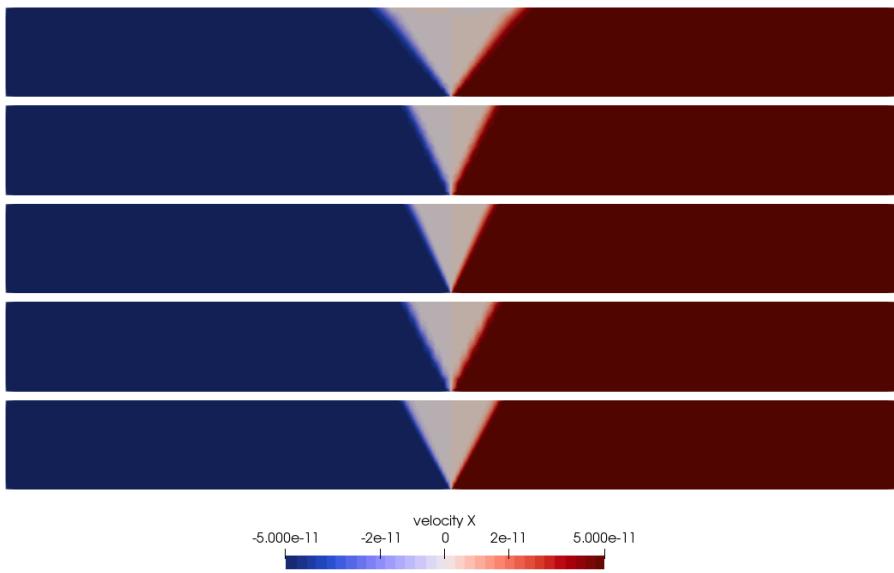
The setup is similar to the one in [337]. It is a 2D Cartesian domain filled with a single rigid-plastic material characterised by a cohesion  $c = 10\text{MPa}$ , an angle of friction  $\phi$ , a dilation angle  $\psi$  and a density  $\rho = 2800\text{kg/m}^3$ . Extensional boundary conditions are as follows:

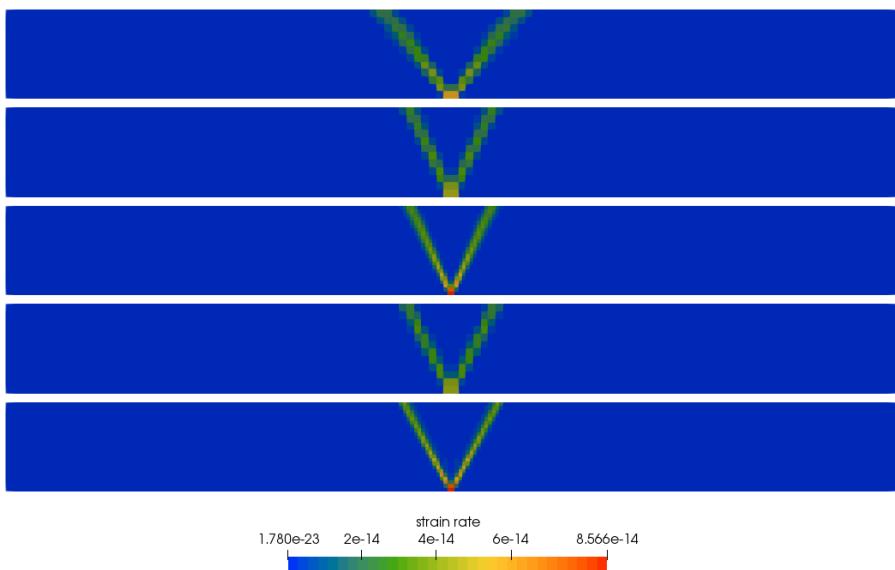
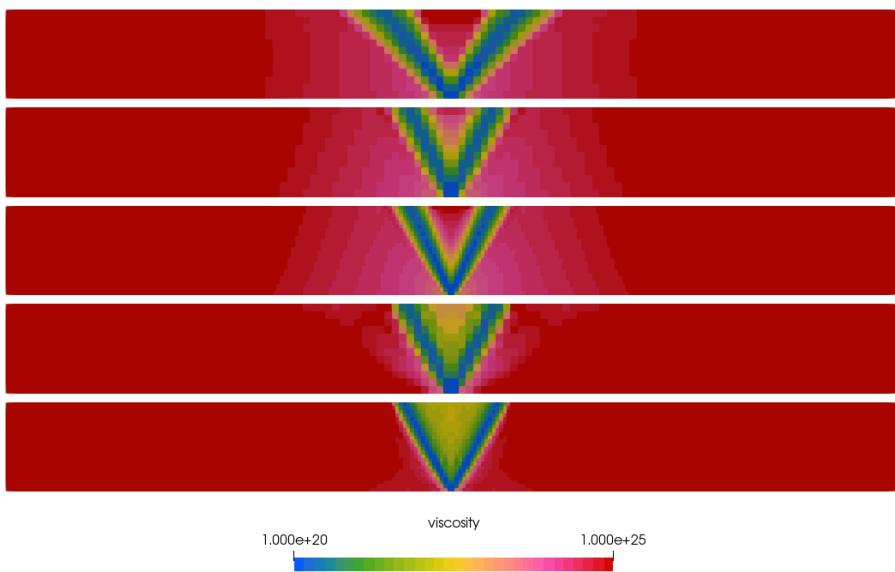
- left boundary:  $u = -v_{bc}$ ;
- right boundary:  $u = +v_{bc}$ ;
- bottom boundary:  $v = 0$ ,  $u = -v_{bc}$  for  $x < L_x/2$ ,  $u = +v_{bc}$  for  $x > L_x/2$ , and  $u = 0$  if  $x = L_x/2$ ;
- top boundary: zero traction.

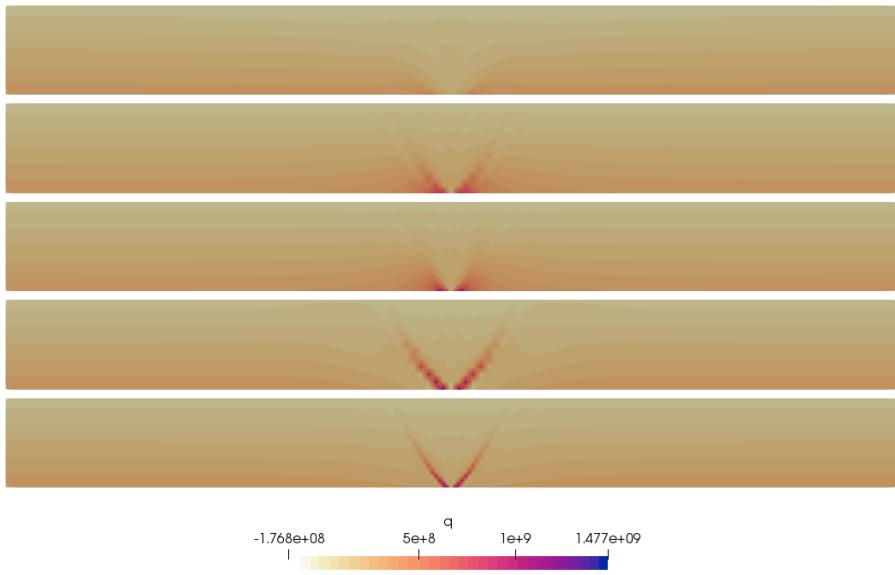
For compressional boundary conditions the signs of all horizontal velocities should be reversed. The nonlinear tolerance is set to  $\text{tol} = 10^{-6}$ . Nonlinear iterations stop when maximum of the normalised nonlinear residual reaches the desired tolerance.

Following Choi & Petersen [130], we run the experiment with an associative ( $\phi = \psi$ ) plasticity and a non associative one ( $\psi = 0$ , i.e.  $R = 0$ ). This second approach is essentially what many codes do.

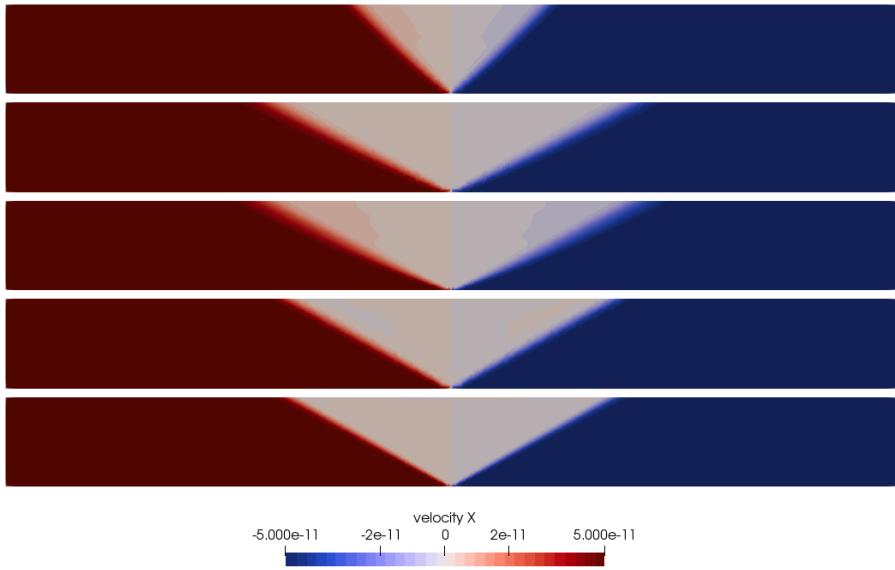
The velocity, pressure, strain rate, dilation rate, and velocity divergence are shown hereunder both in extension and compression.

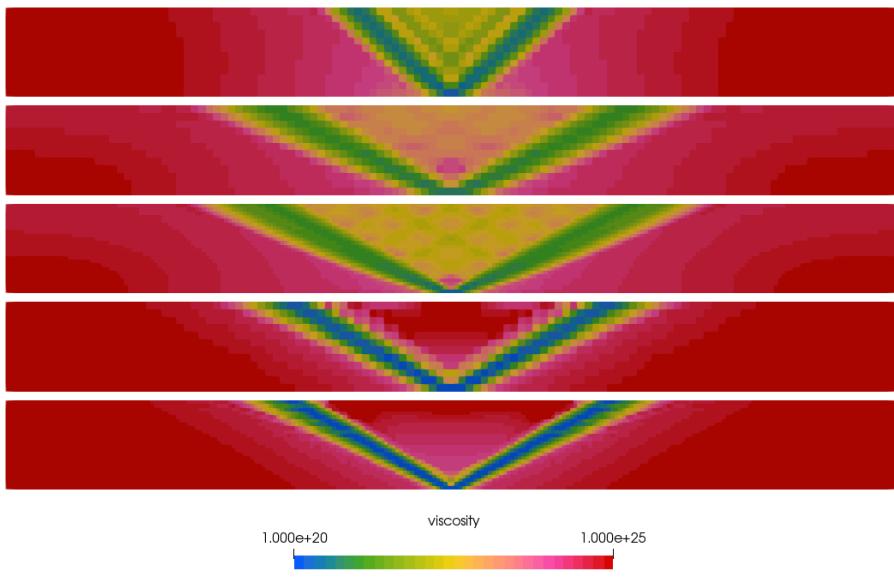
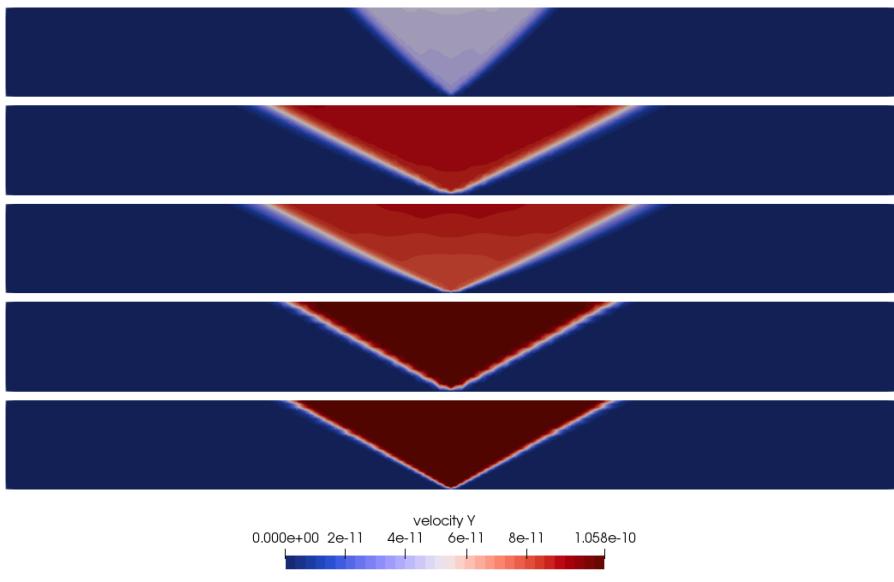


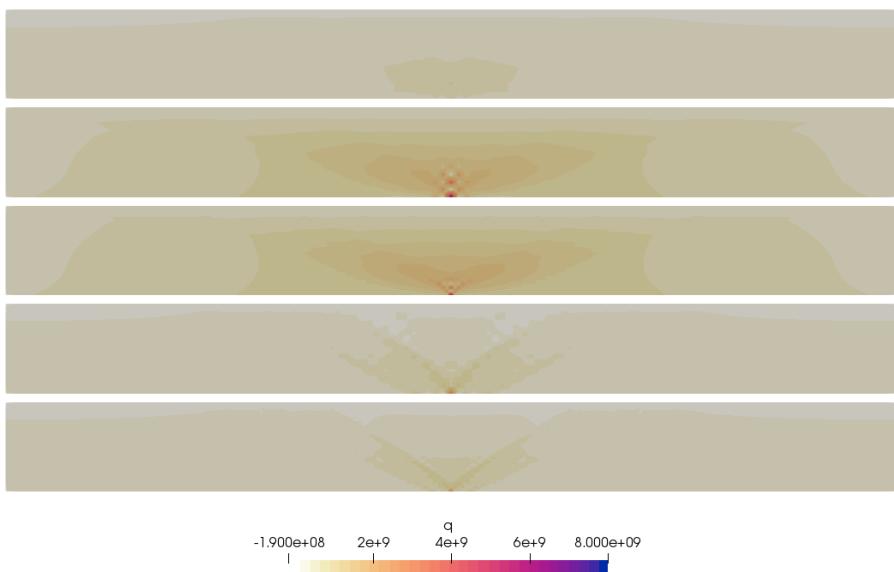
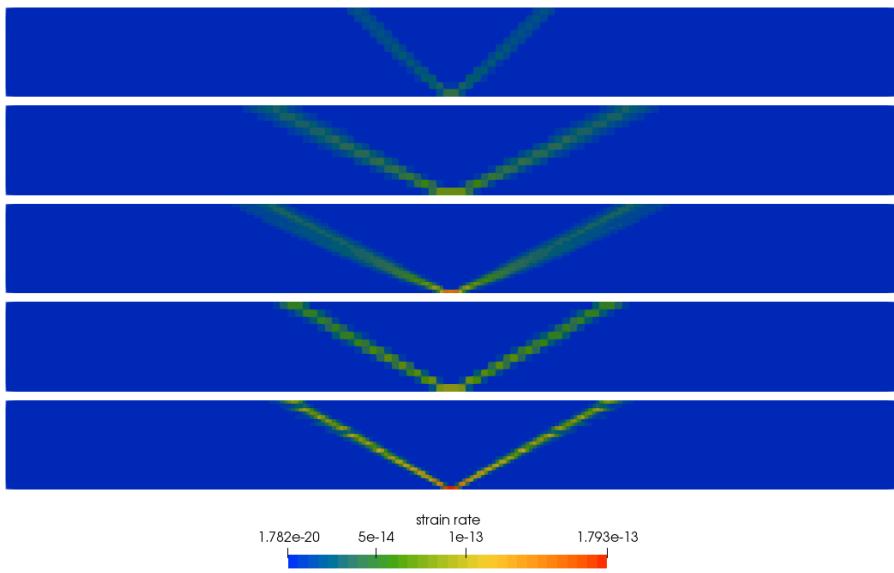




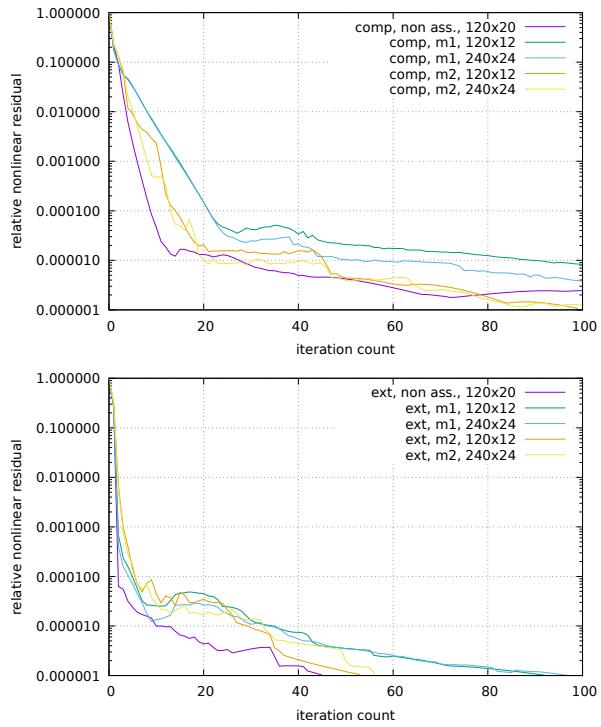
Extension. 1st row: Non-associative plasticity; 2nd and 3rd row: associative plasticity ( $\psi = \phi$ ) with method 1 for two resolutions 120x12 and 240; 4th and 5th row: associative plasticity ( $\psi = \phi$ ) with method 2 for two resolutions 120x12 and 240



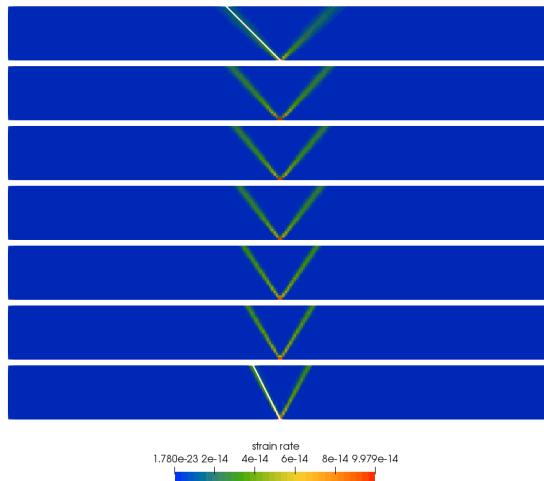


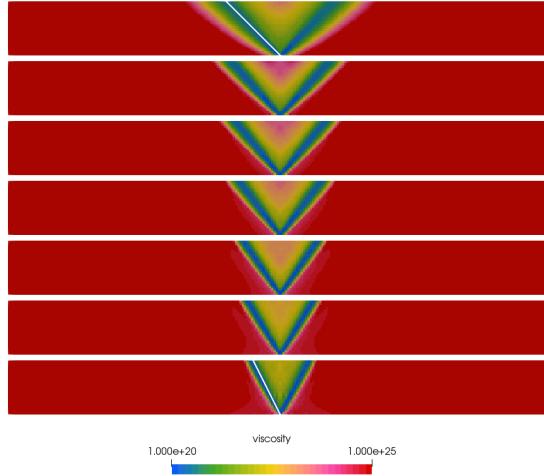


Compression. 1st row: Non-associative plasticity; 2nd and 3rd row: associative plasticity ( $\psi = \phi$ ) with method 1 for two resolutions 120x12 and 240; 4th and 5th row: associative plasticity ( $\psi = \phi$ ) with method 2 for two resolutions 120x12 and 240



One can also run the extension model for  $\phi = \psi = 0, 5, 10, 15, 20, 25, 30^\circ$





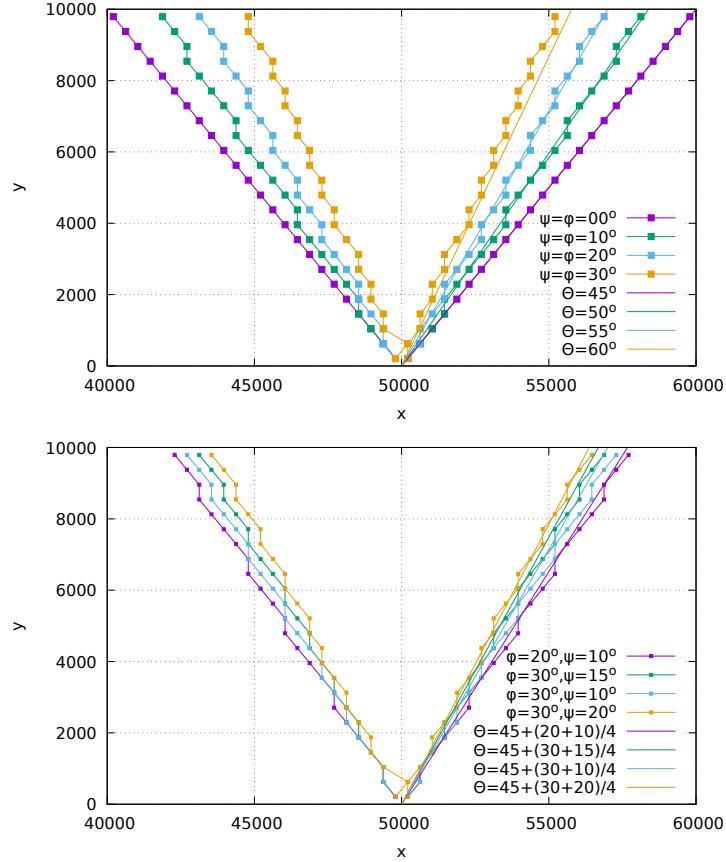
Three angles are mechanically stable (e.g. [337]):

$$\theta = \frac{\pi}{4} \pm \frac{\psi}{2} \quad \text{Roscoe angle}$$

$$\theta = \frac{\pi}{4} \pm \frac{\phi}{2} \quad \text{Coulomb angle}$$

$$\theta = \frac{\pi}{4} \pm \frac{\phi + \psi}{4} \quad \text{Arthur angle}$$

In the case of associative plasticity,  $\phi = \psi$ , so that all three angles are the same. Per row of elements, and per half of the domain (left and right) we find the element with the highest strain-rate and record their center coordinates on the figure hereunder. These elements are shown for  $\phi = \psi = \{0, 10, 20, 30\}^\circ$  alongside a line corresponding to the expected analytical shear band angle value.



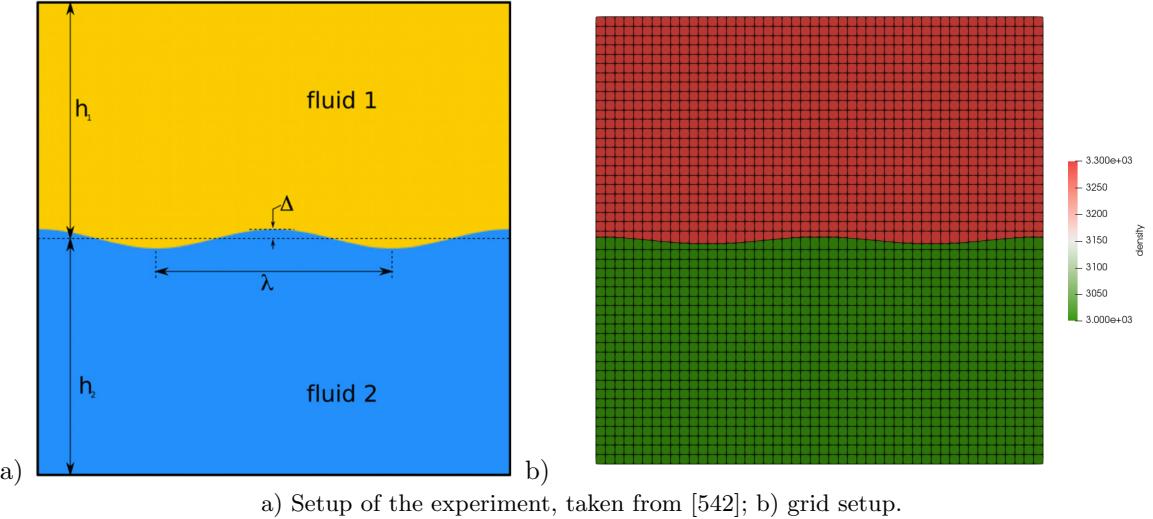
Results obtained on a 240x24 grid, max 50 nl iterations.

Note that benchmarking this is not easy. One solution Timo and I found was to add a velocity field  $\vec{v} = (x, y, z)$  (with  $\nabla \cdot \vec{v} = 3$ ) to an existing analytical problem, e.g. the Burstedde benchamrk.

## 47 fieldstone\_40: Rayleigh-Taylor instability

This benchmark is carried out in [162, 241, 542] and is based on the analytical solution by Ramberg (1968). It consists of a two-layer system driven by gravity. Free slip are imposed on the sides while no-slip boundary conditions are imposed on the top and the bottom of the box.

Fluid 1 ( $\rho_1, \eta_1$ ) of thickness  $h_1$  overlays fluid 2 ( $\rho_2, \eta_2$ ) of thickness  $h_2$  (with  $h_1 + h_2 = L_y$ ). An initial sinusoidal disturbance of the interface between these layers is introduced and is characterised by an amplitude  $\Delta$  and a wavelength  $\lambda = L_x/2$  as shown in Figure ??.



a) Setup of the experiment, taken from [542]; b) grid setup.

Under this condition, the velocity of the diapiric growth  $v_y$  is given by the relation

$$\frac{v_y}{\Delta} = -K \frac{\rho_1 - \rho_2}{2\eta_2} h_2 g$$

with the dimensionless growth factor  $K$  being

$$K = \frac{-d_{12}}{c_{11}j_{22} - d_{12}i_{21}}$$

and

$$c_{11} = \frac{\eta_1 2\phi_1^2}{\eta_2(\cosh 2\phi_1 - 1 - 2\phi_1^2)} - \frac{2\phi_2^2}{\cosh 2\phi_2 - 1 - 2\phi_2^2} \quad (559)$$

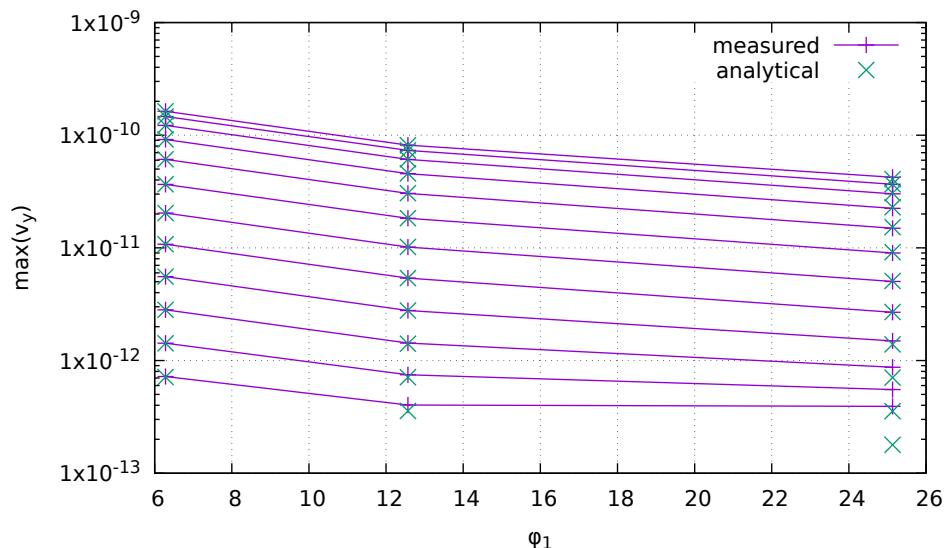
$$d_{12} = \frac{\eta_1(\sinh 2\phi_1 - 2\phi_1)}{\eta_2(\cosh 2\phi_1 - 1 - 2\phi_1^2)} + \frac{\sinh 2\phi_2 - 2\phi_2}{\cosh 2\phi_2 - 1 - 2\phi_2^2} \quad (560)$$

$$i_{21} = \frac{\eta_1 \phi_2 (\sinh 2\phi_1 + 2\phi_1)}{\eta_2(\cosh 2\phi_1 - 1 - 2\phi_1^2)} + \frac{\phi_2 (\sinh 2\phi_2 + 2\phi_2)}{\cosh 2\phi_2 - 1 - 2\phi_2^2} \quad (561)$$

$$j_{22} = \frac{\eta_1 2\phi_1^2 \phi_2}{\eta_2(\cosh 2\phi_1 - 1 - 2\phi_1^2)} - \frac{2\phi_2^3}{\cosh 2\phi_2 - 1 - 2\phi_2^2} \quad (562)$$

$$\phi_1 = \frac{2\pi h_1}{\lambda} \quad (563)$$

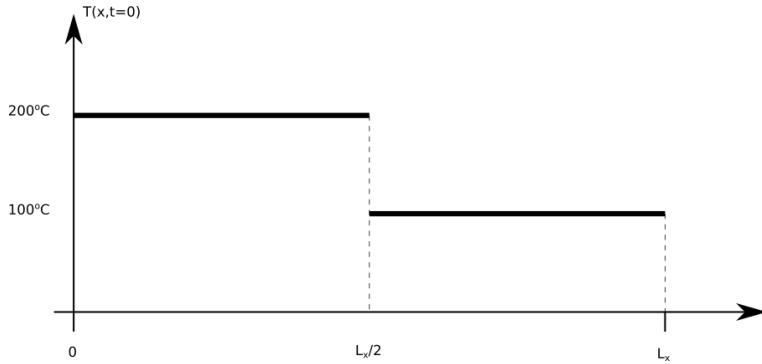
$$\phi_2 = \frac{2\pi h_2}{\lambda} \quad (564)$$



Note that in [542] I fixed  $\lambda = L_x/2$  and varied  $L_x$ . Here I keep  $L_x$  fixed and vary  $\lambda = L_x/2, L_x, 4, L_x/8$ . Each line corresponds to a different value of the viscosity  $\eta_2$ .

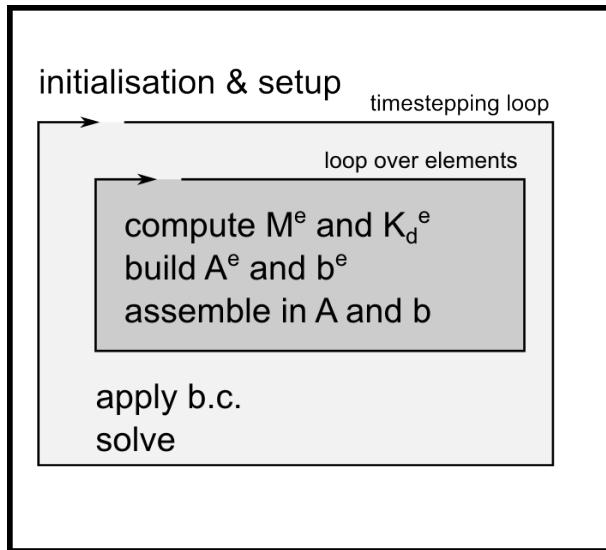
## 48 fieldstone\_42: 1D diffusion

This is the simplest case for a FE code: a 1D (temperature) diffusion problem. It puts into practice what is presented in section 5.1. The initial temperature profile 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 properties of the material are as follows:  $\rho = 3000$ ,  $k = 3$ ,  $C_p = 1000$  and the domain size is  $L_x = 100\text{km}$ . Boundary conditions are  $T(t, x = 0) = 200^\circ\text{C}$  and  $T(t, x = L_x) = 100^\circ\text{C}$ . There are `nelt` elements and `nnx` nodes. All elements are `hx` long. The code will carry out `nstep` timesteps of length `dt` or will stop before that when steady state is reached. The code structure is summarised hereunder:



## 49 fieldstone\_43: the rotating cone

This benchmark originates in [165]. It considers the advection of a product-cosine hill in a prescribed velocity field. The initial temperature is:

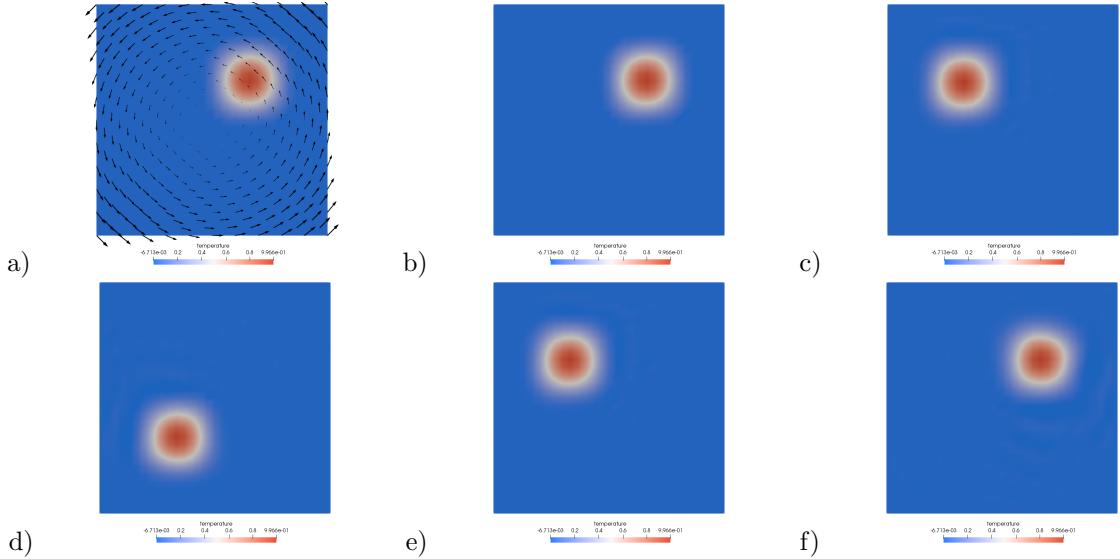
$$T_0(x, y) = \begin{cases} \frac{1}{4} \left(1 + \cos \pi \frac{x-x_c}{\sigma}\right) \left(1 + \cos \pi \frac{y-y_c}{\sigma}\right) & \text{if } (x - x_c)^2 + (y - y_c)^2 \leq \sigma^2 \\ 0 & \text{otherwise} \end{cases} \quad (565)$$

The boundary conditions are  $T(x, y) = 0$  on all four sides of the unit square domain. In what follows we set  $x_c = y_c = 1/6$  and  $\sigma = 0.2$ . The velocity field is analytically prescribed:  $\vec{v} = (-(y - y_c), +(x - x_c))$ .

In what follows we test the time integration scheme by setting  $\alpha_T = 1$  (fully implicit formulation),  $\alpha = 0$  (fully explicit formulation) and  $\alpha_T = 1/2$  (Crank-Nicolson). The timestep is set to  $\delta t = 2\pi/200$ . The density and heat capacity values are set to 1. We monitor the minimum and maximum value of the temperature field, as well as the total thermal energy  $E_T$  in the system during the 200 time steps ( $2\pi$  rotation of the cone):

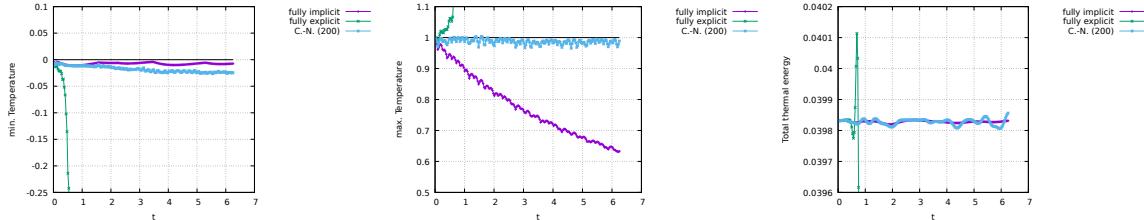
$$E_T = \int_{\Omega} \rho_0 c_p T dV = \int_{\Omega} T dV = |\Omega| \langle T \rangle \quad \text{where} \quad \langle T \rangle = \frac{1}{|\Omega|} \int_{\Omega} T dV$$

The time evolution of the temperature with the Crank-Nicolson algorithm is shown hereunder:



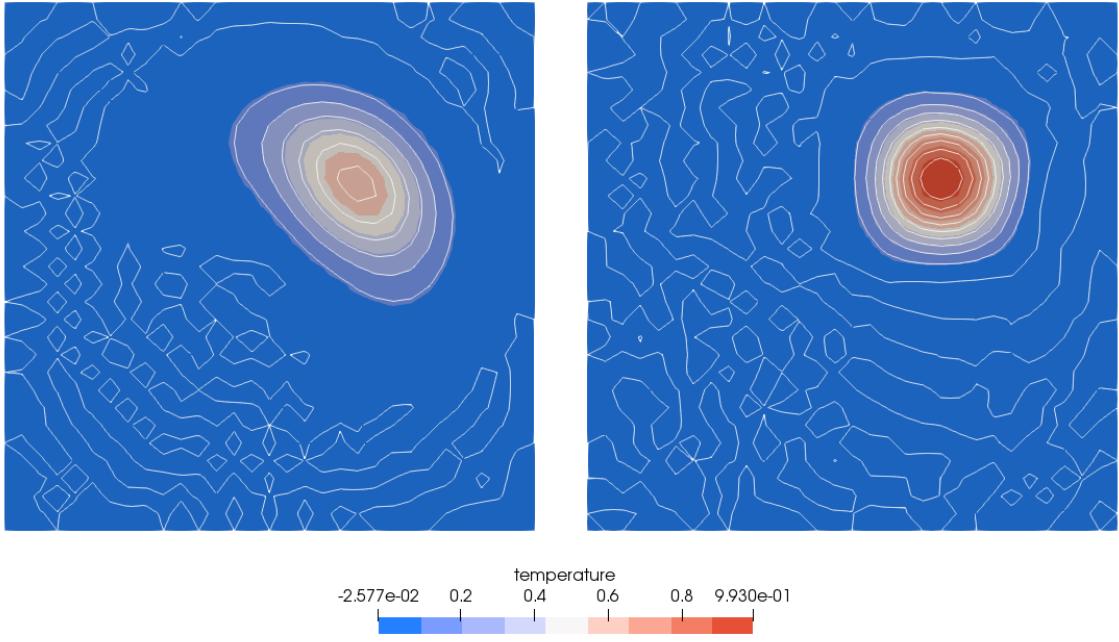
a) Velocity field and initial temperature; b,c,d,e,f) Temperature field at timesteps 0,50,100,150,199.

Turning now to the statistics, we plot  $\min(T)$ ,  $\max(T)$  and  $E_T$  as a function of time:



Time evolution of the min and max temperature and the total energy

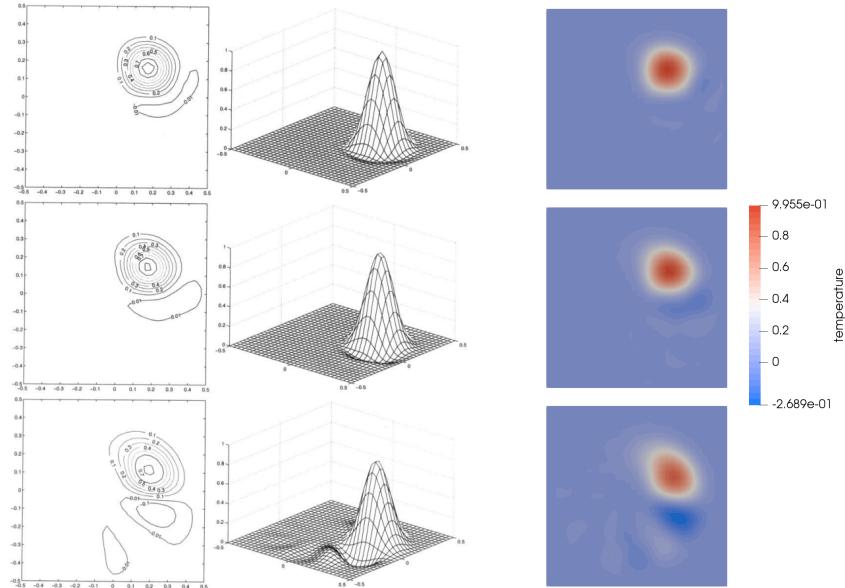
The conclusions are clear: the explicit method diverges quickly and is unusable. The fully implicit and Crank-Nicolson method yield similar energy conservation but the fully-implicit showcases a clear loss in maximum temperature as shown in the following figure:



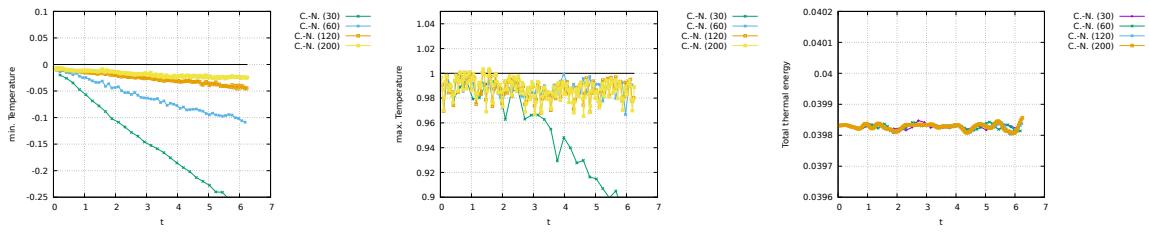
Temperature field after a full rotation with isocontours every 0.1 value.

Left: Fully-implicit; Right: Crank-Nicolson

Finally we can run the experiment (still a  $2\pi$  rotation) with three different time steps ( $\delta t = 2\pi/30, 2\pi/60, 2\pi/120$ ) and we recover very similar results to those presented in [165]:



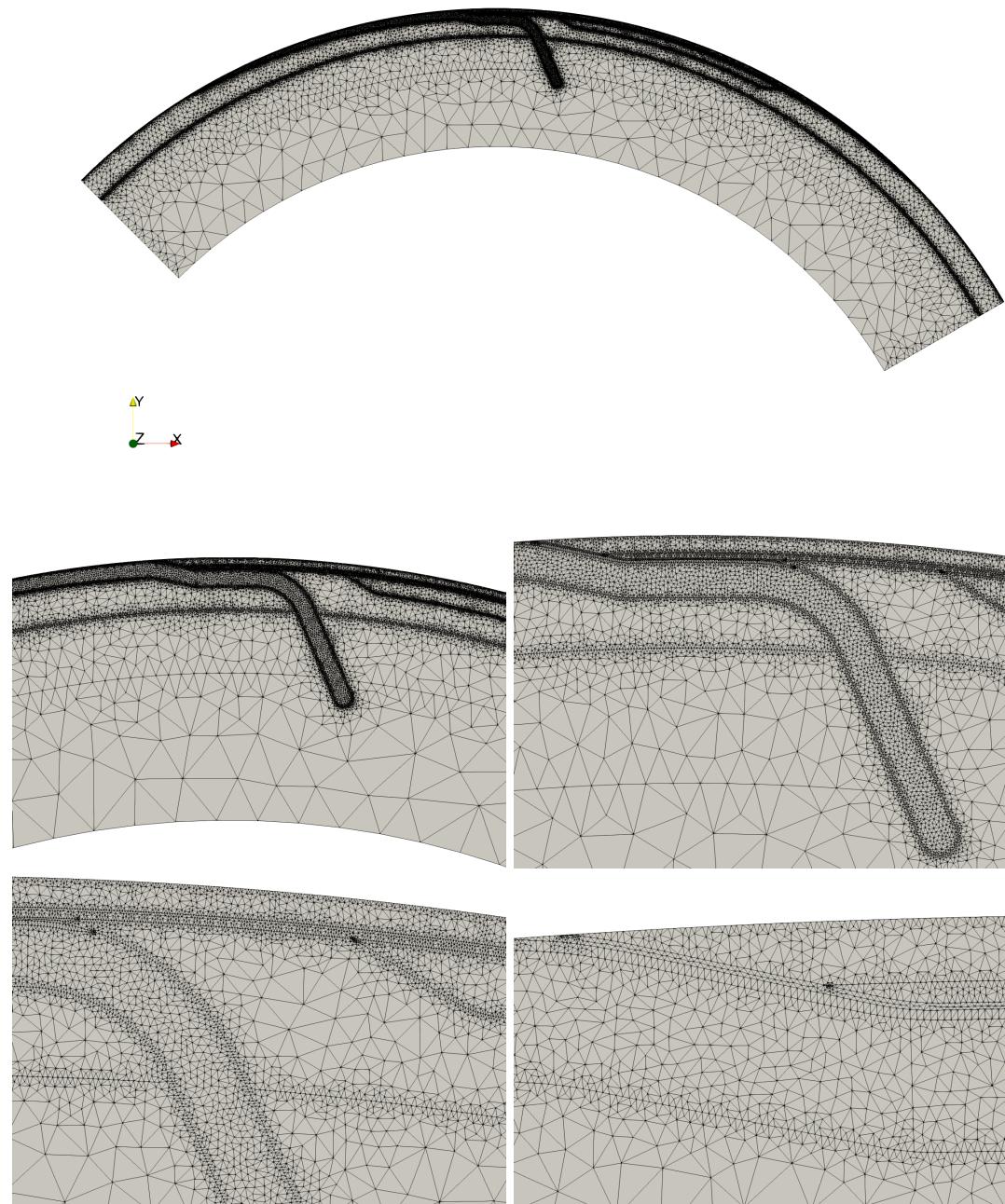
From top to bottom:  $\delta t = 2\pi/120, 2\pi/60, 2\pi/30$  with Crank-Nicolson. Left panel is taken from donea & Huerta [165]



Time evolution of the min and max temperature and the total energy obtained with the Crank-Nicolson algorithm for 4 values of the timestep as indicated by the number between parenthesis.

## 50 fieldstone\_44: the flat slab

I need a list of nodes on the boundary I need a GCOORD.txt file with more decimals I need an even lower resolution grid I need the scaling factors for rho,eta, ...

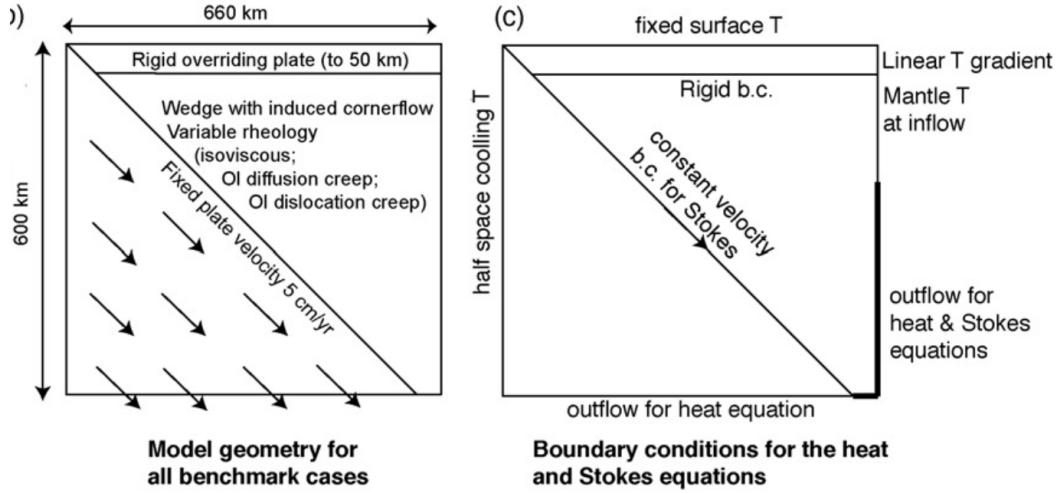


## 51 fieldstone\_45: the corner flow

VERY MUCH WORK IN PROGRESS

This experiment is based on the benchmark paper by van Keken et al, 2008 [570]. It shares similarities with the time dehydration processes in subduction zones work by Magni et al., 2014 [390] and the 3D corner flow study of Plunder et al, 2018[452].

The domain is 660km × 600km. Note that in the original paper the origin of the coordinate system is at the top left while it is at the lower left corner in our code.



As shown in the figure above, the inflow boundaries (at both wedge and trench sides) and top of the model have prescribed temperature. The wedge is assumed to be an incompressible fluid that is driven only by the kinematic forcing of the slab. The wedge is confined by the top of the slab and the base of the rigid overriding plate (located at a depth of 50km). The boundary conditions for the wedge are no-slip below the overriding plate and constant velocity along the top of the slab. The velocity boundary conditions for the boundaries of the wedge are either provided by the Batchelor cornerflow solution (cases 1a and 1b) or based on free inflow/outflow boundaries. The velocity field is discontinuous between the slab and the overriding plate. The velocity in the slab is constant (5cm/yr) and it dips at a 45° angle. There is no radiogenic of shear heating.

The flow is assumed to be incompressible and buoyancy effects are neglected. All the experiments shown in the paper are at steady state, i.e. the temperature field satisfies:

$$\rho C_p \vec{v} \cdot \vec{\nabla} T = \vec{\nabla} \cdot (k \vec{\nabla} T) \quad (566)$$

In the paper a simplified diffusion creep formulation is adopted and the effective diffusion creep viscosity is computed as follows:

$$\eta_{\text{diff}} = A_{\text{diff}} \exp \frac{Q_{\text{diff}}}{RT}$$

The dislocation creep effective viscosity is given by

$$\eta_{\text{disl}} = A_{\text{disl}} \dot{\varepsilon}^{(1-n)/n} \exp \frac{Q_{\text{disl}}}{nRT}$$

Note that in both the activation volume has been set to zero, which decouples pressure from the effective viscosities. Both effective viscosities are limited with a maximum viscosity as follows:

$$\eta_{\text{diff}}^* = \left( \frac{1}{\eta_{\text{diff}}} + \frac{1}{\eta_{\text{max}}} \right)^{-1} \quad \eta_{\text{disl}}^* = \left( \frac{1}{\eta_{\text{disl}}} + \frac{1}{\eta_{\text{max}}} \right)^{-1}$$

The top boundary condition is  $T_{\text{top}} = T(y = L_y) = 273K$ . At the inflow boundary of the wedge (i.e. where  $u < 0$ )<sup>36</sup> temperature is fixed at  $T_0 = 1573K$  and a linear geotherm is used at the left hand

<sup>36</sup>Think about it: it makes little sense to prescribe a temperature where the fluid is leaving the domain

boundary of the overriding plate from 0 to 50 km depth. The temperature at the slab inflow boundary is described by an error-function solution for half-space cooling for 50 Myr:

$$T(x = 0, y) = T_{top} + (T_0 - T_{top}) \operatorname{erf} \frac{L_y - y}{2\sqrt{\kappa t_{50}}}$$

where  $t_{50}$  is the age of the slab.

At the slab and wedge outflow boundaries we prescribe the natural boundary condition (zero curvature) for the heat equation.

the original paper considers multiple cases:

- Case 1a: analytical cornerflow model. The wedge flow is prescribed by the analytical expression for cornerflow [36], so that we do not need to solve for the Stokes equations, only the energy equation.
- Case 1b: dynamical flow in isoviscous wedge I This case is the same as 1a, except that the solution for the wedge flow is determined by solving the Stokes equations while the Batchelor solution is imposed on the inflow and outflow boundaries. This case tests the ability of the numerical method to accurately reproduce the corner flow solution.
- Case 1c: dynamical flow in isoviscous wedge II. Same as case 1b, but with stress-free boundary conditions on the mantle wedge.
- Case 2a: dynamical flow with diffusion creep
- Case 2b: dynamical flow with dislocation creep

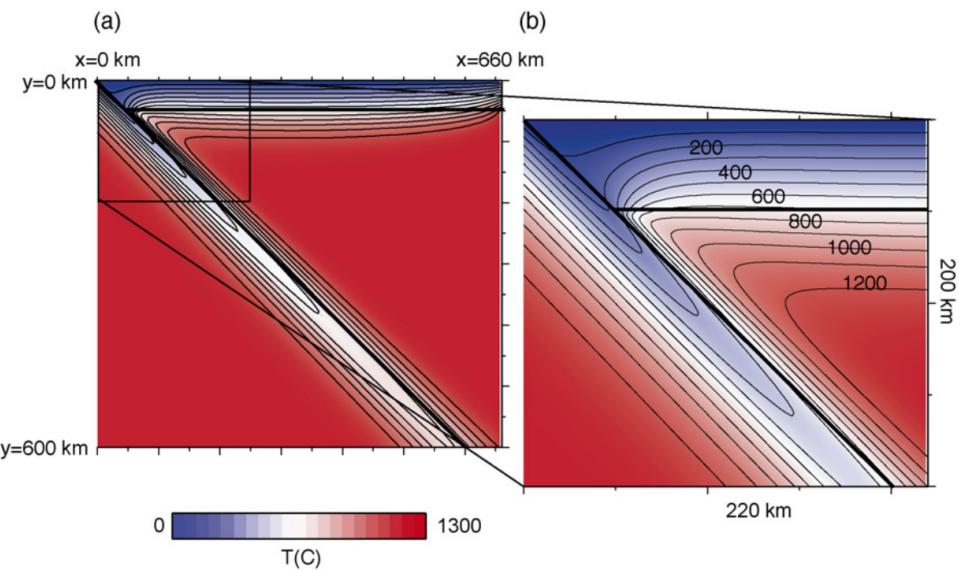
The temperature field as discreted values  $T_{ij}$  on an equidistant grid with 6km spacing, which is a  $111 \times 101$  matrix stored row-wise starting in the top left corner. From this grid the following measurements are extracted for direct comparison:

1. the temperature  $T_{11,11}$  which is at coordinates (60, 60km) and just down-stream from the corner point. This provides therefore one of the most critical tests of accuracy of the numerical codes;
2. the L2 norm of the slab-wedge interface temperature between 0 and 210 km depth defined by

$$T_{\text{slab}} = \sqrt{\frac{1}{36} \sum_{i=1}^{36} T_{ii}^2}$$

3. the L2 norm of the temperature in the triangular part of the tip of the wedge, between 54 and 120 km depth:

$$T_{\text{wedge}} = \sqrt{\frac{1}{78} \sum_{i=10}^{21} \sum_{j=10}^i T_{ij}^2}$$



(a) Temperature prediction for case 1a. The bold lines indicate the top of the slab and base of the overriding plate. (b) Close up of the top left part of the model. Figure taken from [570].

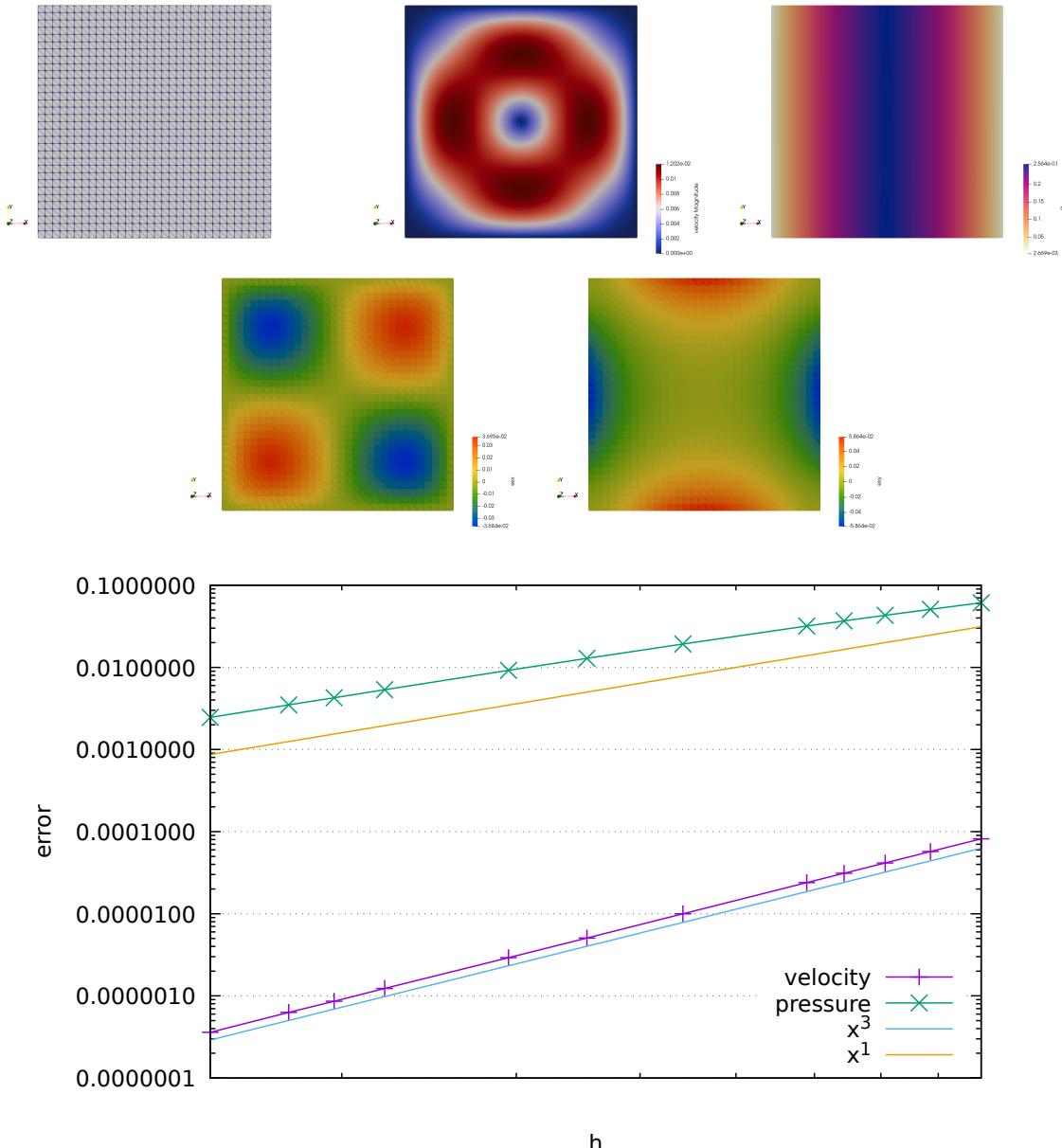
#### QUESTIONS

are the nodes on the 45 degree diagonal in the slab or in the mantle ?!

## 52 fieldstone\_46: MMS1 with Crouzeix-Raviart elements

This stone showcases the Crouzeix-Raviart element (see Section 6.2.9) used to solve the analytical problem "Donea & Huerta" (see Section 7.4.1).

Out of convenience the pressure is set to zero at location  $(x, y) = (1, 1)$ , so that the analytical solution is now  $p(x, y) = x(1 - x)$ .



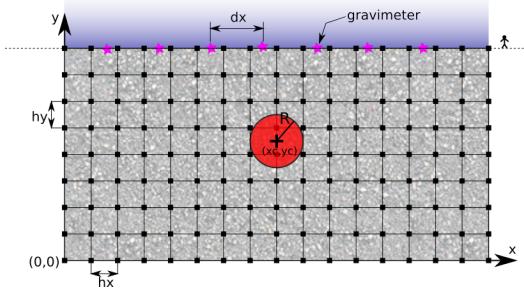
## 53 fieldstone: Gravity: buried sphere

Before you proceed further, please read :

[http://en.wikipedia.org/wiki/Gravity\\_anomaly](http://en.wikipedia.org/wiki/Gravity_anomaly)

<http://en.wikipedia.org/wiki/Gravimeter>

Let us consider a vertical domain  $Lx \times Ly$  where  $L_x = 1000\text{km}$  and  $L_y = 500\text{km}$ . This domain is discretised by means of a grid which counts  $nnp = nnx \times nny$  nodes. This grid then counts  $nel = nelx \times nely = (nnx - 1) \times (nny - 1)$  cells. The horizontal spacing between nodes is  $hx$  and the vertical spacing is  $hy$ .



Assume that this domain is filled with a rock type which mass density is given by  $\rho_{medium} = 3000\text{kg/m}^3$ , and that there is a circular inclusion of another rock type ( $\rho_{sphere} = 3200\text{kg/m}^3$ ) at location  $(xsphere, ysphere)$  of radius  $rsphere$ . The density in the system is then given by

$$\rho(x, y) = \begin{cases} \rho_{sphere} & \text{inside the circle} \\ \rho_{medium} & \text{outside the circle} \end{cases}$$

Let us now assume that we place  $nsurf$  gravimeters at the surface of the model. These are placed equidistantly between coordinates  $x = 0$  and coordinates  $x = Lx$ . We will use the arrays  $xsurf$  and  $ysurf$  to store the coordinates of these locations. The spacing between the gravimeters is  $\delta_x = Lx/(nsurf - 1)$ .

At any given point  $(x_i, y_i)$  in a 2D space, one can show that the gravity anomaly due to the presence of a circular inclusion can be computed as follows:

$$g(x_i, y_i) = 2\pi G(\rho_{sphere} - \rho_0)R^2 \frac{y_i - ysphere}{(x_i - xsphere)^2 + (y_i - ysphere)^2} \quad (567)$$

where  $r_{sphere}$  is the radius of the inclusion,  $(xsphere, ysphere)$  are the coordinates of the center of the inclusion, and  $\rho_0$  is a reference density.

However, the general formula to compute the gravity anomaly at a given point  $(x_i, y_i)$  in space due to a density anomaly of any shape is given by:

$$g(x_i, y_i) = 2G \int \int_{\Omega} \frac{\Delta\rho(x, y)(y - y_i)}{(x - x_i)^2 + (y - y_i)^2} dxdy \quad (568)$$

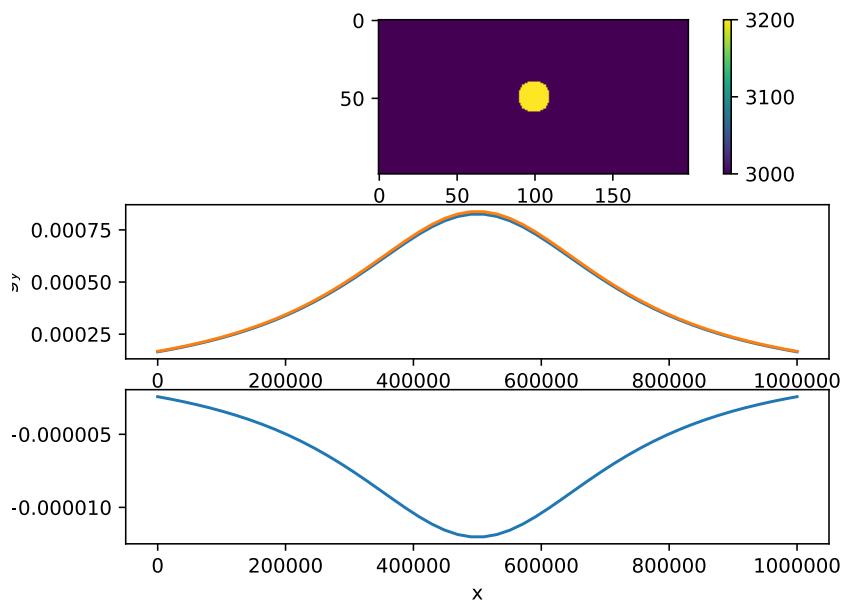
where  $\Omega$  is the area of the domain on which the integration is to be carried out. Furthermore the density anomaly can be written :  $\Delta\rho(x, y) = \rho(x, y) - \rho_0$ . We can then carry out the integration for each cell and sum their contributions:

$$g(x_i, y_i) = 2G \sum_{ic=1}^{nel} \int \int_{\Omega_e} \frac{(\rho(x, y) - \rho_0)(y - y_i)}{(x - x_i)^2 + (y - y_i)^2} dxdy \quad (569)$$

where  $\Omega_e$  is now the area of a single cell. Finally, one can assume the density to be constant within each cell so that  $\rho(x, y) \rightarrow \rho(ic)$  and  $\int \int_{\Omega_e} dxdy \rightarrow hx \times hy$  and then

$$g(x_i, y_i) = 2G \sum_{ic=1}^{nel} \frac{(\rho(ic) - \rho_0)(y(ic) - y_i)}{(x(ic) - x_i)^2 + (y(ic) - y_i)^2} s_x s_y \quad (570)$$

We will then use the array  $gsurf$  to store the value of the gravity anomaly measured at each gravimeter at the surface.



### To go further

- explore the effect of the size of the inclusion on the gravity profile.
- explore the effect of the  $\rho_0$  value.
- explore the effect of the grid resolution.
- measure the time that is required to compute the gravity. How does this time vary with nsurf ? how does it vary when the grid resolution is doubled ?
- Assume now that  $\rho_2 < \rho_1$ . What does the gravity profile look like ?
- what happens when the gravimeters are no more at the surface of the Earth but in a satellite ?
- if you feel brave, redo the whole exercise in 3D...

## 54 Problems, to do list and projects for students

- Darcy flow. redo WAFLE (see <http://cedrichieulot.net/wafle.html>)
- carry out critical Rayleigh experiments for various geometries/aspect ratios. Use Arie's notes.
- Newton solver
- elasticity with markers
- Indentor/punch with stress b.c. ?
- chunk grid
- read in crust 1.0 in 2D on chunk
- compute gravity based on tetrahedra
- compare Q2 with Q2-serendipity
- NS a la [http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible\\_Ex2.htm](http://ww2.lacan.upc.edu/huerta/exercises/Incompressible/Incompressible_Ex2.htm)
- produce example of mckenzie slab temperature
- write about impose bc on el matrix
- constraints
- discontinuous galerkin
- formatting of code style
- navier-stokes ? (LUKAS) use dohu matlab code
- nonlinear poiseuille
- Finish nonlinear cavity case5.
- write about mappings
- write about stream functions
- free surface
- zaleski disk advection
- better yet simple matrix storage ?
- write Scott about matching compressible2 setup with his paper
- deal with large matrices.
- compositions, marker chain
- free-slip bc on annulus and sphere . See for example p540 Gresho and Sani book. find book [161]. also check [183] !!
- non-linear rheologies (two layer brick spmw16, tosn15)
- Picard vs Newton
- including phase changes (w. R. Myhill)
- compute strainrate in middle of element or at quad point for punch?
- GEO1442 code

- GEO1442 indenter setup in plane ?
- in/out flow on sides for lith modelling
- Fehlberg RK advection
- redo puth17 2 layer experiment
- create stone for layeredflow (see folder one up)
- SIMPLE a la p667 [330]
- implement mms5 7.4.5
- implement mms7 7.4.7
- implement/monitor div v
- shape fct, trial fct, basis fct vs test fct doc
- Delaunay triangulation, Voronoi, stripack
- symmetric vs gradient formulation of Stokes
- write/draw the whole FEM process for a 4x3 grid for compgeo
- Sphere drag in a visco-plastic fluid [159]
- mention Lattice-Boltzmann in geosciences [307]
- lukas' 2D and 3D benchmark
- ROTATING disc
- cylindrical footing on (elasto)-viscous medium - analytical solution, Haskell, etc ...
- lev hager 2008 RT instability with anisotropic visc

`iCells; iDataArray type=Int32 Name=connectivity .../i; iDataArray type=Int32 Name=offsets .../i;  
iDataArray type=UInt8 Name=types .../i; i/Cells;`

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

open questions: what does it mean to have a negative pressure ? should we threshold it when computing yield strength ?

## A Three-dimensional applications

In the following table I list many papers which showcase high-resolution models of geodynamical processes (subduction, rifting, mantle flow, plume transport, ...). Given the yearly output of our community and the long list of journal in which research can be disseminated, this list can not be exhaustive.

Ref.	topic	resolution
[32]	Small-scale sublithospheric convection in the Pacific	$448 \times 56 \times 64$
[521]	Migration and morphology of subducted slabs in the upper mantle	$50 \times 50 \times 25$
[466]	Subduction scissor across the South Island of New Zealand	$17 \times 9 \times 9$
[404]	Influence of a buoyant oceanic plateau on subduction zones	$80 \times 40 \times 80$
[120]	Subduction dynamics, origin of Andean orogeny and the Bolivian orocline	$96 \times 96 \times 64$
[181]	Feedback between rifting and diapirism, ultrahigh-pressure rocks exhumation	$100 \times 64 \times 20$
[14]	Numerical modeling of upper crustal extensional systems	$160 \times 160 \times 12$
[15]	Rift interaction in brittle-ductile coupled systems	$160 \times 160 \times 23$
[362]	Kinematic interpretation of the 3D shapes of metamorphic core complexes	$67 \times 67 \times 33$
[323]	Role of rheology and slab shape on rapid mantle flow: the Alaska slab edge	$960 \times 648 \times 160$
[119]	Complex mantle flow around heterogeneous subducting oceanic plates	$96 \times 96 \times 64$
[87]	Oblique rifting and continental break-up	$150 \times 50 \times 30$
[55]	Influence of mantle plume head on dynamics of a retreating subduction zone	$80 \times 40 \times 80$
[86]	Rift to break-up evolution of the Gulf of Aden	$83 \times 83 \times 40$
[88]	Thermo-mechanical impact of plume arrival on continental break-up	$100 \times 70 \times 20$
[121]	Subduction and slab breakoff controls on Asian indentation tectonics	$96 \times 96 \times 64$
[190]	Modeling of upper mantle deformation and SKS splitting calculations	$96 \times 64 \times 96$
[494]	Backarc extension/shortening, slab induced toroidal/poloidal mantle flow	$352 \times 80 \times 64$
[391]	Sediment transports in the context of oblique subduction modelling	$500 \times 164 \times 100$
[631]	Crustal growth at active continental margins	$404 \times 164 \times 100$
[364]	Dynamics of India-Asia collision	$257 \times 257 \times 33$
[596]	Strain-partitioning in the Himalaya	$256 \times 256 \times 40$
[374]	Collision of continental corner from 3-D numerical modeling	$500 \times 340 \times 164$
[460]	Dependence of mid-ocean ridge morphology on spreading rate	$196 \times 196 \times 100$
[189]	Mid mantle seismic anisotropy around subduction zones	$197 \times 197 \times 53$
[298]	Oblique rifting of the Equatorial Atlantic	$120 \times 80 \times 20$
[415]	Dynamics of continental accretion	$256 \times 96 \times 96$
[485]	Thrust wedges: infl. of decollement strength on transfer zones	$309 \times 85 \times 149$
[107]	Asymmetric three-dimensional topography over mantle plumes	$500 \times 500 \times 217$
[547]	modelled crustal systems undergoing orogeny and subjected to surface processes	$96 \times 32 \times 14$
[377]	Thermo-mechanical modeling ontinental rifting and seafloor spreading	$197 \times 197 \times 197$

## B Codes in geodynamics

In what follows I make a quick inventory of the main codes of computational geodynamics, for crust, lithosphere and/or mantle modelling.

in order to find all CIG-codes citations go to: <https://geodynamics.org/cig/news/publications-refbase/>

- ABAQUS [233] [359] [397] [424] [444]
- ADELI [295] [582] [66] [67] [231] [257] [590] [123] [124]
- ASPECT

This code is hosted by CIG at <https://geodynamics.org/cig/software/aspect/>  
[34] [358] [26] [550] [154] [223] [621] [297] [152] [299] [483] [484] [25] [543] [83] [433] [530] [618] [153]  
[435] [259] [300] [224] [443] [457] [82] [39] [523] [143] [383]

- CHIC [430]
- CitcomS and CITCOMCU

These codes are hosted by CIG at <https://geodynamics.org/cig/software/citcomcu/> and <https://geodynamics.org/cig/software/citcoms/>.

[516] [418] [624] [569] [627] [287] [59] [533] [568] [140] [60] [514] [61] [46] [449] [531] [49] [141] [62]  
[626] [395] [33] [480] [420] [139] [258] [222] [625] [305] [376] [21] [620] [18] [203] [103] [98] [578] [32]  
[104] [619] [47] [367] [565] [371] [381] [323] [63] [69] [318] [629] [515] [304] [74] [75] [324] [470] [432]  
[23] [142] [202] [97] [335] [24] [589] [503] [7] [390] [632] [73] [70] [506] [155] [566] [591] [588] [587] [294]  
[535] [372] [595] [594] [389] [210] [296] [302] [345] [405] [212]

[cross check with CIG database](#)

- CONMAN This code is hosted by CIG at <https://geodynamics.org/cig/software/conman/>  
[349] [321] [346] [347] [322] [48] [348]
- CONVRS [615] [614]
- DOUAR  
[80] [546] [608] [81] [386] [421] [596] [427]
- DYNEARTHSSOL [131]
- M-DOODS [609] [607]
- FENICS [12]
- GAIA
- GALE

This code is hosted by CIG at <https://geodynamics.org/cig/software/gale/>  
[195] [58] [147] [362] [23]

ADD:

Cruz, L.; Malinski, J.; Hernandez, M.; Take, A.; Hilly, G. Erosional control of the kinematics of the Aconcagua fold-and-thrust belt from numerical simulations and physical experiments 2011 Geology

Goyette, S.; Takatsuka, M.; Clark, S.; Mller, R.D.; Rey, P.; Stegman, D.R. Increasing the usability and accessibility of geodynamic modelling tools to the geoscience community: UnderworldGUI 2008 Visual Geosciences

Li, Y.; Qi, J. Salt-related Contractual Structure and Its Main Controlling Factors of Kelasu Structural Zone in Kuqa Depression: Insights from Physical and Numerical Experiments 2012 Procedia Engineering

- GTECTON [267] [268] [93] [94] [269] [207] [385] [29] [30] [398]

- ELEFANT  
[550] [388] [96] [361] [543] [452]
- ELLIPSIS  
[416] [434] [417] [172] [436] [466] [371] [368]
- FANTOM  
[542] [14] [15] [16] [187] [547] [185] [186]  
FDCON [184] [214] [213]
- FLUIDITY [157] [221]
- IFISS: Incompressible Flow Iterative Solution Solver is a MATLAB package that is a very useful tool for people interested in learning about solving PDEs. IFISS includes built-in software for 2D versions of: the Poisson equation, the convection-diffusion equation, the Stokes equations and the Navier-Stokes equations.  
<https://personalpages.manchester.ac.uk/staff/david.silvester/ifiss/>
- the I2(3)E(L)VIS code  
[250][251][249] [252][253][247] [102] [91][239][264][243] [240][263] [497][242][555][191][630] [244] [237][428] [169][171][375][238][246] [144][168][617] [374][423][391][562][561][631][173][245][402] [170][460][485][584][31][378][524][3] [107][265][31][563] [167][556][486][248][486] [4][392][201] [354]
- I3MG [189]
- LAMEM [497] [338] [363] [406] [364] [138] [197] [196] [458] [198] [137] [339]
- LAPEX2D [511] [91] [27] [497] [512]
- LITMOD [215] [5] [6] [216]
- MARC [426] [425]
- MILAMIN [151] [610] [228] [387] [229] [540] [327] [401]
- PARAVOZ/FLAMAR [453] [109] [105] [28] [136] [291] [234] [286] [236] [551] [108] [605] [110] [606] [232] [608] [20] [231] [283] [230] [192] [208] [220] [106] [603] [403] [227] [163]
- PINK3D [585]
- PLASTI [217]
- PTATIN [445] [3] [407] [333] [332]
- RHEA [111] [518] [11] [113]
- SAMOVAR [178]
- SEPRAN [557] [583] [134] [567] [574] [575] [576] [577] [502] [379] [380] [573] [78] [77] [451] [581] [558] [527] [571] [559] [50] [127] [19] [128] [419] [560] [622]
- SLIM3D  
[455] [474] [87] [88] [86] [85] [298] [353] [135]
- SLOMO [336]  
SNAC [132]
- SOPALE  
[599] [40] [218] [179] [44] [598] [463] [45] [42] [315] [464] [314] [579] [600] [468] [43] [465] [462] [180] [226] [225] [317] [461] [505] [316] [149] [414] [504] [592] [593] [341] [41] [95] [271] [510] [10] [9] [272] [467] [148] [118] [313] [273] [274] [352] [270] [325] [116] [126] [199] [200] [266] [275] [351] [429] [326] [261] [13] [117] [301] [382] [115]

- STAGYY [482] [613] [145]
- SUBMAR [400] [399] [481]
- SULEC SULEC is a finite element code that solves the incompressible Navier-Stokes equations for slow creeping flows. The code is developed by Susan Ellis (GNS Sciences, NZ) and Susanne Buitier (NGU).  
[471] [181] [92] [537] [144] [280] [254] [255] [472] [422] [638] [538]
- TERRA [101] [100] [446] [602] [601] [156] [564]
- YACC [550] [549]
- UNDERWORLD 1&2 [520] [417] [493] [366] [439] [122] [404] [521] [519] [194] [120] [119] [55] [494] [190] [193] [469] [56] [495] [508] [509] [440] [350] [411] [490] [612]
- VEMAN [57]

## C Matrix properties

### C.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  $\mathbf{M} = \mathbf{LL}^T$  (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.

### C.2 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  $\mathbf{D}$  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  $\vec{y}$ . 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**.

$$\mathbb{S}_p \cdot \vec{p} = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \vec{f} \quad \text{with} \quad \mathbb{S}_p = \mathbb{G}^T \cdot \mathbb{K}^{-1} \cdot \mathbb{G} + \mathbb{C}$$

## D Dont 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) CITCOM [1], DOUAR [8], FANTOM [2], IELVIS [5], LaMEM [3], pTatin [4], SLIM3D [10], SOPALE [7], StaggYY [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 owns 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 taylored 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 ones 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

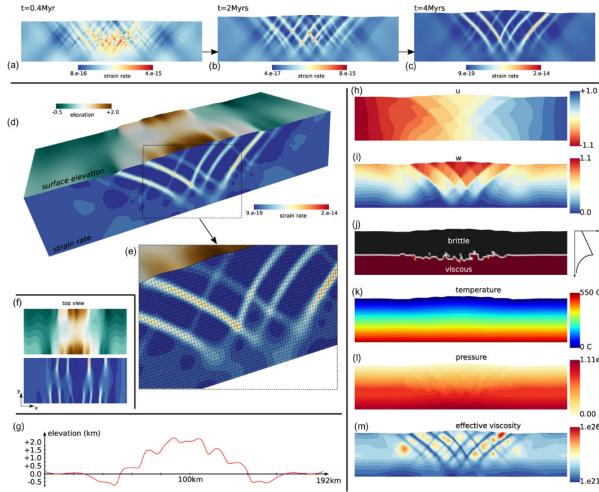
#### References

- [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

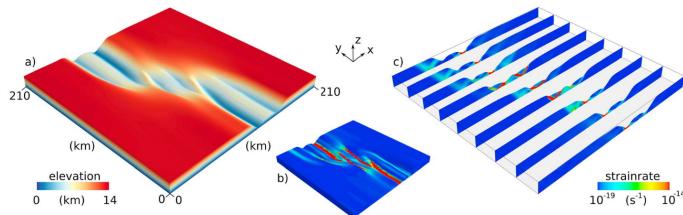
## E A FANTOM, an ELEFANT and a GHOST

While a post-doctoral researcher at Bergen University I developed the FANTOM code. Here is what other people and I have published with it:

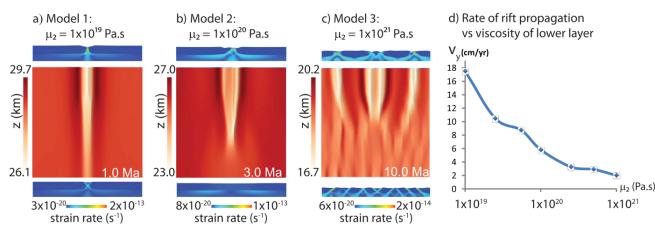
- *FANTOM : two- and three-dimensional numerical modelling of creeping flows for the solution of geological problems*, C. Thieulot, Physics of the Earth and Planetary Interiors, 188, 2011.



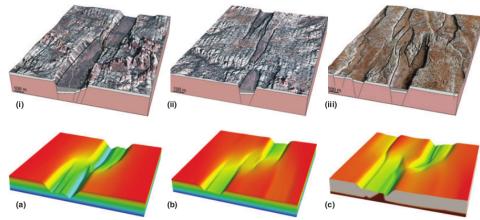
- *Three-dimensional numerical modeling of upper crustal extensional systems*, V. Allken, R.S. Huismans and C. Thieulot, JGR 116, 2011. <https://doi.org/10.1029/2011JB008319>



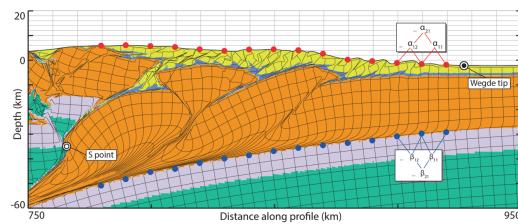
- *Factors controlling the mode of rift interaction in brittle-ductile coupled systems: A 3D numerical study*, V. Allken, R.S. Huismans and C. Thieulot, Geochem. Geophys. Geosyst. 13(5), 2012. <https://doi.org/10.1029/2012GC004077>



- *3D numerical modelling of graben interaction and linkage: a case study of the Canyonlands grabens, Utah*, V. Allken, R.S. Huismans, Haakon Fossen and C. Thieulot, Basin Research, 25, 1-14, 2013. <https://doi.org/10.1111/bre.12010>

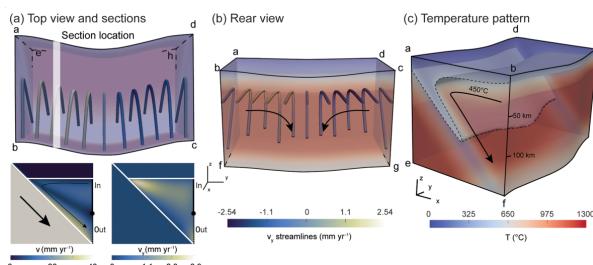


- Three-dimensional numerical simulations of crustal systems undergoing orogeny and subjected to surface processes, C. Thieulot, P. Steer and R.S. Huismans, *Geochem. Geophys. Geosyst.*, 15, 2014. doi:10.1002/2014GC005490
- Extensional inheritance and surface processes as controlling factors of mountain belt structure, Z. Erdős, R.S. Huismans, P. van der Beek, and C. Thieulot, *J. Geophys. Res. Solid Earth*, 119, 2014. doi:10.1002/2014JB011408
- First-order control of syntectonic sedimentation on crustal-scale structure of mountain belts, Z. Erdős, R.S. Huismans, P. van der Beek, *J. Geophys. Res. Solid Earth*, 120, 5362-5377, 2015. doi:10.1002/2014JB011785
- Control of increased sedimentation on orogenic fold-and-thrust belt structure - insights into the evolution of the Western Alps, Z. Erdős, R.S. Huismans and P. van der Beek, *Solid Earth*, 10, 391-404, 2019. <https://doi.org/10.5194/se-10-391-2019>

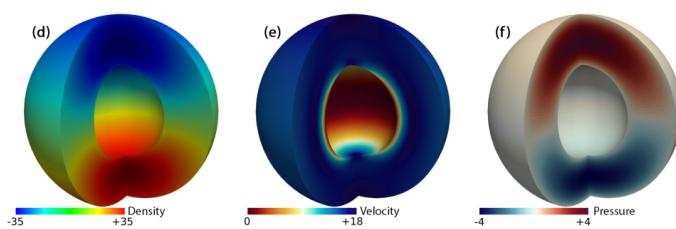


Upon my arrival at Utrecht University in 2012 I started working on a more flexible code, called ELEFANT, which has since very much diverged from FANTOM.

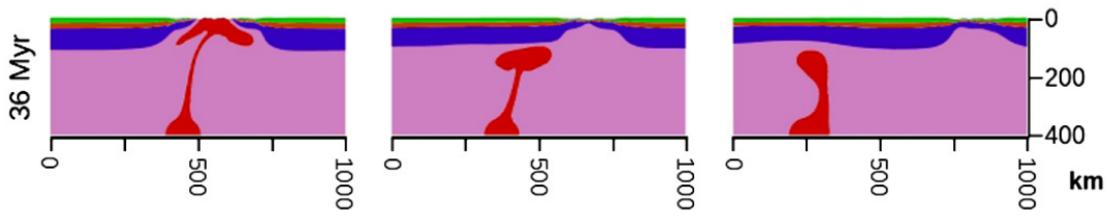
- The effect of obliquity on temperature in subduction zones: insights from 3-D numerical modeling, A. Plunder, C. Thieulot and D.J.J. van Hinsbergen, *Solid Earth* 9, 759-776, 2018. <https://doi.org/10.5194/se-9-759-2018>



- Analytical solution for viscous incompressible Stokes flow in a spherical shell, C. Thieulot, *Solid Earth* 8, 1181-1191, 2017. <https://doi.org/10.5194/se-8-1181-2017>



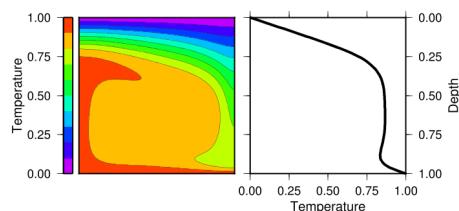
- *Lithosphere erosion and continental breakup: interaction of extension, plume upwelling and melting*, A. Lavecchia, C. Thieulot, F. Beekman, S. Cloetingh and S. Clark, E.P.S.L. 467, p89-98, 2017.



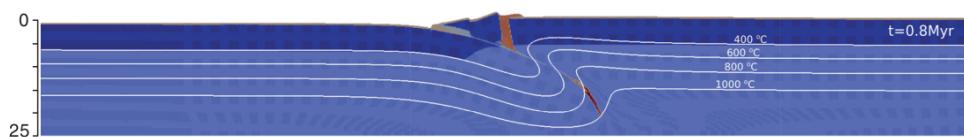
- *Benchmarking numerical models of brittle thrust wedges*, Susanne J.H. Buiter, Guido Schreurs, Markus Albertz, Taras V. Gerya, Boris Kaus, Walter Landry, Laetitia le Pourhiet, Yury Mishin, David L. Egholm, Michele Cooke, Bertrand Maillot, Cedric Thieulot, Tony Crook, Dave May, Pauline Souloumiac, Christopher Beaumont Journal of Structural Geology 92, p140-177, 2016. <https://doi:10.1016/j.jsg.2016.03.003>



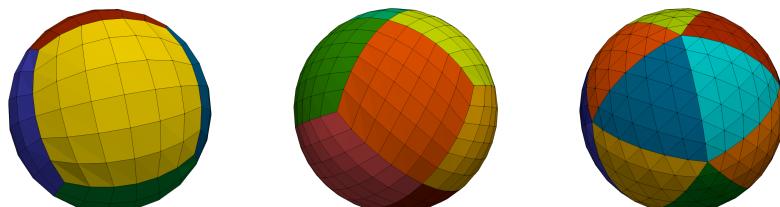
- *A community benchmark for viscoplastic thermal convection in a 2-D square box*, N. Tosi, C. Stein, L. Noack, C. Huettig, P. Maierova, H. Samuel, D.R. Davies, C.R. Wilson, S.C. Kramer, C. Thieulot, A. Glerum, M. Fraters, W. Spakman, A. Rozel, P.J. Tackley, Geochem. Geophys. Geosyst. 16, doi:10.1002/2015GC005807, 2015.



- *Dynamics of intraoceanic subduction initiation: 1. Oceanic detachment fault inversion and the formation of supra-subduction zone ophiolites*, M. Maffione, C. Thieulot, D.J.J. van Hinsbergen, A. Morris, O. Pluempert and W. Spakman, Geochem. Geophys. Geosyst. 16, p1753-1770, 2015.



- *GHOST: Geoscientific Hollow Sphere Tessellation*, C. Thieulot, Solid Earth, 9, 11691177, 2018. <https://doi.org/10.5194/se-9-1169-2018>



## F Some useful Python commands

### F.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)
```

### F.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))
```

## G Some useful maths

### G.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

$$x = \frac{1}{\det(\mathbf{A})} \begin{vmatrix} b_1 & a_{21} \\ b_2 & a_{22} \end{vmatrix} \quad y = \frac{1}{\det(\mathbf{A})} \begin{vmatrix} a_{11} & b_1 \\ a_{21} & b_2 \end{vmatrix}$$

### G.2 Inverse of a 3x3 matrix

Let us consider the 3x3 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 2x2 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

$$x = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} b_1 & a_{12} & a_{13} \\ b_2 & a_{22} & a_{23} \\ b_3 & a_{32} & a_{33} \end{vmatrix} \quad y = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} a_{11} & b_1 & a_{13} \\ a_{21} & b_2 & a_{23} \\ a_{31} & b_3 & a_{33} \end{vmatrix} \quad z = \frac{1}{\det(\mathbf{M})} \begin{vmatrix} a_{11} & a_{12} & b_1 \\ a_{21} & a_{22} & b_2 \\ a_{31} & a_{32} & b_3 \end{vmatrix}$$

## References

- [1] PhD thesis.
- [2]
- [3] *pTatin3D: High-Performance Methods for Long-Term Lithospheric Dynamics*, 2014.
- [4] A. Koptev aand E. Burov, E. Calais, S. Leroy, T. Gerya, L. Guillou-Frottier, and S. Cloetingh. Contrasted continental rifting via plume-craton interaction: Applications to Central East African Rift. *Geoscience Frontiers*, 7:221–236, 2016.
- [5] J. C. Afonso, M. Fernandez, G. Ranalli, W.L. Griffin, and J.A.D. Connolly. Integrated geophysical-petrological modeling of the lithosphere and sublithospheric upper mantle: Methodology and applications. *Geochem. Geophys. Geosyst.*, 9(5):doi:10.1029/2007GC001834, 2008.
- [6] J.C. Afonso, G. Ranalli, and M. Fernandez. Density structure and buoyancy of the oceanic lithosphere revisited. *Geophys. Res. Lett.*, 34:L10302, 2007.
- [7] R. Agrusta, J. van Hunen, and S. Goes. The effect of metastable pyroxene on the slab dynamics. *Geophys. Res. Lett.*, 41:8800–8808, 2014.
- [8] M. Albers. A local mesh refinement multigrid method for 3D convection problems with strongly variable viscosity. *J. Comp. Phys.*, 160:126–150, 2000.
- [9] M. Albertz and C. Beaumont. An investigation of salt tectonic structural styles in the Scotian Basin, offshore Atlantic Canada: 2. Comparison of observations with geometrically complex numerical models. *Tectonics*, 29(TC4018), 2010.
- [10] M. Albertz, C. Beaumont, J.W. Shimeld, S.J. Ingsand, and S. Gradmann. An investigation of salt tectonic structural styles in the Scotian Basin, offshore Atlantic Canada: Part 1, comparison of observations with geometrically simple numerical models. *Tectonics*, 29, 2010.
- [11] L. Alisic, M. Gurnis, G. Stadler, C. Burstedde, and O. Ghattas. Multi-scale dynamics and rheology of mantle flow with plates. *J. Geophys. Res.*, 117, 2012.
- [12] L. Alisic, J.F. Rudge, R.F. Katz, G.N. Wells, and S. Rhebergen. Compaction around a rigid, circular inclusion in partially molten rock. *J. Geophys. Res.*, 119:5903–5920, 2014.
- [13] J. Allen and C. Beaumont. Continental Margin Syn-Rift Salt Tectonics at Intermediate Width Margins. *Basin Research*, page doi: 10.1111/bre.12123, 2014.
- [14] V. Allken, R. Huismans, and C. Thieulot. Three dimensional numerical modelling of upper crustal extensional systems. *J. Geophys. Res.*, 116:B10409, 2011.
- [15] V. Allken, R. Huismans, and C. Thieulot. Factors controlling the mode of rift interaction in brittle-ductile coupled systems: a 3d numerical study. *Geochem. Geophys. Geosyst.*, 13(5):Q05010, 2012.
- [16] 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. *Basin Research*, 25:1–14, 2013.
- [17] J.D. Anderson. *Computational Fluid Dynamics*. McGraw-Hill, 1995.
- [18] E.R. Andrews and M.I. Billen. Rheologic controls on the dynamics of slab detachment. *Tectonophysics*, 464:60–69, 2009.
- [19] A. Androvicova, H. Čížková, and A. van den Berg. The effects of rheological decoupling on slab deformation in the Earth’s upper mantle. *Stud. Geophys. Geod.*, 57:460–481, 2013.
- [20] S. Angiboust, S. Wolf, E. Burov, P. Agard, and P. Yamato. Effect of fluid circulation on subduction interface tectonic processes: Insights from thermo-mechanical numerical modelling. *Earth Planet. Sci. Lett.*, 357-358:238–248, 2012.

- [21] J.J. Armitage, T.J. Henstock, T.A. Minshull, and J.R. Hopper. Lithospheric controls on melt production during continental breakup at slow rates of extension: Application to the North Atlantic. *Geochem. Geophys. Geosyst.*, 10(6), 2009.
- [22] D.N. Arnold, F. Brezzi, and M. Fortin. A stable finite element for the Stokes equation. *Calcolo*, XXI(IV):337–344, 1984.
- [23] P.-A. Arrial and M.I. Billen. Influence of geometry and eclogitization on oceanic plateau subduction . *Earth Planet. Sci. Lett.*, 363:34–43, 2013.
- [24] 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. *Geosci. Model Dev.*, 7:2065–2076, 2014.
- [25] J. Austermann, J. X. Mitrovica, P. Huybers, and A. Rovere. Detection of a dynamic topography signal in last interglacial sea-level records. *Science Advances*, 3(7):1700457, 2017.
- [26] J. Austermann, D. Pollard, J. X. Mitrovica, R. Moucha, A. M. Forte, R. M. DeConto, D. B. Rowley, and M. E. Raymo. The impact of dynamic topography change on antarctic ice sheet stability during the mid-pliocene warm period. *Geology*, 43(10):927–930, 2015.
- [27] A. Babeyko and S. Sobolev. High-resolution numerical modeling of stress distribution in visco-elasto-plastic subducting slabs. *Lithos*, 103:205–216, 2008.
- [28] A.Yu. Babeyko, S.V. Sobolev, R.B. Trumbull, O. Oncken, and L.L. Lavie. Numerical models of crustal scale convection and partial melting beneath the Altiplano-Puna plateau. *Earth Planet. Sci. Lett.*, 199:373–388, 2002.
- [29] M. Baes, R. Govers, and R. Wortel. Subduction initiation along the inherited weakness zone at the edge of a slab: Insights from numerical models. *Geophys. J. Int.*, 184:991–1008, 2011.
- [30] M. Baes, R. Govers, and R. Wortel. Switching between alternative responses of the lithosphere to continental collision. *Geophys. J. Int.*, 2011.
- [31] B. Baitsch-Ghirardello, Taras V. Gerya, and J.-P. Burg. Geodynamic regimes of intra-oceanic subduction: Implications for arc extension vs. shortening processes. *Gondwana Research*, 25:546–560, 2014.
- [32] M.D. Ballmer, G. Ito, J. van Hunen, and P.J. Tackley. Small-scale sublithospheric convection reconciles geochemistry and geochronology of 'Superplume' volcanism in th western and south pacific. *Earth Planet. Sci. Lett.*, 290:224–232, 2010.
- [33] M.D. Ballmer, J. van Hunen, G. Ito, P.J. Tackley, and T.A. Bianco. Non-hotspot volcano chains originating from small-scale sublithospheric convection. *Geophys. Res. Lett.*, 34(L23310):doi:10.1029/2007GL031636, 2007.
- [34] W. Bangerth, R. Hartmann, and G. Kanschat. deal.II - a general purpose object oriented finite element library. *ACM Transaction on mathematical software*, 33(4), 2007.
- [35] R. Barrett, M. Berry, T.F. Chan, J. Demmel, J.M. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. Van der Vorst. *Templates for the solution of linear systems: building blocks for iterative methods*. SIAM, 1994.
- [36] G.K. Batchelor. *An introduction to fluid dynamics*. Cambridge University Press, 1967.
- [37] K.-J. Bathe. *Finite Element Procedures in Engineering Analysis*. Prentice-Hall, 1982.
- [38] L. Battaglia, M.A. Storti, and J. D’Elia. An interface capturing finite element approach for free surface flows using unstructured grids. *Mecanica Computational*, XXVII:33–48, 2008.
- [39] A. Bauville and T. S. Baumann. geomio: an open-source matlab toolbox to create the initial configuration of 2d/3d thermo-mechanical simulations from 2d vector drawings. *Geochemistry, Geophysics, Geosystems*, 2019.

- [40] C. Beaumont, P. Fullsack, and J. Hamilton. Styles of crustal deformation in compressional orogens caused by subduction of the underlying lithosphere. *Tectonophysics*, 232:119–132, 1994.
- [41] C. Beaumont, R.A. Jamieson, J.P. Butler, and C.J. Warren. Crustal structure: A key constraint on the mechanism of ultra-high-pressure rock exhumation. *Earth Planet. Sci. Lett.*, 287:116–129, 2009.
- [42] C. Beaumont, R.A. Jamieson, M.H. Nguyen, and B. Lee. Himalayan tectonics explained by extrusion of a low-viscosity crustal channel coupled to focused surface denudation. *Nature*, 414:738–742, 2001.
- [43] C. Beaumont, R.A. Jamieson, M.H. Nguyen, and S. Medvedev. Crustal channel flows: 1. Numerical models with applications to the tectonics of the Himalayan-Tibetan orogen. *J. Geophys. Res.*, 109(B06406), 2004.
- [44] C. Beaumont, P.J. Kamp, J. Hamilton, and P. Fullsack. The continental collision zone, south island, new zealand: comparison of geodynamical models and observations. *J. Geophys. Res.*, 101:3333–3359, 1996.
- [45] 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. *J. Geophys. Res.*, 105:8121–8145, 2000.
- [46] T.W. Becker. On the effect of temperature and strain-rate dependent viscosity on global mantle flow, net rotation, and plate-driving forces. *Geophys. J. Int.*, 167:943–957, 2006.
- [47] T.W. Becker and C. Faccenna. Mantle conveyor beneath the Tethyan collisional belt. *Earth Planet. Sci. Lett.*, 310:453–461, 2011.
- [48] T.W. Becker, C. Faccenna, R. O’Connell, and D. Giardini. The development of slabs in the upper mantle: Insights from numerical and laboratory experiments. *J. Geophys. Res.*, 104(B7):15,207–15,226, 1999.
- [49] T.W. Becker, V. Schulte-Pelkum, D.K. Blackman, J.B. Kellogg, and R.J. O’Connell. Mantle flow under the western United States from shear wave splitting. *Earth Planet. Sci. Lett.*, 247:235–251, 2006.
- [50] A.K. Bengtson and P.E. van Keken. Three-dimensional thermal structure of subduction zones: effects of obliquity and curvature. *Solid Earth*, 3:365–373, 2012.
- [51] M. Benzi, G.H. Golub, and J. Liesen. Numerical solution of saddle point problems. *Acta Numerica*, 14:1–137, 2005.
- [52] D. Bercovici, G. Schubert, and G.A. Glatzmaier. Three-dimensional convection of an infinite Prandtl-number compressible fluid in a basally heated spherical shell. *J. Fluid Mech.*, 239:683–719, 1992.
- [53] M. Bercovier and M. Engelman. A finite-element for the numerical solution of viscous incompressible flows. *J. Comp. Phys.*, 30:181–201, 1979.
- [54] M. Bercovier and M. Engelman. A finite-element method for incompressible Non-Newtonian flows. *J. Comp. Phys.*, 36:313–326, 1980.
- [55] P.G. Betts, W.G. Mason, and L. Moresi. The influence of a mantle plume head on the dynamics of a retreating subduction zone. *Geology*, 40(8):739–742, 2012.
- [56] P.G. Betts, L. Moresi, M.S. Miller, and D. Willis. Geodynamics of oceanic plateau and plume head accretion and their role in Phanerozoic orogenic systems of China. *Geoscience Frontiers*, 6:49–59, 2015.
- [57] M.J. Beuchert and Y.Y. Podladchikov. Viscoelastic mantle convection and lithospheric stresses. *GJI*, 183:35–63, 2010.

- [58] E. Beutel, J. van Wijk, C. Ebinger, D. Keir, and A. Agostini. Formation and stability of magmatic segments in the Main Ethiopian and Afar rifts. *Earth Planet. Sci. Lett.*, 293:225–235, 2010.
- [59] M.I. Billen and M. Gurnis. A low wedge in subduction zones. *Earth Planet. Sci. Lett.*, 193:227–236, 2001.
- [60] M.I. Billen and M. Gurnis. Comparison of dynamic flow models for the Central Aleutian and Tonga-Kermadec subduction zones. *Geochem. Geophys. Geosyst.*, 4(4), 2003.
- [61] M.I. Billen and G. Hirth. Newtonian versus non-Newtonian upper mantle viscosity: Implications for subduction initiation. *Geophys. Res. Lett.*, 32(L19304):doi:10.1029/2005GL023457, 2005.
- [62] M.I. Billen and G. Hirth. Rheologic controls on slab dynamics. *Geochem. Geophys. Geosyst.*, 8(8):doi:10.1029/2007GC001597, 2007.
- [63] M.I. Billen and M. Jadamec. Origin of localized fast mantle f low velocity in numerical models of subduction. *Geochem. Geophys. Geosyst.*, 13(1):doi:10.1029/2011GC003856, 2012.
- [64] B. Blankenbach, F. Busse, U. Christensen, L. Cserepes, D. Gunkel, U. Hansen, H. Harder, G. Jarvis, M. Koch, G. Marquart, D. Moore, P. Olson, H. Schmeling, and T. Schnaubelt. A benchmark comparison for mantle convection codes. *Geophys. J. Int.*, 98:23–38, 1989.
- [65] P. B. Bochev, C. R. Dohrmann, and M. D. Gunzburger. Stabilization of low-order mixed finite elements for the stokes equations. *SIAM Journal on Numerical Analysis*, 44(1):82–101, 2006.
- [66] M.-A. Bonnardot, R. Hassani, and E. Tric. Numerical modelling of lithosphereasthenosphere interaction in a subduction zone. *Earth Planet. Sci. Lett.*, 272:698–708, 2008.
- [67] M.-A. Bonnardot, R. Hassani, E. Tric, E. Ruellan, and M. Regnier. Effect of margin curvature on plate deformation in a 3-D numerical model of subduction zones. *Geophy. J. Int.*, 173:1084–1094, 2008.
- [68] O. Botella and R. Peyret. Benchmark spectral results on the lid-driven cavity flow. *Computers and Fluids*, 27(4):421–433, 1998.
- [69] A.D. Bottrill, J. van Hunen, and M.B. Allen. Insight into collision zone dynamics from topography: numerical modelling results and observations. *Solid Earth*, 3:387–399, 2012.
- [70] P. Bouilhol, V. Magni, J. van Hunen, and L. Kaislaniemi. A numerical approach to melting in warm subduction zones. *Earth Planet. Sci. Lett.*, 411:37–44, 2015.
- [71] L. Bourgouin, H.-B. Mühlhaus, A.J. Hale, and A. Arsac. Towards realistic simulations of lava dome growth using the level set method. *Acta Geotechnica*, 1:225–236, 2006.
- [72] 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. *Geophy. J. Int.*, 170:1431–1438, 2007.
- [73] D.J. Bower, M. Gurnis, and N. Flament. Assimilating lithosphere and slab history in 4-D Earth models. *Phys. Earth. Planet. Inter.*, 238:8–22, 2015.
- [74] D.J. Bower, M. Gurnis, and M. Seton. Lower mantle structure from paleogeographically constrained dynamic Earth models. *Geochem. Geophys. Geosyst.*, 14(1):44–63, 2012.
- [75] D.J. Bower, M. Gurnis, and D. Sun. Dynamic origins of seismic wavespeed variation in D”. *Phys. Earth. Planet. Inter.*, 214:74–86, 2013.
- [76] D. Braess. *Finite Elements*. Cambridge, 2007.
- [77] J.P. Brandenburg and P.E. van Keken. Deep storage of oceanic crust in a vigorously convecting mantle. *J. Geophys. Res.*, 112(B06403), 2007.
- [78] J.P. Brandenburg and P.E. van Keken. Methods for thermochemical convection in Earths mantle with force-balanced plates. *Geochem. Geophys. Geosyst.*, 8(11), 2007.

- [79] 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. *Computers and Geosciences*, 29:787–794, 2003.
- [80] 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. *Phys. Earth. Planet. Inter.*, 171:76–91, 2008.
- [81] J. Braun and P. Yamato. Structural evolution of a three-dimensional, finite-width crustal wedge. *Tectonophysics*, 484:181–192, 2009.
- [82] E. Bredow and B. Steinberger. Variable melt production rate of the kerguelen hotspot due to long-term plume-ridge interaction. *Geophysical Research Letters*, 45(1):126–136, 2018.
- [83] E. Bredow, B. Steinberger, R. Gassmöller, and J. Dannberg. How plume-ridge interaction shapes the crustal thickness pattern of the réunion hotspot track. *Geochemistry, Geophysics, Geosystems*, 2017.
- [84] 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. *Computer Methods in Applied Mechanics and Engineering*, 32:199–259, 1982.
- [85] S. Brune. Evolution of stress and fault patterns in oblique rift systems: 3-D numerical lithospheric-scale experiments from rift to breakup. *Geochem. Geophys. Geosyst.*, 15:3392–3415, 2014.
- [86] S. Brune and J. Autin. The rift to break-up evolution of the Gulf of Aden: Insights from 3D numerical lithospheric-scale modelling. *Tectonophysics*, 607:65–79, 2013.
- [87] S. Brune, A.A. Popov, and S. Sobolev. Modeling suggests that oblique extension facilitates rifting and continental break-up. *J. Geophys. Res.*, 117(B08402), 2012.
- [88] S. Brune, A.A. Popov, and S. Sobolev. Quantifying the thermo-mechanical impact of plume arrival on continental break-up. *Tectonophysics*, 604:51–59, 2013.
- [89] C.-H. Bruneau and M. Saad. The 2D lid-driven cavity problem revisited. *Computers & Fluids*, 35:326–348, 2006.
- [90] H.H. Bui, R. Fukugawa, K. Sako, and S. Ohno. Lagrangian meshfree particles method (SPH) for large deformation and failure flows of geomaterial using elasticplastic soil constitutive model. *Int. J. Numer. Anal. Geomech.*, 32(12):1537–1570, 2008.
- [91] S. Buiter, A.Y. Babeyko, S. Ellis, T.V. Gerya, B.J.P. Kaus, A. Kellner, G. Schreurs, and Y. Yamada. The numerical sandbox: comparison of model results for a shortening and an extension experiment. *Analogue and Numerical Modelling of Crustal-Scale Processes. Geological Society, London. Special Publications*, 253:29–64, 2006.
- [92] S.J.H. Buiter. A review of brittle compressional wedge models. *Tectonophysics*, 530:1–17, 2012.
- [93] S.J.H. Buiter, R. Govers, and M.J.R. Wortel. A modelling study of vertical surface displacements at convergent plate margins. *Geophys. J. Int.*, 147:415–427, 2001.
- [94] S.J.H. Buiter, R. Govers, and M.J.R. Wortel. Two-dimensional simulations of surface deformation caused by slab detachment. *Tectonophysics*, 354:195–210, 2002.
- [95] S.J.H. Buiter, O.A. Pfiffner, and C. Beaumont. Inversion of extensional sedimentary basins: A numerical evaluation of the localisation of shortening. *Earth Planet. Sci. Lett.*, 288:492–504, 2009.
- [96] S.J.H. Buiter, G. Schreurs, M. Albertz, T.V. Gerya, B. Kaus, W. Landry, L. le Pourhiet, Y. Mishin, D.L. Egholm, M. Cooke, B. Maillot, C. Thieulot, T. Crook, D. May, P. Souloumiac, and C. Beaumont. Benchmarking numerical models of brittle thrust wedges. *Journal of Structural Geology*, 92:140–177, 2016.

- [97] A.L. Bull, M. Domeier, and T.H. Torsvik. The effect of plate motion history on the longevity of deep mantle heterogeneities. *Earth Planet. Sci. Lett.*, 401:172–182, 2014.
- [98] 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. *Phys. Earth. Planet. Inter.*, 182:129–138, 2010.
- [99] P.S. Bullen. *Handbook of Means and Their Inequalities*. Springer; 2nd edition, 2003.
- [100] 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. *Science*, 280:91–95, 1998.
- [101] 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 . *J. Geophys. Res.*, 102(B6):11,991–12,007, 1997.
- [102] J.-P. Burg and T.V. Gerya. The role of viscous heating in Barrovian metamorphism of collisional orogens: thermomechanical models and application to the Lepontine Dome in the Central Alps. *J. Metamorphic Geology*, 23:75–95, 2005.
- [103] E.R. Burkett and M.I. Billen. Dynamics and implications of slab detachment due to ridge-trench collision. *J. Geophys. Res.*, 114(B12402), 2009.
- [104] E.R. Burkett and M.I. Billen. Three-dimensionality of slab detachment due to ridge-trench collision: Laterally simultaneous boudinage versus tear propagation. *Geochem. Geophys. Geosyst.*, 11(11), 2010.
- [105] 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. *J. Geophys. Res.*, 106(B8):16,461–16,481, 2001.
- [106] E. Burov, T. Francois, P. Agard, L. Le Pourhiet, B. Meyer, C. Tirel, S. Lebedev, P. Yamato, and J.-P. Brun. Rheological and geodynamic controls on the mechanisms of subduction and HP/UHP exhumation of crustal rocks during continental collision: Insights from numerical models. *Tectonophysics*, 2014.
- [107] E. Burov and T. Gerya. Asymmetric three-dimensional topography over mantle plumes. *Nature*, 513:doi:10.1038/nature13703, 2014.
- [108] E. Burov and L. Guillou-Frottier. The plume head-continental lithosphere interaction using a tectonically realistic formulation for the lithosphere. *Geophys. J. Int.*, 161:469–490, 2005.
- [109] E. Burov, L. Jolivet, L. Le Pourhiet, and A. Poliakov. A thermomechanical model of exhumation of high pressure (HP) and ultra-high pressure (UHP) metamorphic rocks in Alpine-type collision belts. *Tectonophysics*, 342:113–136, 2001.
- [110] E. Burov and G. Toussaint. Surface processes and tectonics: Forcing of continental subduction and deep processes. *Global and Planetary Change*, 58:141–164, 2007.
- [111] C. Burstedde, O. Ghattas, M. Gurnis, G. Stadler, E. Tan, T. Tu, L.C. Wilcox, and S. Zhong. Scalable Adaptive Mantle Convection Simulation on Petascale Supercomputers. *ACM/IEEE SC Conference Series*, 2008, 2008.
- [112] C. Burstedde, O. Ghattas, G. Stadler, T. Tu, and L.C. Wilcox. Parallel scalable adjoint-based adaptive solution of variable-viscosity Stokes flow problems. *Computer Methods in Applied Mechanics and Engineering*, 198:1691–1700, 2009.
- [113] C. Burstedde, G. Stadler, L. Alisic, L.C. Wilcox, E. Tan, M. Gurnis, and O. Ghattas. Large-scale adaptive mantle convection simulation. *Geophys. J. Int.*, 192:889–906, 2013.

- [114] F.H. Busse, U. Christensen, R. Clever, L. Cserepes, C. Gable, E. Giannandrea, L. Guillou, G. Houseman, H.-C. Nataf, M. Ogawa, M. Parmentier, C. Sotin, and B. Travis. 3D convection at infinite Prandtl number in Cartesian geometry - a benchmark comparison. *Geophys. Astrophys. Fluid Dynamics*, 75:39–59, 1993.
- [115] J.P. Butler and C. Beaumont. Subduction zone decoupling/retreat modeling explains south Tibet (Xigaze) and other supra-subduction zone ophiolites and their UHP mineral phases. *Earth Planet. Sci. Lett.*, 463:101–117, 2017.
- [116] J.P. Butler, C. Beaumont, and R.A. Jamieson. The Alps 1: A working geodynamic model for burial and exhumation of (ultra)high-pressure rocks in Alpine-type orogens. *Earth Planet. Sci. Lett.*, 337–378:114–131, 2013.
- [117] J.P. Butler, C. Beaumont, and R.A. Jamieson. Paradigm lost: Buoyancy thwarted by the strength of the Western Gneiss Region (ultra)high-pressure terrane, Norway. *Lithosphere*, page doi:10.1130/L426.1, 2015.
- [118] J.P. Butler, C. Beaumont, and R.A. Jamieson. Crustal emplacement of exhuming (ultra)high-pressure rocks: Will that be pro- or retro-side? *Geology*, 39:635–638, 2011.
- [119] F.A. Capitanio and M. Faccenda. Complex mantle flow around heterogeneous subducting oceanic plates. *Earth Planet. Sci. Lett.*, 353–354:29–37, 2012.
- [120] F.A. Capitanio, C. Faccenna, S. Zlotnik, and D.R. Stegman. Subduction dynamics and the origin of Andean orogeny and the Bolivian orocline. *Nature*, 480:doi:10.1038/nature10596, 2011.
- [121] F.A. Capitanio and A. Replumaz. Subduction and slab breakoff controls on Asian indentation tectonics and Himalayan western syntaxis formation. *Geochem. Geophys. Geosyst.*, 14(9):doi:10.1002/ggge.20171, 2013.
- [122] F.A. Capitanio, D.R. Stegman, L.N. Moresi, and W. Sharples. Upper plate controls on deep subduction, trench migrations and deformations at convergent margins. *Tectonophysics*, 483:80–92, 2010.
- [123] N.G. Cerpa, R. Araya, M. Gerbault, and R. Hassani. Relationship between slab dip and topography segmentation in an oblique subduction zone: Insights from numerical modeling. *Geophys. Res. Lett.*, 41:10.1002/2015GL064047, 2015.
- [124] N.G. Cerpa, B. Guillaume, and J. Martinod. The interplay between overriding plate kinematics, slab dip and tectonics. *Geophys. J. Int.*, 215:1789–1802, 2018.
- [125] J.S. Chen, C. Pan, and T.Y.P. Chang. On the control of pressure oscillation in bilinear-displacement constant-pressure element. *Comput. Methods Appl. Mech. Engrg.*, 128:137–152, 1995.
- [126] P. Chenin and C. Beaumont. Influence of offset weak zones on the development of rift basins: Activation and abandonment during continental extension and breakup. *J. Geophys. Res.*, 118:1–23, 2013.
- [127] M.V. Chertova, T. Geenen, A. van den Berg, and W. Spakman. Using open sidewalls for modelling self-consistent lithosphere subduction dynamics. *Solid Earth*, 3:313–326, 2012.
- [128] M.V. Chertova, W. Spakman, T. Geenen, A.P. van den Berg, and D.J.J. van Hinsbergen. Underpinning tectonic reconstructions of the western Mediterranean region with dynamic slab evolution from 3-D numerical modeling. *J. Geophys. Res.*, 119:10.1002/2014JB011150, 2014.
- [129] S. Chiw Webster, E.J. Hinch, and J.R. Lister. Very viscous horizontal convection. *J. Fluid Mech.*, 611:395–426, 2008.
- [130] E. Choi and K.D. Petersen. Making Coulomb angle-oriented shear bands in numerical tectonic models. *Tectonophysics*, 657:94–101, 2015.

- [131] E. Choi, E. Tan, L.L. Lavier, and V.M. Calo. DynEarthSol2D: An efficient unstructured finite element method to study long-term tectonic deformation. *J. Geophys. Res.*, 118:1–16, 2013.
- [132] Eunseo Choi, Luc Lavier, and Michael Gurnis. Thermomechanics of mid-ocean ridge segmentation. *Phys. Earth Planet. Interiors*, 171:374–386, 2008.
- [133] Edmund Christiansen and Knud D. Andersen. Computation of collapse states with von mises type yield condition. *International Journal for Numerical Methods in Engineering*, 46:1185–1202, 1999.
- [134] H. Ciskova, J. van Hunen, A.P. van den Berg, and N.J. Vlaar. The influence of rheological weakening and yield stress on the interaction of slabs with the 670 km discontinuity. *Earth Planet. Sci. Lett.*, 199:447–457, 2002.
- [135] P.D. Clift, S. Brune, and J. Quinteros. Climate changes control offshore crustal structure at South China Sea continental margin. *Earth Planet. Sci. Lett.*, 420:66–72, 2015.
- [136] S. Cloetingh, E. Burov, F. Beekman, B. Andeweg, P.A.M. Andriessen, D. Garcia-Castellanos, G. de Vicente, and R. Vegas. Lithospheric folding in Iberia. *Tectonics*, 21(5):10.1029/2001TC901031, 2002.
- [137] M. Collignon, N. Fernandez, and B.J.P. Kaus. Influence of surface processes and initial topography on lateral fold growth and fold linkage mode. *Tectonics*, 34:1622–1645, 2015.
- [138] M. Collignon, B.J.P. Kaus, D.A. May, and N. fernandez. Influences of surface processes on fold growth during 3-D detachment folding. *Geochem. Geophys. Geosyst.*, 15:doi:10.1002/2014GC005450, 2014.
- [139] C.P. Conrad, M.D. Behn, and P.G. Silver. Global mantle flow and the development of seismic anisotropy: Differences between the oceanic and continental upper mantle. *J. Geophys. Res.*, 112(B07317), 2007.
- [140] C.P. Conrad and M. Gurnis. Seismic tomography, surface uplift, and the breakup of Gondwanaland: Integrating mantle convection backwards in time. *Geochem. Geophys. Geosyst.*, 4(3), 2003.
- [141] C.P. Conrad and C. Lithgow-Bertelloni. Influence of continental roots and asthenosphere on plate-mantle coupling. *Geophys. Res. Lett.*, 33(L05312), 2006.
- [142] C.P. Conrad, B. Steinberger, and T.H. Torsvik. Stability of active mantle upwelling revealed by net characteristics of plate tectonics. *Nature*, 498:479, 2013.
- [143] G. Corti, R. Cioni, Z. Franceschini, F. Sani, Stéphane Scaillet, P. Molin, I. Isola, F. Mazzarini, S. Brune, D. Keir, A. Erbello, A. Muluneh, F. Illsley-Kemp, and A. Glerum. Aborted propagation of the ethiopian rift caused by linkage with the kenyan rift. *Nature Communications*, 10, 2019.
- [144] F. Crameri, H. Schmeling, G.J. Golabek, T. Duretz, R. Orendt, S.J.H. Buiter, D.A. May, B.J.P. Kaus, T.V. Gerya, and P.J. Tackley. A comparison of numerical surface topography calculations in geodynamic modelling: an evaluation of the 'sticky air' method. *Geophy. J. Int.*, 189:38–54, 2012.
- [145] F. Crameri and P.J. Tackley. Spontaneous development of arcuate single-sided subduction in global 3-D mantle convection models with a free surface. *J. Geophys. Res.*, 119:doi:10.1002/2014JB010939, 2014.
- [146] M. Crouzeix and P.-A. Raviart. Conforming and nonconforming finite element methods for solving the stationary Stokes equations I. *R.A.I.R.O.*, 7(3):33–75, 1973.
- [147] L. Cruz, J. Malinski, A. Wilson, W.A. Take, and G. Hilley. Erosional control of kinematics and geometry of fold-and-thrust belts imaged in a physical and numerical sandbox. *J. Geophys. Res.*, 115(B09404):doi:10.1029/2010JB007472, 2010.
- [148] C.A. Currie and C. Beaumont. Are diamond-bearing Cretaceous kimberlites related to low-angle subduction beneath western North America. *Earth Planet. Sci. Lett.*, 303:59–70, 2011.

- [149] C.A. Currie, C. Beaumont, and R.S. Huismans. The fate of subducted sediments: a case for backarc intrusion and underplating. *Geology*, 35(12):1111–1114, 2007.
- [150] C. Cuvelier, A. Segal, and A.A. van Steenhoven. *Finite Element Methods and Navier-Stokes Equations*. D. Reidel Publishing Company, 1986.
- [151] M. Dabrowski, M. Krotkiewski, and D.W. Schmid. Milamin: Matlab based finite element solver for large problems. *Geochem. Geophys. Geosyst.*, 9(4):Q04030, 2008.
- [152] J. Dannberg, Z. Eilon, U. Faul, R. Gassmoeller, P. Moulik, and R. Myhill. The importance of grain size to mantle dynamics and seismological observations. *Geochem. Geophys. Geosyst.*, 18:3034–3061, 2017.
- [153] J. Dannberg and R. Gassmöller. Chemical trends in ocean islands explained by plume-slab interaction. *PNAS*, 115(17):4351–4356, 2018.
- [154] J. Dannberg and T. Heister. Compressible magma/mantle dynamics: 3-D, adaptive simulations in ASPECT. *Geophy. J. Int.*, 207:1343–1366, 2016.
- [155] J. Dannberg and S.V. Sobolev. Low-buoyancy thermochemical plumes resolve controversy of classical mantle plume concept. *Nature Communications*, 6(6960):doi:10.1038/ncomms7960, 2015.
- [156] 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. *Geosci. Model Dev.*, 6:1095–1107, 2013.
- [157] D.R. Davies, C.R. Wilson, and S.C. Kramer. Fluidity: A fully unstructured anisotropic adaptive mesh computational modeling framework for geodynamics. *Geochem. Geophys. Geosyst.*, 12(6), 2011.
- [158] P. Davy and P. Cobbold. Indentation tectonics in nature and experiment. 1. experiments scaled for gravity. *Bulletin of the Geological Institutions of Uppsala*, 14:129–141, 1988.
- [159] B. Deglo de Besses, A. Magnin, and P. Jay. Sphere drag in a viscoplastic fluid. *AIChE Journal*, 50(10):2627–2629, 2004.
- [160] J. de Frutos, V. John, and J. Novo. Projection methods for incompressible flow problems with WENO finite difference schemes. *J. Comp. Phys.*, 309:368–386, 2016.
- [161] C.S. Desai and J.F. Abel. *Introduction to the Finite Element Method: A Numerical Method for Engineering Analysis*. Van Nostrand Reinhold, 1972.
- [162] 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. *Phys. Earth Planet. Interiors*, 171:92–111, 2008.
- [163] A.E. Svartman Dias, L.L. Lavie, and N.W. Hayman. Conjugate rifted margins width and asymmetry: The interplay between lithospheric strength and thermomechanical processes. *J. Geophys. Res.*, 120:8672–8700, 2015.
- [164] C.R. Dohrmann and P.B. Bochev. A stabilized finite element method for the Stokes problem based on polynomial pressure projections. *Int. J. Num. Meth. Fluids*, 46:183–201, 2004.
- [165] Jean Donea and Antonio Huerta. *Finite Element Methods for Flow Problems*. John Wiley & Sons, 2003.
- [166] D.A. Dunavant. High-degree efficient symmetrical Gaussian quadrature rules for the triangle. *Int. J. Num. Meth. Eng.*, 21:1129–1148, 1985.
- [167] T. Duretz, Ph. Agard, Ph. Yamato, C. Ducassou, E.B. Burov, and T.V. Gerya. Thermo-mechanical modeling of the obduction process based on the Oman Ophiolite case. *Gondwana Research*, 2015.

- [168] T. Duretz, T.V. Gerya, B.J.P. Kaus, and T.B. Andersen. Thermomechanical modeling of slab exhumation. *J. Geophys. Res.*, 117(B08411), 2012.
- [169] T. Duretz, T.V. Gerya, and D.A. May. Numerical modelling of spontaneous slab breakoff and subsequent topographic response. *Tectonophysics*, 502:244–256, 2011.
- [170] T. Duretz, T.V. Gerya, and W. Spakman. Slab detachment in laterally varying subduction zones: 3-D numerical modeling. *Geophys. Res. Lett.*, 41:1951–1956, 2014.
- [171] 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. *Geochem. Geophys. Geosyst.*, 12(Q07004), 2011.
- [172] S. Dyksterhuis, P. Rey, R.D. Mueller, and L. Moresi. Effects of initial weakness on rift architecture. *Geological Society, London, Special Publications*, 282:443–455, 2007.
- [173] D. Dymkova and T. Gerya. Porous fluid flow enables oceanic subduction initiation on Earth. *Geophys. Res. Lett.*, 2013.
- [174] David L. Egholm. A new strategy for discrete element numerical models: 1. Theory. *J. Geophys. Res.*, 112:B05203, doi:10.1029/2006JB004557, 2007.
- [175] David L. Egholm, Mike Sandiford, Ole R. Clausen, and Søren B. Nielsen. A new strategy for discrete element numerical models: 2. Sandbox applications. *J. Geophys. Res.*, 112:B05204, doi:10.1029/2006JB004558, 2007.
- [176] R. Eid. Higher order isoparametric finite element solution of Stokes flow . *Applied Mathematics and Computation*, 162:1083–1101, 2005.
- [177] V. Eijkhout. *Introduction to High Performance Scientific Computing*. Creative Commons, 2013.
- [178] 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. *Arabian Journal of Geosciences*, 3:477–497, 2010.
- [179] 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. *Geophys. J. Int.*, 120:24–44, 1995.
- [180] S. Ellis, G. Schreurs, and M. Panien. Comparisons between analogue and numerical models of thrust wedge development. *Journal of Structural Geology*, 26:1649–1675, 2004.
- [181] S.M. Ellis, T.A. Little, L.M. Wallace, B.R. Hacker, and S.J.H. Buiter. Feedback between rifting and diapirism can exhume ultrahigh-pressure rocks. *Earth Planet. Sci. Lett.*, 311:427–438, 2011.
- [182] H. Elman, D. Silvester, and A. Wathen. *Finite Elements and Fast Iterative Solvers*. Oxford Science Publications, 2014.
- [183] M.S. Engelman and R.L. Sani. The implementation of normal and/or tangential boundary conditions in finite element codes for incompressible fluid flow. *Int. J. Num. Meth. Fluids*, 2:225–238, 1982.
- [184] A. Enns, T.W. Becker, and H. Schmeling. The dynamics of subduction and trench migration for viscosity stratification. *Geophys. J. Int.*, 160:761–775, 2005.
- [185] Z. Erdos, R.S. huismans, and P. van der Beek. First-order control of syntectonic sedimentation on crustal-scale structure of mountain belts . *J. Geophys. Res.*, 120:doi:10.1002/2014JB011785, 2015.
- [186] Z. Erdos, R.S. Huismans, and P. van der Beek. Control of increased sedimentation on orogenic fold-and-thrust belt structure - insights into the evolution of the Western Alps. *Solid Earth*, 10:391–404, 2019.

- [187] Z. Erdos, R.S. huismans, P. van der Beek, and C. Thieulot. Extensional inheritance and surface processes as controlling factors of mountain belt structure. *J. Geophys. Res.*, 119:doi:10.1002/2014JB011408, 2014.
- [188] E. Erturk. Discussions on Driven Cavity Flow. *Int. J. Num. Meth. Fluids*, 60:275–294, 2009.
- [189] M. Faccenda. Mid mantle seismic anisotropy around subduction zones. *Phys. Earth. Planet. Inter.*, 227:1–19, 2014.
- [190] M. Faccenda and F.A. Capitanio. Seismic anisotropy around subduction zones: Insights from three-dimensional modeling of upper mantle deformation and SKS splitting calculations. *Geochem. Geophys. Geosyst.*, 14(1):doi:10.1029/2012GC004451, 2013.
- [191] M. Faccenda, T.V. Gerya, and S. Chakraborty. Styles of post-subduction collisional orogeny: Influence of convergence velocity, crustal rheology and radiogenic heat production. *Lithos*, 103:257–287, 2008.
- [192] 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. *Geochem. Geophys. Geosyst.*, 13(1):doi:10.1029/2011GC003860, 2012.
- [193] R.J. Farrington, L.-N. Moresi, and F.A. Capitanio. The role of viscoelasticity in subducting plates. *Geochem. Geophys. Geosyst.*, 15:4291–4304, 2014.
- [194] R.J. Farrington, D.R. Stegman, L.N. Moresi, M. Sandiford, and D.A. May. Interactions of 3D mantle flow and continental lithosphere near passive margins. *Tectonophysics*, 483:20–28, 2010.
- [195] N.P. Fay, R.A. Bennett, J.C. Spinler, and E.D. Humphreys. Small-scale upper mantle convection and crustal dynamics in southern California. *Geochem. Geophys. Geosyst.*, 9(8), 2008.
- [196] N. Fernandez and B. Kaus. Influence of pre-existing salt diapirs on 3D folding patterns. *Tectonophysics*, 637:354–369, 2014.
- [197] N. Fernandez and B. Kaus. Fold interaction and wavelength selection in 3D models of multilayer detachment folding. *Tectonophysics*, 632:199–217, 2014.
- [198] N. Fernandez and B. Kaus. Pattern formation in 3-D numerical models of down-built diapirs initiated by a RayleighTaylor instability. *Geophy. J. Int.*, 202:1253–1270, 2015.
- [199] C. Fillon, R.S. Huismans, and P. van der Beek. Syntectonic sedimentation effects on the growth of fold-and-thrust belts. *Geology*, 41(1):83–86, 2013.
- [200] C. Fillon, R.S. Huismans, P. van der Beek, and J.A. Mu noz. Syntectonic sedimentation controls on the evolution of the southern Pyrenean fold-and-thrust belt: Inferences from coupled tectonic-surface processes models. *J. Geophys. Res.*, 118:5665–5680, 2013.
- [201] R. Fischer and T. Gerya. Early Earth plume-lid tectonics: A high-resolution 3D numerical modelling approach. *Journal of Geodynamics*, 100:198–214, 2016.
- [202] N. Flament, M. Gurnis, S. Williams, M. Seton, J. Skogseid, C. Heine, and D. Müller. Topographic asymmetry of the South Atlantic from global models of mantle flow and lithospheric stretching. *Earth Planet. Sci. Lett.*, 387:107–119, 2014.
- [203] B.J. Foley and T.W. Becker. Generation of plate-like behavior and mantle heterogeneity from a spherical, viscoplastic convection model. *Geochem. Geophys. Geosyst.*, 10(8):doi:10.1029/2009GC002378, 2009.
- [204] M. Fortin. Old and new finite elements for incompressible flows. *Int. J. Num. Meth. Fluids*, 1:347–364, 1981.
- [205] M. Fortin and A. Fortin. Experiments with several elements for viscous incompressible flows. *Int. J. Num. Meth. Fluids*, 5:911–928, 1985.

- [206] L.P. Franca and S.P. Oliveira. Pressure bubbles stabilization features in the Stokes problem. *Computer Methods in Applied Mechanics and Engineering*, 192:1929–1937, 2003.
- [207] R. De Franco, R. Govers, and R. Wortel. Numerical comparison of different convergent plate contacts: subduction channel and subduction fault. *Geophys. J. Int.*, pages doi: 10.1111/j.1365–246X.2006.03498.x, 2006.
- [208] T. Francois, E. Burov, P. Agard, and B. Meyer. Buildup of a dynamically supported orogenic plateau: Numerical modeling of the Zagros/Central Iran case study. *Geochem. Geophys. Geosyst.*, 15:doi:10.1002/ 2013GC005223, 2014.
- [209] 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. *Geophys. J. Int.*, 2019.
- [210] R. Freeburn, P. Bouilhol, B. Maunder, V. Magni, and J. van Hunen. Numerical models of the magmatic processes induced by slab breakoff. *Earth Planet. Sci. Lett.*, 478:203–213, 2017.
- [211] P.J. Frey and P.-L. George. *Mesh generation*. Hermes Science, 2000.
- [212] L. Fuchs and Th.W. Becker. Role of strain-dependent weakening memory on the style of mantle convection and plate boundary stability. *Geophys. J. Int.*, 218:601–618, 2019.
- [213] L. Fuchs, H. Koyi, and H. Schmeling. Numerical modeling of the effect of composite rheology on internal deformation in down-built diapirs. *Tectonophysics*, 646:79–95, 2015.
- [214] 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. *Geophys. J. Int.*, 194(2):623–639, 2013.
- [215] J. Fullea, J.C. Afonso, J.A.D. Connolly, M. Fernandez, D. Garcia-Castellanos, and H. Zeyen. LitMod3D: An interactive 3-D software to model the thermal, compositional, density, seismological, and rheological structure of the lithosphere and sublithospheric upper mantle. *Geochem. Geophys. Geosyst.*, 10(8):doi:10.1029/2009GC002391, 2009.
- [216] J. Fullea, M. Fernandez, J.C. Afonso, J. Verges, and H. Zeyen. iThe structure and evolution of the lithosphereasthenosphere boundary beneath the AtlanticMediterranean Transition Region. *Lithos*, 2010.
- [217] Ch.W. Fuller, S.D. Willett, and M.T. Brandon. Formation of forearc basins and their influence on subduction zone earthquakes. *Geology*, 34(2):65–68, 2006.
- [218] P. Fullsack. An arbitrary Lagrangian-Eulerian formulation for creeping flows and its application in tectonic models. *Geophys. J. Int.*, 120:1–23, 1995.
- [219] M. Furuichi and D. Nishiura. Robust coupled fluid-particle simulation scheme in Stokes-flow regime: Toward the geodynamic simulation including granular media. *Geochem. Geophys. Geosyst.*, 15:2865–2882, 2014.
- [220] J. Ganne, M. Gerbault, and S. Block. Thermo-mechanical modeling of lower crust exhumationConstraints from the metamorphic record of the Palaeoproterozoic Eburnean orogeny, West African Craton. *Precambrian Research*, 243:88–109, 2014.
- [221] 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. *Geochem. Geophys. Geosyst.*, 15(1739–1765):doi:10.1002/2014GC005257, 2014.
- [222] E.J. Garnero and A.K. McNamara. Structure and Dynamics of Earths Lower Mantle. *Science*, 320:626–628, 2008.
- [223] R. Gassmöller, J. Dannberg, E. Bredow, B. Steinberger, and T. H. Torsvik. Major influence of plume-ridge interaction, lithosphere thickness variations, and global mantle flow on hotspot volcanism-the example of tristan. *Geochem. Geophys. Geosyst.*, 17(4):1454–1479, 2016.

- [224] R. Gassmöller, 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. *Geochem. Geophys. Geosyst.*, 19(9):3596–3604, 2018.
- [225] L. Gemmer, C. Beaumont, and S. Ings. Dynamic modelling of passive margin salt tectonics: effects of water loading, sediment properties and sedimentation patterns. *Basin Research*, 17:383–402, 2005.
- [226] L. Gemmer, S.J. Ings, S. Medvedev, and C. Beaumont. Salt tectonics driven by differential sediment loading: stability analysis and finite-element experiments. *Basin Research*, 16:199–218, 2004.
- [227] L. Geoffroy, E.B. Burov, and P. Werner. Volcanic passive margins: another way to break up continents. *Scientific Reports*, 5:DOi:10.1038/srep14828, 2015.
- [228] 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. *Geochem. Geophys. Geosyst.*, 13(4):Q04001, doi:10.1029/2011GC003934, 2012.
- [229] M. Gérault, L. Husson, M.S. Miller, and E.D. Humphreys. Flat-slab subduction, topography, and mantle dynamics in southwestern Mexico. *Tectonics*, 34:10.1002/2015TC003908, 2015.
- [230] M. Gerbault. Pressure conditions for shear and tensile failure around a circular magma chamber; insight from elasto-plastic modelling. *Geological Society, London, Special Publications*, 367:111–130, 2012.
- [231] M. Gerbault, F. Cappa, and R. Hassani. Elasto-plastic and hydromechanical models of failure around an infinitely long magma chamber. *Geochem. Geophys. Geosyst.*, 13(3):doi:10.1029/2011GC003917, 2012.
- [232] M. Gerbault, J. Cembrano, C. Mpodozis, M. Farias, and M. Pardo. Continental margin deformation along the Andean subduction zone: Thermo-mechanical models. *Phys. Earth. Planet. Inter.*, 177:180–205, 2009.
- [233] M. Gerbault, F. Davey, and S. Henrys. Three-dimensional lateral crustal thickening in continental oblique collision: an example from the Southern Alps, New Zealand. *Geophys. J. Int.*, 150:770–779, 2002.
- [234] M. Gerbault, S. Henrys, and F. Davey. Numerical models of lithospheric deformation forming the Southern Alps of New Zealand. *J. Geophys. Res.*, 108(B7), 2003.
- [235] M. Gerbault, A.N.B. Poliakov, and M. Daignieres. Prediction of faulting from the theories of elasticity and plasticity: what are the limits? *Journal of Structural Geology*, 20:301–320, 1998.
- [236] M. Gerbault and W. Willingshofer. Lower crust indentation or horizontal ductile flow during continental collision? *Tectonophysics*, 387:169–187, 2004.
- [237] T. Gerya. Dynamical instability produces transform faults at mid-ocean ridges. *Science*, 329:1047–1050, 2010.
- [238] T. Gerya. Future directions in subduction modeling. *Journal of Geodynamics*, 52:344–378, 2011.
- [239] T. Gerya and B. Stöckhert. Two-dimensional numerical modeling of tectonic and metamorphic histories at active continental margins. *Int J Earth Sci (Geol Rundsch)*, 95:250–274, 2006.
- [240] T. Gerya and D.A. Yuen. Robust characteristics method for modelling multiphase visco-elasto-plastic thermo-mechanical problems. *Phys. Earth. Planet. Inter.*, 163:83–105, 2007.
- [241] Taras Gerya. *Numerical Geodynamic Modelling*. Cambridge University Press, 2010.
- [242] T.V. Gerya, J.A.D. Connolly, and D.A. Yuen. Why is terrestrial subduction one-sided ? *Geology*, 36(1):43–46, 2008.

- [243] T.V. Gerya, J.A.D. Connolly, D.A. Yuen, W. Gorczyk, and A.M. Capel. Seismic implications of mantle wedge plumes. *Phys. Earth. Planet. Inter.*, 156:59–74, 2006.
- [244] T.V. Gerya, D. Fossati, C. Cantieni, and D. Seward. Dynamic effects of aseismic ridge subduction: numerical modelling. *Eur. J. Mineral.*, 21:649–661, 2009.
- [245] 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. *Geochem. Geophys. Geosyst.*, 14(4), 2013.
- [246] T.V. Gerya and F.I. Meilick. Geodynamic regimes of subduction under an active margin: effects of rheological weakening by fluids and melts. *Journal of Metamorphic Geology*, 29:7–31, 2011.
- [247] T.V. Gerya, L.L. Perchuk, W.V. Maresch, and A.P. Willner. Inherent gravitational instability of hot continental crust: Implications for doming and diapirism in granulite facies terrains. *Geological Society of America*, 380:97–115, 2004.
- [248] T.V. Gerya, R.J. Stern, M.Baes, S.V. Sobolev, and S.A. Whattam. Plate tectonics on the Earth triggered by plume-induced subduction initiation. *Nature*, 527, 2015.
- [249] T.V. Gerya, R. Uken, J. Reinhardt, M. Watkeys, W.V. Maresch, and B.M. Clarke. Cold fingers in a hot magma: Numerical modeling of country-rock diapirs in the Bushveld Complex, South Africa. *Geology*, 31(9):753–756, 2003.
- [250] 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. *Phys. Earth. Planet. Inter.*, 140:293–318, 2003.
- [251] T.V. Gerya and D.A. Yuen. Rayleigh-Taylor instabilities from hydration and melting propels ‘cold plumes’ at subduction zones. *Earth Planet. Sci. Lett.*, 212:47–62, 2003.
- [252] T.V. Gerya, D.A. Yuen, and W.V. Maresch. Thermomechanical modelling of slab detachment. *Earth Planet. Sci. Lett.*, 226:101–116, 2004.
- [253] T.V. Gerya, D.A. Yuen, and E.O.D. Sevre. Dynamical causes for incipient magma chambers above slabs. *Geology*, 32(1):89–92, 2004.
- [254] R.K. Ghazian and S.J.H. Buitter. A numerical investigation of continental collision styles. *GJI*, 2013.
- [255] R.K. Ghazian and S.J.H. Buitter. Numerical modelling of the role of salt in continental collision: An application to the southeast Zagros fold-and-thrust belt. *Tectonophysics*, in press, 2014.
- [256] U. Ghia, K.N. Ghia, and C.T. Shin. High-Re Solutions for incompressible flow using the Navier-Stokes equations and a multigrid method. *J. Comp. Phys.*, 48:387–411, 1982.
- [257] G. Gibert, M. Gerbault, R. Hassani, and E. Tric. Dependency of slab geometry on absolute velocities and conditions for cyclicity: insights from numerical modelling. *Geophys. J. Int.*, 189:747–760, 2012.
- [258] E. Di Giuseppe, J. van Hunen, F. Funiciello, C. Faccenna, and D. Giardini. Slab stiffness control of trench motion: Insights from numerical models. *Geochem. Geophys. Geosyst.*, 9(2), 2008.
- [259] A. Glerum, C. Thieulot, M. Fraters, C. Blom, and W. Spakman. Nonlinear viscoplasticity in ASPECT: benchmarking and applications to subduction. *Solid Earth*, 9(2):267–294, 2018.
- [260] R. Glowinski. *Handbook of Numerical Analysis, vol IX: Numerical methods for fluids*. North-Holland, 2003.
- [261] Oguz H. Gogus. Rifting and subsidence following lithospheric removal in continental back arcs. *Geology*, page doi:10.1130/G36305.1, 2014.
- [262] G.H. Golub and C.F. van Loan. *Matrix Computations, 4th edition*. John Hopkins University Press, 2013.

- [263] W. Gorczyk, T.V. Gerya, J.A.D. Connolly, and D.A. Yuen. Growth and mixing dynamics of mantle wedge plumes. *Geology*, 35(7):587–590, 2007.
- [264] W. Gorczyk, T.V. Gerya, J.A.D. Connolly, D.A. Yuen, and M. Rudolph. Large-scale rigid-body rotation in the mantle wedge and its implications for seismic tomography . *Geochem. Geophys. Geosyst.*, 7(5):10.1029/2005GC001075, 2006.
- [265] W. Gorczyk, H. Smithies, F. Korhonen, H. Howard, and R. Quentin De Gromard. Ultra-hot Mesoproterozoic evolution of intracontinental central Australia. *Geoscience Frontiers*, 6(1):23–37, 2014.
- [266] R. Goteti, C. Beaumont, and S.J. Ings. Factors controlling early stage salt tectonics at rifted continental margins and their thermal consequences. *J. Geophys. Res.*, 117:1–31, 2013.
- [267] R. Govers and M.J.R. Wortel. Extension of stable continental lithosphere and the initiation of lithospheric scale faults. *Tectonics*, 14(4):1041–1055, 1995.
- [268] R. Govers and M.J.R. Wortel. Some remarks on the relation between vertical motions of the lithosphere during extension and the necking depth parameter inferred from kinematic modeling studies. *J. Geophys. Res.*, 104:23,245–23,253, 1999.
- [269] R. Govers and M.J.R. Wortel. Lithosphere tearing at STEP faults: Response to edges of subduction zones . *Earth Planet. Sci. Lett.*, 236:505–523, 2005.
- [270] S. Gradmann and C. Beaumont. Coupled fluid flow and sediment deformation in margin-scale salt-tectonic systems: 2. Layered sediment models and application to the northwestern Gulf of Mexico. *Tectonics*, 31(TC4011), 2012.
- [271] S. Gradmann, C. Beaumont, and M. Albertz. Factors controlling the evolution of the Perdido Fold Belt, northwestern Gulf of Mexico, determined from numerical models. *Tectonics*, 28(TC2002), 2009.
- [272] R. Gray and R.N. Pysklywec. Geodynamic models of Archean continental collision and the formation of mantle lithosphere keels. *Geophys. Res. Lett.*, 37(L19301), 2010.
- [273] R. Gray and R.N. Pysklywec. Geodynamic models of mature continental collision: Evolution of an orogen from lithospheric subduction to continental retreat/delamination. *J. Geophys. Res.*, 117(B03408), 2012.
- [274] R. Gray and R.N. Pysklywec. Influence of sediment deposition on deep lithospheric tectonics. *Geophys. Res. Lett.*, 39(L11312), 2012.
- [275] R. Gray and R.N. Pysklywec. Influence of viscosity pressure dependence on deep lithospheric tectonics during continental collision. *J. Geophys. Res.*, 118, 2013.
- [276] P.M. Gresho, S.T. Chan, M.A. Christon, and A.C. Hindmarsch. A little more on stabilised  $Q_1 Q_1$  for transient viscous incompressible flow. *Int. J. Num. Meth. Fluids*, 21:837–856, 1995.
- [277] P.M. Gresho and R.L. Sani. *Incompressible flow and the Finite Element Method, vol II*. John Wiley and Sons, Ltd, 2000.
- [278] M. Griebel, T. Dornseifer, and T. Neunhoeffer. *Numerical simulation in Fluid Dynamics*. SIAM.
- [279] D. Griffiths and D. Silvester. Unstable modes of the q1-p0 element. Technical Report 257, University of MAnchester/UMIST, 1994.
- [280] S.G. Grigull, S.M. Ellis, T.A. Little, M.P. Hill, and S.J.H. Buiter. Rheological constraints on quartz derived from scaling relationships and numerical models of sheared brittle-ductile quartz veins, central Southern alps, New Zealand. *Journal of Structural Geology*, 37:200–222, 2012.
- [281] L. Gross, L. Bourgouin, A. Hale, and H.-B. Mühlhaus. Interface modeling in incompressible media using level sets in Escript. *Phys. Earth. Planet. Inter.*, 163:23–34, 2007.

- [282] J.-L. Guermond, R. Pasquetti, and Bojan Popov. Entropy viscosity method for nonlinear conservation laws. *J. Comp. Phys.*, page doi:10.1016/j.jcp.2010.11.043, 2011.
- [283] L. Guillou-Frottier, E. Burov, S. Cloetingh, E. Le Goff, Y. Deschamps, B. Huet, and V. Bouchet. Plume-induced dynamic instabilities near cratonic blocks: Implications for PT paths and metallogeny. *Global and Planetary Change*, 90-91:37–50, 2012.
- [284] M. Gunzburger. *Finite Element Methods for Viscous Incompressible Flows: A Guide to Theory, Practice and Algorithms*. Academic, Boston, 1989.
- [285] M. Gunzburger, P. Bochev, and R. Lehoucq. On stabilized finite element methods for the Stokes problem in the small time step limit. 2006.
- [286] M. Gurnis, C. Hall, and L. Lavier. Evolving force balance during incipient subduction. *Geochem. Geophys. Geosyst.*, 5(7), 2004.
- [287] M. Gurnis, J.X. Mitrovica, J. Ritsema, and H.-J. van Heijst. Constraining mantle density structure using geological evidence of surface uplift rates: The case of the African superplume. *Geochem. Geophys. Geosyst.*, 1, 2000.
- [288] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*. Springer-Verlag, Berlin.
- [289] A.J. Hale, L. Bourgouin, and H.B. Muehlhaus. Using the level set method to model endogenous lava dome growth . *J. Geophys. Res.*, 112(B03213), 2007.
- [290] A.J. Hale, K.-D. Gottschaldt, G. Rosenbaum, L. Bourgouin, M. Bauchy, and Hans Mühlhaus. Dynamics of slab tear faults: Insights from numerical modelling. *Tectonophysics*, 483:58–70, 2010.
- [291] C.E. Hall, M. Gurnis, M. Sdrolias, L.L. Lavier, and R.D. Mueller. Catastrophic initiation of subduction following forced convergence across fracture zones. *Earth Planet. Sci. Lett.*, 212:15–30, 2003.
- [292] D.L. Hansen. A meshless formulation for geodynamic modeling. *J. Geophys. Res.*, 108:doi:10.1029/2003JB002460, 2003.
- [293] F.H. Harlow and J.E. Welch. Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface. *The physics of fluids*, 8(12):2182, 1965.
- [294] R. Hassan, N. Flament, M. Gurnis, D.J. Bower, and D. Müller. Provenance of plumes in global convection models. *Geochem. Geophys. Geosyst.*, 16:1465–1489, 2015.
- [295] R. Hassani, D. Jongmans, and Jean Chéry. Study of plate deformation and stress in subduction processes using two-dimensional numerical models. *J. Geophys. Res.*, 102(B8):17,951–17,96, 1997.
- [296] K.L. Haynie and M.A. Jadamec. Tectonic drivers of the Wrangell block: Insights on fore-arc sliver processes from 3-D geodynamic models of Alaska. *Tectonics*, 36, 2017.
- [297] 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. *Phys. Earth. Planet. Inter.*, 263:23–37, 2017.
- [298] C. Heine and S. Brune. Oblique rifting of the Equatorial Atlantic: Why there is no Saharan Atlantic Ocean. *Geology*, 42(3):211–214, 2014.
- [299] T. Heister, J. Dannberg, R. Gassmöller, and W. Bangerth. High Accuracy Mantle Convection Simulation through Modern Numerical Methods. II: Realistic Models and Problems. *Geophys. J. Int.*, 210(2):833–851, 2017.
- [300] P. J. Heron, R. N. Pysklywec, R. Stephenson, and J. van Hunen. Deformation driven by deep and distant structures: Influence of a mantle lithosphere suture in the ouachita orogeny, southeastern united states. *Geology*, 2018.

- [301] P.J. Heron, R.N. Pysklywec, and R. Stephenson. Intraplate orogenesis within accreted and scarred lithosphere: Example of the Eurekan Orogeny, Ellesmere Island. *Tectonophysics*, 2015.
- [302] B.H. Heyn, C.P. Conrad, and R.G. Tronnes. Stabilizing Effect of Compositional Viscosity Contrasts on Thermochemical Piles. *Geophys. Res. Lett.*, 45:7523–7532, 2018.
- [303] 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. *Solid Earth*, 5(2):1087–1098, 2014.
- [304] J.M. Hines and M.I. Billen. Sensitivity of the short- to intermediate-wavelength geoid to rheologic structure in subduction zones. *J. Geophys. Res.*, 117(B05410), 2012.
- [305] T. Höink and A. Lenardic. Three-dimensional mantle convection simulations with a low-viscosity asthenosphere and the relationship between heat flow and the horizontal length scale of convection. *Geophys. Res. Lett.*, 35(L10304), 2008.
- [306] P. Hood and C. Taylor. NavierStokes equations using mixed interpolation. In J.T. Oden, R.H. Gallagher, O.C. Zienkiewicz, and C. Taylor, editors, *Finite element methods in flow problems*. Huntsville Press, University of Alabama, 1974.
- [307] Ch. Huber, A. Parmigiani, B. Chopard, M. Manga, and O. Bachmann. Lattice Boltzmann model for melting with natural convection. *International Journal of Heat and Fluid Flow*, 29:1469–1480, 200.
- [308] T.J.R. Hughes. *The Finite Element Method. Linear Static and Dynamic Finite Element Analysis*. Dover Publications, Inc., 2000.
- [309] T.J.R. Hughes and A. Brooks. A theoretical framework for petrov-galerkin methods with discontinuous weighting functions: application to the streamline-upwind procedure. *Finite Elements in Fluids*, 4:47–65, 1982.
- [310] T.J.R. Hughes, L.P. Franca, and M. Balestra. A new finite element formulation for computational fluid dynamics: V. Circumventing the Babuка-Brezzi condition: A stable Petrov-Galerkin formulation of the Stokes problem accommodating equal-order interpolations. *Computer Methods in Applied Mechanics and Engineering*, 59(1):85–99, 1986.
- [311] T.J.R. Hughes, W.K. Liu, and A. Brooks. Finite element analysis of Incompressible viscous flows by the penalty function formulation. *J. Comp. Phys.*, 30:1–60, 1979.
- [312] Hoon Huh, Choong Ho Lee, and Wei H. Yang. A general algorithm for plastic flow simulation by finite element limit analysis. *International Journal of Solids and Structures*, 36:1193–1207, 1999.
- [313] R. Huismans and C. Beaumont. Depth-dependent extension, two-stage breakup and cratonic underplating at rifted margins. *Nature*, 473:74–79, 2011.
- [314] R. S. Huismans and C. Beaumont. Symmetric and asymmetric lithospheric extension: Relative effects of frictional-plastic and viscous strain softening. *J. Geophys. Res.*, 108 (B10)(2496), 2003.
- [315] R.S. Huismans and C. Beaumont. Complex rifted continental margins explained by dynamical models of depth-dependent lithospheric extension. *Geology*, 30(3):211–214, 2002.
- [316] 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*, volume 282, pages 111–138. Geological Society, London, Special Publications, 2007.
- [317] R.S. Huismans, S.J.H. Buiter, and C. Beaumont. Effect of plastic-viscous layering and strain softening on mode selection during lithospheric extension. *J. Geophys. Res.*, 110:B02406, 2005.
- [318] L. Husson, C.P. Conrad, and C. Faccenna. Plate motions, Andean orogeny, and volcanism above the South Atlantic convection cell. *Earth Planet. Sci. Lett.*, 317-318:126–135, 2012.

- [319] F. Ilinca and D. Pelletier. Computation of accurate nodal derivatives of finite element solutions: The finite node displacement method. *Int. J. Num. Meth. Eng.*, 71:1181–1207, 2007.
- [320] Alik Ismail-Zadeh and Paul Tackley. *Computational Methods for Geodynamics*. Cambridge University Press, 2010.
- [321] 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. *J. Geophys. Res.*, 99(B8):15,919–15,938, 1994.
- [322] J. Ita and S.D. King. The influence of thermodynamic formulation on simulations of subduction zone geometry and history. *Geophys. Res. Lett.*, 25(9):1463–1466, 1998.
- [323] M.A. Jadamec and M.I. Billen. The role of rheology and slab shape on rapid mantle flow: Three-dimensional numerical models of the Alaska slab edge. *J. Geophys. Res.*, 117(B02304), 2012.
- [324] M.A. Jadamec, M.I. Billen, and S.M. Roeske. Three-dimensional numerical models of flat slab subduction and the Denali fault driving deformation in south-central Alaska. *Earth Planet. Sci. Lett.*, 376:29–42, 2013.
- [325] S. Jammes and R.S. Huismans. Structural styles of mountain building: Controls of lithospheric rheologic stratification and extensional inheritance. *J. Geophys. Res.*, 117, 2012.
- [326] S. Jammes, R.S. Huismans, and J.A. Muñoz. Lateral variation in structural style of mountain building: controls of rheological and rift inheritance. *Terra Nova*, 0:doi:10.1111/ter.12087, 2013.
- [327] Y. Jaquet, Th. Duretz, and S.M. Schmalholz. Dramatic effect of elasticity on thermal softening and strain localization during lithospheric shortening. *Geophys. J. Int.*, 204:780–784, 2016.
- [328] G.T. Jarvis. Effects of curvature on two-dimensional models of mantle convection: cylindrical polar coordinates. *J. Geophys. Res.*, 98(B3):4477–4485, 1993.
- [329] C. Jaupart and J.-C. Mareschal. *Heat Generation and Transport in the Earth*. Cambridge, 2011.
- [330] V. John. *Finite Element Methods for Incompressible Flow Problems*. Springer, 2016.
- [331] L. Jolivet, P. Davy, and P. Cobbold. Right-lateral shear along the Northwest Pacific margin and the India-Eurasia collision. *Tectonics*, 9(6):1409–1419, 1990.
- [332] A. Jourdon, L. Le Pourhiet, Mouthereau F, and E. Masini. Role of rift maturity on the architecture and shortening distribution in mountain belts. *Earth Planet. Sci. Lett.*, 512:89–99, 2019.
- [333] A. Jourdon, L. Le Pourhiet, C. Petit, and Y. Rolland. Impact of range-parallel sediment transport on 2D thermo-mechanical models of mountain belts: Application to the Kyrgyz Tien Shan. *Terra Nova*, 30:279–288, 2018.
- [334] L.M. Kachanov. *Fundamentals of the Theory of Plasticity*. Dover Publications, Inc., 2004.
- [335] L. Kaislaniemi and J. van Hunen. Dynamics of lithospheric thinning and mantle melting by edge-driven convection: Application to Moroccan Atlas mountains. *Geochem. Geophys. Geosyst.*, 15:3175–3189, 2014.
- [336] B. Kaus. *Modelling approaches to geodynamic processes, PhD thesis*. PhD thesis, ETH Zurich, 2005.
- [337] B.J.P. Kaus. Factors that control the angle of shear bands in geodynamic numerical models of brittle deformation. *Tectonophysics*, 484:36–47, 2010.
- [338] B.J.P. Kaus, H. Mühlhaus, and D.A. May. A stabilization algorithm for geodynamic numerical simulations with a free surface. *Phys. Earth. Planet. Inter.*, 181:12–20, 2010.
- [339] B.J.P. Kaus, A.A. Popov, T.S. Baumann, A.E. Pusok, A. Bauville, N. Fernandez, and M. Collignon. Forward and Inverse Modelling of Lithospheric Deformation on Geological Timescales. *NIC Symposium 2016*, pages 299–307, 2016.

- [340] M. Kawaguti. Numerical solution of the Navier-Stokes equations for the flow in a two-dimensional cavity. *Journal of the Physical Society of Japan*, 16(12):2307–2315, 1961.
- [341] D.F. Keppe, C.A. Currie, and C. Warren. Subduction erosion modes: comparing finite element numerical models with the geological record. *Earth Planet. Sci. Lett.*, 287:241–254, 2009.
- [342] M. Kimmritz and M. Braack. iDiscretization of the hydrostatic Stokes system by stabilized finite elements of equal order.
- [343] S. King, C. Lee, P. van Keeken, W. Leng, S. Zhong, E. Tan, N. Tosi, and M. Kameyama. A community benchmark for 2D Cartesian compressible convection in the Earths mantle. *Geophys. J. Int.*, 180:7387, 2010.
- [344] S.D. King. On topography and geoid from 2-D stagnant lid convection calculations. *Geochem. Geophys. Geosyst.*, 10(3), 2009.
- [345] S.D. King. Venus Resurfacing Constrained by Geoid and Topography. *JGR: Planets*, 123:1041–1060, 2018.
- [346] S.D. King and D.L. Anderson. An alternative mechanism of flood basalt formation. *Earth Planet. Sci. Lett.*, 136:269–279, 1995.
- [347] S.D. King and D.L. Anderson. Edge-driven convection. *Earth Planet. Sci. Lett.*, 160:289–296, 1998.
- [348] S.D. King, D.J. Frost, and D.C. Rubie. Why cold slabs stagnate in the transition zone. *Geology*, page doi:10.1130/G36320.1, 2015.
- [349] S.D. King, A. Raefsky, and B.H. Hager. ConMan: Vectorizing a finite element code for incompressible two-dimensional convection in the Earths mantle. *Phys. Earth. Planet. Inter.*, 59:195–208, 1990.
- [350] A. Kiraly, F.A. Capitanio, F. Funiciello, and C. Faccenna. Subduction zone interaction: Controls on arcuate belts. *Geology*, 2016.
- [351] 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. *Geochem. Geophys. Geosyst.*, 2013.
- [352] T. Komut, R. Gray, R. Pysklywec, and O. Gogus. Mantle flow uplift of western Anatolia and the Aegean: Interpretation from geophysical analyses and geodynamic modeling. *J. Geophys. Res.*, 117(B11412), 2012.
- [353] H. Koopmann, S. Brune, D. Franke, and S. Breuer. Linking rift propagation barriers to excess magmatism at volcanic rifted margins. *Geology*, page doi:10.1130/G36085.1, 2014.
- [354] A. Koptev, A. Beniest, L. Jolivet, and S. Leroy. PlumeInduced Breakup of a Subducting Plate: Microcontinent Formation Without Cessation of the Subduction Process. *Geophys. Res. Lett.*, 46:3663–3675, 2019.
- [355] J.R. Koseff and R.L. Street. The Lid-Driven Cavity Flow: A Synthesis of Qualitative and Quantitative Observations. *J. Fluids Eng*, 106:390–398, 1984.
- [356] L.I.G. Kovasznay. Laminar flow behind a two-dimensional grid. *Mathematical Proceedings of the Cambridge Philosophical Society*, 44(1):58–62, 1948.
- [357] Peter Kovesi. Good colour maps: How to design them. *CoRR*, abs/1509.03700, 2015.
- [358] M. Kronbichler, T. Heister, and W. Bangerth. High accuracy mantle convection simulation through modern numerical methods . *Geophys. J. Int.*, 191:12–29, 2012.
- [359] D. Kurfess and O. Heidbach. CASQUS: A new simulation tool for coupled 3D finite element modeling of tectonic and surface processes based on ABAQUS and CASCADE. *Computers and Geosciences*, 35:1959–1967, 2009.

- [360] W. Landry, L. Hodkinson, and S. Kientz. Gale user manual. Technical report, CIG, VPAC, 2011.
- [361] A. Lavecchia, C. Thieulot, F. Beekman, S. Cloetingh, and S. Clark. Lithosphere erosion and continental breakup: Interaction of extension, plume upwelling and melting. *Earth Planet. Sci. Lett.*, 467:89–98, 2017.
- [362] L. Le Pourhiet, B. Huet, D.A. May, L. Labrousse, and L. Jolivet. Kinematic interpretation of the 3D shapes of metamorphic core complexes. *Geochem. Geophys. Geosyst.*, 13(Q09002), 2012.
- [363] S.M. Lechmann, D.A. May, B.J.P. Kaus, and S.M. Schmalholz. Comparing thin-sheet models with 3-D multilayer models for continental collision. *Geophys. J. Int.*, 187:10–33, 2011.
- [364] S.M. Lechmann, S.M. Schmalholz, G. Hetenyi, D.A. May, and B.J.P. Kaus. Quantifying the impact of mechanical layering and underthrusting on the dynamics of the modern India-Asia collisional system with 3-D numerical models. *J. Geophys. Res.*, 119:doi:10.1002/2012JB009748, 2014.
- [365] R. Lee, P. Gresho, and R. Sani. Smoothing techniques for certain primitive variable solutions of the Navier-Stokes equations. . *Int. J. Num. Meth. Eng.*, 14:1785–1804, 1979.
- [366] V. Lemiale, H.-B. Mühlhaus, L. Moresi, and J. Stafford. Shear banding analysis of plastic models formulated for incompressible viscous flows. *Phys. Earth. Planet. Inter.*, 171:177–186, 2008.
- [367] A. Lenardic, L. Moresi, A.M. Jellinek, C.J. O'Neill, C.M. Cooper, and C.T. Lee. Continents, supercontinents, mantle thermal mixing, and mantle thermal isolation: Theory, numerical simulations, and laboratory experiments. *Geochem. Geophys. Geosyst.*, 12(10), 2011.
- [368] W. Leng, M. Gurnis, and P. Asimov. From basalts to boninites: The geodynamics of volcanic expression during induced subduction initiation. *Initiation and Termination of Subduction: Rock Re- cord, Geodynamic Models, Modern Plate Boundaries*, 2012.
- [369] W. Leng and S. Zhong. Viscous heating, adiabatic heating and energetic consistency in compressible mantle convection. *Geophys. J. Int.*, 173:693–702, 2008.
- [370] W. Leng and S. Zhong. Implementation and application of adaptive mesh refinement for thermochemical mantle convection studies. *Geochem. Geophys. Geosyst.*, 12(4), 2011.
- [371] Wei Leng and Michael Gurnis. Dynamics of subduction initiation with different evolutionary pathways. *Geochem. Geophys. Geosyst.*, 12(12):Q12018, doi:10.1029/2011GC003877, 2011.
- [372] Wei Leng and Michael Gurnis. Subduction initiation at relic arcs. *Geophys. Res. Lett.*, 42:7014–7021, 2015.
- [373] 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. *Computing*, 86:37–51, 2009.
- [374] Z.-H. Li, Z. Xu, T. Gerya, and J.-P. Burg. Collision of continental corner from 3-D numerical modeling. *Earth Planet. Sci. Lett.*, 380:98–111, 2013.
- [375] Z.H. Li, Z.Q. Xu, and T.V. Gerya. Flat versus steep subduction: Contrasting modes for the formation and exhumation of high- to ultrahigh-pressure rocks in continental collision zones. *Earth Planet. Sci. Lett.*, 301:65–77, 2011.
- [376] Z.-X. Lia and S. Zhong. Supercontinentsuperplume coupling, true polar wander and plume mobility: Plate dominance in whole-mantle tectonics. *Phys. Earth. Planet. Inter.*, 176:143–156, 2009.
- [377] J. Liao and T. Gerya. From continental rifting to seafloor spreading: Insight from 3D thermo-mechanical modeling. *Gondwana Research*, 2014.
- [378] J. Liao and T. Gerya. Influence of lithospheric mantle stratification on craton extension: Insight from two-dimensional thermo-mechanical modeling. *Tectonophysics*, page doi:10.1016/j.tecto.2014.01.020, 2014.

- [379] S.-C. Lin and P.E. van Keken. Dynamics of thermochemical plumes: 1. Plume formation and entrainment of a dense layer. *Geochem. Geophys. Geosyst.*, 7(2), 2006.
- [380] S.-C. Lin and P.E. van Keken. Dynamics of thermochemical plumes: 2. Complexity of plume structures and its implications for mapping mantle plumes . *Geochem. Geophys. Geosyst.*, 7(3), 2006.
- [381] L. Liu and D.R. Stegman. Segmentation of the Farallon slab. *Earth Planet. Sci. Lett.*, 311:1–10, 2011.
- [382] S. Liu and C.A. Currie. Farallon plate dynamics prior to the Laramide orogeny: Numerical models of flat subduction. *Tectonophysics*, 666:33–47, 2016.
- [383] S. Liu and S.D. King. A benchmark study of incompressible Stokes flow in a 3-D spherical shell using ASPECT. *Geophys. J. Int.*, 217:650–667, 2019.
- [384] 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. *Geophys. J. Int.*, 194:125–144, 2013.
- [385] Z. Liu and P. Bird. Two-dimensional and three-dimensional finite element modelling of mantle processes beneath central South Island, New Zealand. *Geophys. J. Int.*, 165:1003–1028, 2006.
- [386] C. Loiselet, J. Braun, L. Husson, C. Le Carlier de Veslud, C. Thieulot, P. Yamato, and D. Grujic. Subducting slabs: Jellyfishes in the Earth’s mantle. *Geochem. Geophys. Geosyst.*, 11(8):doi:10.1029/2010GC003172, 2010.
- [387] G. Lu, B.J.P. Kaus, L. Zhao, and T. Zheng. Self-consistent subduction initiation induced by mantle flow. *Terra Nova*, page doi: 10.1111/ter.12140, 2015.
- [388] 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. *Geochem. Geophys. Geosyst.*, 16:1753–1770, 2015.
- [389] V. Magni, M.B. Allen, J. van Hunen, and P. Bouilhol. Continental underplating after slab break-off. *Earth Planet. Sci. Lett.*, 474:59–67, 2017.
- [390] V. Magni, P. Bouilhol, and J. van Hunen. Deep water recycling through time. *Geochem. Geophys. Geosyst.*, 15:4203–4216, 2014.
- [391] C. Malatesta, T. Gerya, L. Crispini, L. Federico, and G. Capponi. Oblique subduction modelling indicates along-trench tectonic transport of sediments. *Nature Communications*, 4, 2013.
- [392] C. Malatesta, T. Gerya, L. Crispini, L. Federico, and G. Capponi. Interplate deformation at early-stage oblique subduction: 3-D thermomechanical numerical modeling. *Tectonics*, 35:1610–1625, 2016.
- [393] D.S. Malkus and T.J.R. Hughes. Mixed finite element methods - reduced and selective integration techniques: a unification of concepts. *Comput. Meth. Appl. Mech. Eng.*, 15:63–81, 1978.
- [394] L.E. Malvern. *Introduction to the mechanics of a continuous medium*. Prentice-Hall, Inc., 1969.
- [395] V. Manea and M. Gurnis. Subduction zone evolution and low viscosity wedges and channels. *Earth Planet. Sci. Lett.*, 264:22–45, 2007.
- [396] V.C. Manea, W.P. Leeman, T. Gerya, M. Manea, and G. Zhu. Subduction of fracture zones controls mantle melting and geochemical signature above slabs. *Nature Communications*, page doi:10.1038/ncomms6095, 2014.
- [397] G. Maniatis, D. Kurfess, A. Hampel, and O. Heidbach. Slip acceleration on normal faults due to erosion and sedimentation Results from a new three-dimensional numerical model coupling tectonics and landscape evolution. *Earth Planet. Sci. Lett.*, 284:570–582, 2009.

- [398] G. Marketos, R. Govers, and C.J. Spiers. Ground motions induced by a producing hydrocarbon reservoir that is overlain by a viscoelastic rocksalt layer: a numerical model. *Geophys. J. Int.*, 203:228–242, 2015.
- [399] A.M. Marotta and M.I. Spalla. Permian-Triassic high thermal regime in the Alps: Result of late Variscan collapse or continental rifting? Validation by numerical modeling. *Tectonics*, 26(TC4016):oi:10.1029/2006TC002047, 2007.
- [400] A.M. Marotta, E. Spelta, and C. Rizzetto. Gravity signature of crustal subduction inferred from numerical modelling. *Geophys. J. Int.*, 166:923–938, 2006.
- [401] F.O. Marques and B.J.P. Kaus. Speculations on the impact of catastrophic subduction initiation on the Earth System. *Journal of Geodynamics*, 93:1–16, 2016.
- [402] F.O. Marques, K. Nikolaeva, M. Assumpcao, T.V. Gerya, F.H.R. Bezerra, A.F. do Nascimento, and J.M. Ferreira. Testing the influence of far-field topographic forcing on subduction initiation at a passive margin. *Tectonophysics*, 608:517–524, 2013.
- [403] J. Martinod, V. Regard, Y. Letourmy, H. Henry, R. Hassani, S. Baratchart, and S. Carretier. How do subduction processes contribute to forearc Andean uplift? Insights from numerical models. *Journal of Geodynamics*, 2015.
- [404] W.G. Mason, L. Moresi, P.G. Betts, and M.S. Miller. Three-dimensional numerical models of the influence of a buoyant oceanic plateau on subduction zones. *Tectonophysics*, 483:71–79, 2010.
- [405] B. Maunder, J. van Hunen, P. Bouilhol, and V. Magni. Modeling Slab Temperature: A Reevaluation of the Thermal Parameter. *Geochem. Geophys. Geosyst.*, 2019.
- [406] D.A. May. Volume reconstruction of point cloud data sets derived from computational geodynamic simulations. *Geochem. Geophys. Geosyst.*, 13(5):Q05019, 2012.
- [407] D.A. May, J. Brown, and L. Le Pourhiet. A scalable, matrix-free multigrid preconditioner for finite element discretizations of heterogeneous Stokes flow. *Computer Methods in Applied Mechanics and Engineering*, 290:496–523, 2015.
- [408] S. McKee, M.F. Tome, V.G. Ferreira, J.A. Cuminato, A. Castelo, F.S. Sousa, and N. Mangiavacchi. The MAC method. *Computers and Fluids*, 37:907–930, 2008.
- [409] D.P. McKenzie. Speculations on the consequences and causes of plate motions. *Geophys. J. R. astr. Soc.*, 18:1–32, 1969.
- [410] A.K. McNamara and S. Zhong. Thermochemical structures within a spherical mantle: Superplumes or piles? *J. Geophys. Res.*, 109(B07402), 2004.
- [411] 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. *Phys. Earth. Planet. Inter.*, 283:110–121, 2018.
- [412] A. Mizukami. A mixed Finite Element method for boundary flux computation. *Computer Methods in Applied Mechanics and Engineering*, 57:239–243, 1986.
- [413] 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. *Geology*, 5:212–216, 1977.
- [414] C. Morency, R.S. Huismans, C. Beaumont, and P. Fullsack. A numerical model for coupled fluid flow and matrix deformation with applications to disequilibrium compaction and delta stability. *J. Geophys. Res.*, 112(B10407), 2007.
- [415] L. Moresi, P.G. Betts, M.S. Miller, and R.A. Cayley. Dynamics of continental accretion. *Nature*, 2014.

- [416] L. Moresi, F. Dufour, and H.B. Mühlhaus. A Lagrangian integration point finite element method for large deformation modeling of visco-elastic geomaterials. *J. Comp. Phys.*, 184(2):476–497, 2003.
- [417] 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. *Phys. Earth. Planet. Inter.*, 163:69–82, 2007.
- [418] L. Moresi and V. Solomatov. Mantle convection with a brittle lithosphere: thoughts on the global tectonics styles of the earth and venus. *Geophys. J. Int.*, 133:669–682, 1998.
- [419] M. Morishige and P.E. van Keken. Along-arc variation in the 3-D thermal structure around the junction between the Japan and Kurile arcs. *Geochem. Geophys. Geosyst.*, 15:2225–2240, 2014.
- [420] R. Moucha, A.M. Forte, J.X. Mitrovica, and A. Daradich. Lateral variations in mantle rheology: implications for convection related surface observables and inferred viscosity models. *Geophys. J. Int.*, 169:113–135, 2007.
- [421] M.A. Murphy, M.H. Taylor, J. Gosse, C.R.P. Silver, D.M. Whipp, and C. Beaumont. Limit of strain partitioning in the Himalaya marked by large earthquakes in western Nepal. *Nature Geoscience*, 2013.
- [422] J. Naliboff and S.J.H. Buiter. Rift reactivation and migration during multiphase extension. *Earth Planet. Sci. Lett.*, 421:58–67, 2015.
- [423] J.B. Naliboff, M.I. Billen, T. Gerya, and J. saunders. Dynamics of outer-rise faulting in oceanic-continent subduction systems. *Geochem. Geophys. Geosyst.*, 14(7):10.1002/ggge.20155, 2013.
- [424] J.B. Naliboff, C. Lithgow-Bertelloni, L.J. Ruff, and N. de Koker. The effects of lithospheric thickness and density structure on Earths stress field. *Geophys. J. Int.*, 188:1–17, 2012.
- [425] 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. *Geophys. J. Int.*, 136:261–274, 1999.
- [426] A.M. Negredo, R. Sabadini, and C. Giunchi. Interplay between subduction and continental convergence:a three- dimensional dynamic model for the Central Mediterranean. *Geophys. J. Int.*, 131:F9–F13, 1997.
- [427] 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. *Solid Earth*, 9:1207–1224, 2018.
- [428] K. Nikolaeva, T.V. Gerya, and F.O. Marques. Subduction initiation at passive margins: numerical modeling. *J. Geophys. Res.*, 115(B03406), 2010.
- [429] F. Nilfouroushan, R. Pysklywec, A. Cruden, and H. Koyi. Thermal-mechanical modeling of salt-based mountain belts with pre-existing basement faults: Application to the Zagros fold and thrust belt, southwest Iran. *Tectonics*, 32:1212–1226, 2013.
- [430] L. Noack, A. Rivoldini, and T. Van Hoolst. CHIC Coupling Habitability, Interior and Crust. *INFO-COMP 2015 : The Fifth International Conference on Advanced Communications and Computation*, 2015.
- [431] S. Norburn and D. Silvester. Fourier analysis of stabilized Q1-Q1 mixed finite element approximation. *SIAM J. Numer. Anal.*, 39:817–833, 2001.
- [432] P. Olson, R. Deguen, L.A. Hinov, and S. Zhong. Controls on geomagnetic reversals and core evolution by mantle convection in the Phanerozoic. *Phys. Earth. Planet. Inter.*, 214:87–103, 214.
- [433] C. O'Neill, S. Marchi, S. Zhang, and W. Bottke. Impact-driven subduction on the hadean earth. *Nature Geoscience*, 10(10):793, 2017.

- [434] 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. *Computers and Geosciences*, 32:1769–1779, 2006.
- [435] C. J. O'Neill and S. Zhang. Lateral mixing processes in the hadean. *Journal of geophysical research.*, 123:7074–7089, 2018.
- [436] C.J. O'Neill, A. Lenardic, W.L. Griffin, and S.Y. O'Reilly. Dynamics of cratons in an evolving mantle. *Lithos*, 102:12–24, 2008.
- [437] S. Osher and R. Fedkiw. Level set methods: an overview and some recent results. *J. Comp. Phys.*, 169:463–502, 2001.
- [438] S. Osher and C.-W. Shu. High-order essentially non-oscillatory schemes for HamiltonJacobi equations. *SIAM J. Numer. Anal.*, 28:907–922, 1991.
- [439] M. OzBench, K. Regenauer-Lieb, D.R. Stegman, G. Morra, R. Farrington, A. Hale, D.A. May, J. Freeman, L. Bourgoin, H.-B. Mühlhaus, and L. Moresi. A model comparison study of large-scale mantle-lithosphere dynamics driven by subduction. *Phys. Earth. Planet. Inter.*, 171:224–234, 2008.
- [440] C. O'Neill, A. Lenardic, M. Weller, L. Moresi, S. Quenette, and S. Zhang. A window for plate tectonics in terrestrial planet evolution? *Phys. Earth. Planet. Inter.*, 255(80–92), 2016.
- [441] F. Pan and A. Acrivos. Steady flows in rectangular cavities. *J. Fluid Mech.*, 28(4):643–655, 1967.
- [442] G. Peltzer and P. Tapponnier. Formation and evolution of strike-slip faults, rifts, and basins during the india-asia collision: an experimental approach. *J. Geophys. Res.*, 93(B12):15085–15177, 1988.
- [443] J. Perry-Houts and L. Karlstrom. Anisotropic viscosity and time-evolving lithospheric instabilities due to aligned igneous intrusions. *Geophysical Journal International*, 216(2):794–802, 2018.
- [444] M. Peters, M. Veveakis, T. Poulet, A. Karrech, M. Herwegh, and K. Regenauer-Lieb. Boudinage as a material instability of elasto-visco-plastic rocks. *Journal of Structural Geology*, 78:86–102, 2015.
- [445] Thomas Philippe. *Single viscous layer fold interplay and linkage*. PhD thesis, ETH Zurich, 2013.
- [446] B.R. Phillips, H.-P. Bunge, and K. Schaber. True polar wander in mantle convection models with multiple, mobile continents. *Gondwana Research*, 15:288–296, 2009.
- [447] E. Pichelin and T. Coupez. Finite element solution of the 3D mold filling problem for viscous incompressible fluid. *Computer Methods in Applied Mechanics and Engineering*, 163:359–371, 1998.
- [448] A. Pinelli and A. Vacca. Chebyshev collocation method and multidomain decomposition for the incompressible Navier-Stokes equations. *International Journal for numerical methods in fluids*, 18:781–799, 1994.
- [449] C. Piromallo, T.W. Becker, F. Funiciello, and C. Faccenna. Three-dimensional instantaneous mantle flow induced by subduction. *Geophys. Res. Lett.*, 33(L08304), 2006.
- [450] J. Pitkäranta and T. Saarinen. A Multigrid Version of a Simple Finite Element Method for the Stokes Problem. *Mathematics of Computation*, 45(171):1–14, 1985.
- [451] T. Plank and P.E. van Keken. The ups and downs of sediments. *Nature Geoscience*, 1:17, 2008.
- [452] A. Plunder, C. Thieulot, and D.J.J. van Hinsbergen. The effect of obliquity on temperature in subduction zones: insights from 3D numerical modeling. *Solid Earth*, 9:759–776, 2018.
- [453] 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*, pages 175–195. Kluwer Academic Publishers, 1993.
- [454] A. Poliakov and Y. Podlachikov. Diapirism and topography. *Geophy. J. Int.*, 109:553–564, 1992.

- [455] A.A. Popov and S.V. Sobolev. SLIM3D: a tool for three-dimensional thermomechanical modelling of lithospheric deformation with elasto-visco-plastic rheology. *Phys. Earth. Planet. Inter.*, 171(1):55–75, 2008.
- [456] A. Prosperetti. Motion of two superposed viscous fluids. *Phys. Fluids*, 24(7):1217–1223, 1981.
- [457] 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. *Physics of the Earth and Planetary Interiors*, 276:10–35, 2018.
- [458] A.E. Pusok and B.J.P. Kaus. Development of topography in 3D continental collision models. *Geochem. Geophys. Geosyst.*, page doi:10.1002/2015GC005732, 2015.
- [459] 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. *Pure and Applied Geophysics*, pages doi:10.1007/s00024-016-1431-8, 2016.
- [460] C. Pütthe and T. Gerya. Dependence of mid-ocean ridge morphology on spreading rate in numerical 3-D models. *Gondwana Research*, 25:270–283, 2014.
- [461] R.N. Pysklywec. Surface erosion control on the evolution of the deep lithosphere. *Geology*, 34:225–228, 2006.
- [462] R.N. Pysklywec and C. Beaumont. Intraplate tectonics: feedback between radioactive thermal weakening and crustal deformation driven by mantle lithosphere instabilities. *Earth Planet. Sci. Lett.*, 221:275–292, 2004.
- [463] R.N. Pysklywec, C. Beaumont, and P. Fullsack. Modeling the behavior of continental mantle lithosphere during plate convergence. *Geology*, 28(7):655–658, 2000.
- [464] 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. *J. Geophys. Res.*, 107(B72133), 2002.
- [465] R.N. Pysklywec and A.R. Cruden. Coupled crust-mantle dynamics and intraplate tectonics: Two-dimensional numerical and three-dimensional analogue modeling. *Geochem. Geophys. Geosyst.*, 5(10), 2004.
- [466] R.N. Pysklywec, S.M. Ellis, and A.R. Gorman. Three-dimensional mantle lithosphere deformation at collisional plate boundaries: A subduction scissor across the South Island of New Zealand. *Earth Planet. Sci. Lett.*, 289:334–346, 2010.
- [467] R.N. Pysklywec, O. Gogus, J. Percival, A.R. Cruden, and C. Beaumont. Insights from geodynamical modeling on possible fates of continental mantle lithosphere: collision, removal, and overturn. *Can. J. Earth Sci.*, 47:541–563, 2010.
- [468] R.N. Pysklywec, J.X. Mitrovica, and M. Ishii. Mantle avalanche as a driving force for tectonic reorganization in the southwest Pacific. *Earth Planet. Sci. Lett.*, 209:29–38, 2003.
- [469] S. Quenette, Y. Xi, J. Mansour, L. Moresi, and D. Abramson. Underworld-GT Applied to Guangdong, a Tool to Explore the Geothermal Potential of the Crust. *Journal of Earth Sciences*, 26(1):78–88, 2015.
- [470] S. Quere, J.P. Lowman, J. Arkani-Hamed, J.H. Roberts, and R. Moucha. Subcontinental sinking slab remnants in a spherical geometry mantle model. *J. Geophys. Res.*, 118:1760–1777, 2013.
- [471] Matthieu E.T. Quinquis, Suzanne J.H. Buitier, and Susan Ellis. The role of boundary conditions in numerical models of subduction zone dynamics. *Tectonophysics*, 497:57–70, 2011.
- [472] M.E.T. Quinquis and S.J.H. Buitier. Testing the effects of basic numerical implementations of water migration on models of subduction dynamics. *Solid Earth*, 5:537–555, 2014.

- [473] J. Quinteros, V.A. Ramos, and P.M. Jacobkis. An elasto-visco-plastic model using the finite element method for crustal and lithospheric deformation. *Journal of Geodynamics*, 48:83–94, 2009.
- [474] J. Quinteros, S.V. Sobolev, and A.A. Popov. Viscosity in transition zone and lower mantle: Implications for slab penetration . *Geophys. Res. Lett.*, 37(L09307,), 2010.
- [475] T. Rabczuk, P.M.A. Areias, and T. Belytschko. A simplified mesh-free method for shear bands with cohesive surfaces . *Int. J. Num. Meth. Eng.*, 69:993–1021, 2007.
- [476] J.N. Reddy. On penalty function methods in the finite element analysis of flow problems. *Int. J. Num. Meth. Fluids*, 2:151–171, 1982.
- [477] J.N. Reddy and D.K. Gartling. *The Finite Element Method in Heat Transfer and Fluid Dynamics*. CRC Press, 2010.
- [478] J. Revenaugh and B. Parsons. Dynamic topography and gravity anomalies for fluid layers whose viscosity varies exponentially with depth. *Geophysical Journal of the Royal Astronomical Society*, 90(2):349–368, 1987.
- [479] J.G. Rice and R.J. Schnipke. An equal-order velocity-pressure formulation that does not exhibit spurious pressure modes. *Computer Methods in Applied Mechanics and Engineering*, 58:135–149, 1986.
- [480] J. Ritsema, A.K. McNamara, and A.L. Bull. Tomographic filtering of geodynamic models: Implications for model interpretation and large-scale mantle structure. *J. Geophys. Res.*, 112(B01303), 2007.
- [481] M. Roda, A.M. Marotta, and M.I. Spalla. Numerical simulations of an oceancontinent convergent system: Influence of subduction geometry and mantle wedge hydration on crustal recycling. *Geochem. Geophys. Geosyst.*, 11(5):10.1029/2009GC003015, 2010.
- [482] T. Rolf and P.J. Tackley. Focussing of stress by continents in 3D spherical mantle convection with selfconsistent plate tectonics. *Geophys. Res. Lett.*, 38(L18301):10.1029/2011GL048677, 2011.
- [483] I. Rose, B. Buffet, and T. Heister. Stability and accuracy of free surface time integration in viscous flows. *Phys. Earth. Planet. Inter.*, 262:90–100, 2017.
- [484] I. Rose and B. Buffet. Scaling rates of true polar wander in convecting planets and moons. *Physics of the Earth and Planetary Interiors*, 2017.
- [485] J.B. Ruh, T. Gerya, and J.-P. Burg. High-resolution 3D numerical modeling of thrust wedges: Influence of décollement strength on transfer zones. *Geochem. Geophys. Geosyst.*, 14(4):1131–1155, 2014.
- [486] J.B. Ruh, L. Le Pourhiet, Ph. Agard, E. Burov, and T. Gerya. Tectonic slicing of subducting oceanic crust along plate interfaces: Numerical modeling. *Geochem. Geophys. Geosyst.*, 16:10.1002/2015GC005998, 2015.
- [487] Y. Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.
- [488] Y. Saad and M.H. Schultz. GMRES: A Generalized Minimal Residual Algorithm for Solving Non-symmetric Linear Systems. *SIAM Journal on Scientific and Statistical Computing*, 7(3):856–869, 1986.
- [489] H. Samuel and M. Evonuk. Modeling advection in geophysical flows with particle level sets. *Geochem. Geophys. Geosyst.*, 11(8):doi:10.1029/2010GC003081, 2010.
- [490] D. Sandiford and L. Moresi. Improving subduction interface implementation in dynamic numerical models. *Solid Earth*, 2019.
- [491] 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. *Int. J. Num. Meth. Fluids*, 1:17–43, 1981.

- [492] 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. *Int. J. Num. Meth. Fluids*, 1:171–204, 1981.
- [493] W.P. Schellart, J. Freeman, D.R. Stegman, L. Moresi, and D. May. Evolution and diversity of subduction zones controlled by slab width. *Nature*, 446:308–311, 2007.
- [494] W.P. Schellart and L. Moresi. A new driving mechanism for backarc extension and backarc shortening through slab sinking induced toroidal and poloidal mantle flow: Results from dynamic subduction models with an overriding plate. *J. Geophys. Res.*, 118:1–28, 2013.
- [495] W.P. Schellart and W. Spakman. Australian plate motion and topography linked to fossil New Guinea slab below Lake Eyre. *Earth Planet. Sci. Lett.*, 421:107–116, 2015.
- [496] S.M. Schmalholz. A simple analytical solution for slab detachment. *Earth Planet. Sci. Lett.*, 304:45–54, 2011.
- [497] H. Schmeling, A.Y. Babeyko, A. Enns, C. Faccenna, F. Funiciello, T. Gerya, G.J. Golabek, S. Grigull, B.J.P. Kaus, G. Morra, S.M. Schmalholz, and J. van Hunen. A benchmark comparison of spontaneous subduction models - Towards a free surface. *Phys. Earth. Planet. Inter.*, 171:198–223, 2008.
- [498] D.W. Schmid and Y.Y. Podlachikov. Analytical solutions for deformable elliptical inclusions in general shear. *Geophy. J. Int.*, 155:269–288, 2003.
- [499] 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. *Numerical Heat Transfer*, 1:433–451, 1978.
- [500] G. Schubert, D.L. Turcotte, and P. Olson. *Mantle Convection in the Earth and Planets*. Cambridge University Press, 2001.
- [501] A. Segal. *Finite element methods for the incompressible Navier-Stokes equations*. 2012.
- [502] A. Segal and N.P. Praagman. The sepran fem package. Technical report, Technical Report, Ingenieursbureau Sepra, The Netherlands. <http://ta.twi.tudelft.nl/sepran/sepran.html>, 2005.
- [503] P. Sekhar and S.D. King. 3D spherical models of Martian mantle convection constrained by melting history. *Earth Planet. Sci. Lett.*, 388:27–37, 2014.
- [504] C. Selzer, S.J.H. Buijer, and O.A. Pfiffner. Numerical modeling of frontal and basal accretion at collisional margins. *Tectonics*, 27(TC3001):doi:10.1029/2007TC002169, 2008.
- [505] Cornelia Selzer. *Tectonic accretion styles at convergent margins: A numerical modelling study*. PhD thesis, University of Bern, 2006.
- [506] M. Seton, N. Flament, J. Whittaker, R.D. Müller, M. Gurnis, and D.J. Bower. Ridge subduction sparked reorganization of the Pacific plate-mantle system 6050 million years ago. *Geophys. Res. Lett.*, 42:doi:10.1002/2015GL063057, 2015.
- [507] P.N. Shankar and M.D. Deshpande. Fluid mechanics in the driven cavity. *Annu. Rev. Fluid Mech.*, 32:93–136, 2000.
- [508] W. Sharples, L.-N. Moresi, and M. A. Jadamec andJ. Revote. Styles of rifting and fault spacing in numerical models of crustal extension. *J. Geophys. Res.*, 120:4379–4404, 2015.
- [509] 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. *Phys. Earth. Planet. Inter.*, 252:77–90, 2016.
- [510] K. Simon, R.S. Huismans, and C. Beaumont. Dynamical modelling of lithospheric extension and small-scale convection: implications for magmatism during the formation of volcanic rifted margins. *Geophy. J. Int.*, 176:327–350, 2009.

- [511] S.V. Sobolev, A. Petrunin, Z. Garfunkel, and A.Y. Babeyko. Thermo-mechanical model of the Dead Sea Transform. *Earth Planet. Sci. Lett.*, 238:78–95, 2005.
- [512] S.V. Sobolev, A.V. Sobolev, D.V. Kuzmin, N.A. Krivolutskaya, A.G. Petrunin, N.T. Arndt, V.A. Radko, and Y.R. Vasiliev. Linking mantle plumes, large igneous provinces and environmental catastrophes. *Nature*, 477:312, 2011.
- [513] F. Soboutia, A. Ghodsib, and J. Arkani-Hamed. On the advection of sharp material interfaces in geodynamic problems: entrainment of the D” layer. *Journal of Geodynamics*, 31:459–479, 2001.
- [514] V.S. Solomatov. Initiation of subduction by small-scale convection. *J. Geophys. Res.*, 109(B01412), 2004.
- [515] V.S. Solomatov. Localized subcritical convective cells in temperature-dependent viscosity fluids. *Phys. Earth. Planet. Inter.*, 200-201:63–71, 2012.
- [516] V.S. Solomatov and L.-N. Moresi. Stagnant lid convection on Venus. *J. Geophys. Res.*, 101(E2):4737–4753, 1996.
- [517] M. Spiegelman, D.A. May, and C. Wilson. On the solvability of incompressible Stokes with viscoplastic rheologies in geodynamics. *Geochem. Geophys. Geosyst.*, 17:2213–2238, 2016.
- [518] 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. *Science*, 329:1033–1038, 2010.
- [519] D.R. Stegman, R. Farrington, F.A. Capitanio, and W.P. Schellart. A regime diagram for subduction styles from 3-D numerical models of free subduction. *Tectonophysics*, 483:29–45, 2010.
- [520] D.R. Stegman, J. Freeman, W.P. Schellart, L. Moresi, and D. May. Influence of trench width on subduction hinge retreat rates in 3-D models of slab rollback. *Geochem. Geophys. Geosyst.*, 7(3):doi:10.1029/2005GC001056, 2006.
- [521] D.R. Stegman, W.P. Schellart, and J. Freeman. Competing influences of plate width and far-field boundary conditions on trench migration and morphology of subducted slabs in the upper mantle. *Tectonophysics*, 483:46–57, 2010.
- [522] C. Stein, J. Lowman, and U. Hansen. A comparison of mantle convection models featuring plates. *Geochem. Geophys. Geosyst.*, 15:2689–2698, 2014.
- [523] B. Steinberger, E. Bredow, S. Lebedev, A. Schaeffer, and T. H. Torsvik. Widespread volcanism in the greenland-north atlantic region explained by the iceland plume. *Nature Geoscience*, 12(1):61, 2019.
- [524] P. Sternai, L. Jolivet, A. Menant, and T. Gerya. Driving the upper plate surface deformation by slab rollback and mantle flow. *Earth Planet. Sci. Lett.*, 405:110–118, 2014.
- [525] J. Suckale, B.H. Hager, L.T. ElkinsTanton, and J.Ch. Nave. It takes three to tango: 2. Bubble dynamics in basaltic volcanoes and ramifications for modeling normal Strombolian activity. *J. Geophys. Res.*, 115(B7), 2010.
- [526] 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. *J. Geophys. Res.*, 115(B07409), 2010.
- [527] E.M. Syracuse, P.E. van Keeken, and G.A. Abers. The global range of subduction zone thermal models. *Phys. Earth. Planet. Inter.*, 183:73–90, 2010.
- [528] P. Tackley. *Three-dimensional models of mantle convection: Influence of phase transitions and temperature-dependent viscosity*. PhD thesis, California Institute of Technology, 1994.
- [529] P.J. Tackley and S.D. King. Testing the tracer ratio method for modeling active compositional fields in mantle convection simulations. *Geochem. Geophys. Geosyst.*, 4(4), 2003.

- [530] K. Takeyama, T. R. Saitoh, and J. Makino. Variable inertia method: A novel numerical method for mantle convection simulation. *New Astronomy*, 2017.
- [531] E. Tan, E. Choi, P. Thoutireddy, M. Gurnis, and M. Aivazis. GeoFramework: Coupling multiple models of mantle convection within a computational framework. *Geochem. Geophys. Geosyst.*, 7(6):10.1029/2005GC001155, 2006.
- [532] E. Tan and M. Gurnis. Compressible thermochemical convection and application to lower mantle structures. *J. Geophys. Res.*, 112(B06304), 2007.
- [533] E. Tan, M. Gurnis, and L. Han. Slabs in the lower mantle and their modulation of plume formation. *Geochem. Geophys. Geosyst.*, 3(11), 2002.
- [534] Paul Tapponnier and Peter Molnar. Slip-line field theory and large-scale continental tectonics. *Nature*, 264:319–324, November 1976.
- [535] J.M. Taramon, J. Rodriguez-Gonzalez, A.M. Negredo, and M.I. Billen. Influence of cratonic lithosphere on the formation and evolution of flat slabs: Insights from 3-D time-dependent modeling. *Geochem. Geophys. Geosyst.*, 16:doi:10.1002/2015GC005940, 2015.
- [536] C. Tayloor and P. Hood. A numerical solution of the Navier-Stokes equations using the finite element technique. *Comput. Fluids*, 1:73–100, 1973.
- [537] J.L. Tetraeault and S.J.H. Buiter. Geodynamic models of terrane accretion: Testing the fate of island arcs, oceanic plateaus, and continental fragments in subduction zones. *J. Geophys. Res.*, 2012.
- [538] J.L. Tetraeault and S.J.H. Buiter. The influence of extension rate and crustal rheology on the evolution of passive margins from rifting to break-up. *Tectonophysics*, 2017.
- [539] T.E. Tezduyar, S. Mittal, S.E. Ray, and R. Shih. Incompressible flow computations with stabilized bilinear and linear equal-order-interpolation velocity-pressure elements. *Comput. Methods Appl. Mech. Engrg.*, 95:221–242, 1992.
- [540] M. Thielmann, B.J.P. Kaus, and A.A. Popov. Lithospheric stresses in RayleighBénard convection: effects of a free surface and a viscoelastic Maxwell rheology. *Geophys. J. Int.*, 203:2200–2219, 2015.
- [541] M. Thielmann, D.A. May, and B.J.P. Kaus. Discretization errors in the Hybrid Finite Element Particle-In-Cell Method. *Pure and Applied Geophysics*, 171(9):2164–2184, 2014.
- [542] C. Thieulot. FANTOM: two- and three-dimensional numerical modelling of creeping flows for the solution of geological problems. *Phys. Earth. Planet. Inter.*, 188(1):47–68, 2011.
- [543] C. Thieulot. Analytical solution for viscous incompressible stokes flow in a spherical shell. *Solid Earth Discussions*, 2017:1–19, 2017.
- [544] C. Thieulot. GHOST: Geoscientific Hollow Sphere Tesselation. *Solid Earth*, 9(1–9), 2018.
- [545] C. Thieulot and W. Bangerth. On the use of equal order elements in geodynamics. *Solid Earth*, 2019.
- [546] C. Thieulot, P. Fullsack, and J. Braun. Adaptive octree-based finite element analysis of two- and three-dimensional indentation problems. *J. Geophys. Res.*, 113:B12207, 2008.
- [547] C. Thieulot, P. Steer, and R.S. Huismans. Three-dimensional numerical simulations of crustal systems undergoing orogeny and subjected to surface processes. *Geochem. Geophys. Geosyst.*, 15, 2014.
- [548] J.F. Thompson, B.K. Soni, and N.P. Weatherill. *Handbook of grid generation*. CRC press, 1998.
- [549] N. Tosi, P. Maierová, and D.A. Yuen. Influence of Variable Thermal Expansivity and Conductivity on Deep Subduction. In *Subduction Dynamics: From Mantle Flow to Mega Disasters, Geophysical Monograph 211*, pages 115–133. John Wiley & Sons, Inc., 2016.

- [550] N. Tosi, C. Stein, L. Noack, C. Huettig, P. Maierova, H. Samuel, D.R. Davies, C.R. Wilson, S.C. Kramer, C. Thieulot, A. Glerum, M. Fraters, W. Spakman, A. Rozel, and P.J. Tackley. A community benchmark for viscoplastic thermal convection in a 2-D square box. *Geochem. Geophys. Geosyst.*, 16(7):21752196, 2015.
- [551] G. Toussaint, E. Burov, and J.-P. Avouac. Tectonic evolution of a continental collision zone: A thermomechanical numerical model. *Tectonics*, 23(TC6003):doi:10.1029/2003TC001604, 2004.
- [552] B.J. Travis, C. Anderson, J. Baumgardner, C.W. Gable, B.H. Hager, R.J. O'Connell, P. Olson, A. Raefsky, and G. Schubert. A benchmark comparison of numerical methods for infinite Prandtl number thermal convection in two-dimensional Cartesian geometry. *Geophysical & Astrophysical Fluid Dynamics*, 55(3-4):137–160, 1990.
- [553] R.A. Trompert and U. Hansen. On the Rayleigh number dependence of convection with a strongly temperature-dependent viscosity. *Physics of Fluids*, 10(2):351–360, 1998.
- [554] D.L. Turcotte and G. Schubert. *Geodynamics, 2nd edition*. Cambridge, 2012.
- [555] K. Ueda, T. Gerya, and S.V. Sobolev. Subduction initiation by thermalchemical plumes: Numerical studies. *Phys. Earth. Planet. Inter.*, 171:296–312, 2008.
- [556] K. Ueda, S.D. Willett, T. Gerya, and J. Ruh. Geomorphologicalthermo-mechanical modeling: Application to orogenic wedge dynamics. *Tectonophysics*, 659:12–30, 2015.
- [557] A. van den Berg, P.E. van Keken, and D.A. Yuen. The effects of a composite non-Newtonian and Newtonian rheology on mantle convection. *Geophys. J. Int.*, 115:62–78, 1993.
- [558] A.P. van den Berg, M.V. De Hoop, D.A. Yuen, A. Duchkov, R.D. van der Hilst, and M.H.G. Jacobs. Geodynamical modeling and multiscale seismic expression of thermo-chemical heterogeneity and phase transitions in the lowermost mantle. *Phys. Earth. Planet. Inter.*, 180:244–257, 2010.
- [559] A.P. van den Berg, G. Segal, and D. Yuen. SEPRAN: A Versatile Finite-Element Package for Realistic Problems in Geosciences. *International Workshop of Deep Geothermal Systems, Wuhan, China, June 29-30*, 2012.
- [560] A.P. van den Berg, G. Segal, and D.A. Yuen. SEPRAN: A Versatile Finite-Element Package for a Wide Variety of Problems in Geosciences. *Journal of Earth Science*, 26(1):089–095, 2015.
- [561] Y. van Dinther, T.V. Gerya, L.A. Dalguer, F. Corbi, F. Funiciello, and P.M. Mai. The seismic cycle at subduction thrusts: 2. Dynamic implications of geodynamic simulations validated with laboratory models. *J. Geophys. Res.*, 118:1502–1525, 2013.
- [562] Y. van Dinther, T.V. Gerya, L.A. Dalguer, P.M. Mai, G. Morra, and D. Giardini. The seismic cycle at subduction thrusts: Insights from seismo-thermo-mechanical models. *J. Geophys. Res.*, 118:1–20, 2013.
- [563] Y. van Dinther, P.M. Mai, L.A. Dalguer, and T.V. Gerya. Modeling the seismic cycle in subduction zones: The role and spatiotemporal occurrence of off-megathrust earthquakes. *Geophys. Res. Lett.*, 41:1194–1201, 2014.
- [564] H.J. van Heck, J.H. Davies, T. Elliott, and D. Porcelli. Global-scale modelling of melting and isotopic evolution of Earths mantle: melting modules for TERRA. *Geosci. Model Dev.*, 9:1399–1411, 2016.
- [565] J. van Hunen and M.B. Allen. Continental collision and slab break-off: A comparison of 3-D numerical models with observations. *Earth Planet. Sci. Lett.*, 302:27–37, 2011.
- [566] J. van Hunen and M.S. Miller. Collisional processes and links to episodic changes in subduction zones. *Elements*, 11:119–124, 2015.
- [567] 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. *Tectonophysics*, 352(3-4):317–333, 2002.

- [568] J. van Hunen and S. Zhong. New insight in the Hawaiian plume swell dynamics from scaling laws. *Geophys. Res. Lett.*, 30(15):doi:10.1029/2003GL017646,, 2003.
- [569] P. van Keken and S. Zhong. Mixing in a 3D spherical model of present-day mantle convection. *Earth Planet. Sci. Lett.*, 171:533–547, 1999.
- [570] P.E. van Keken, C. Currie, S.D. King, M.D. Behn, Amandine Cagnioncle, J. Hee, R.F. Katz, S.-C. Lin, E.M. Parmentier, M. Spiegelman, and K. Wang. A community benchmark for subduction zone modelling. *Phys. Earth. Planet. Inter.*, 171:187–197, 2008.
- [571] P.E. van Keken, B.R. Hacker, E.M. Syracuse, and G.A. Abers. Subduction factory: 4. Depth-dependent flux of H<sub>2</sub>O from subducting slabs worldwide. *J. Geophys. Res.*, 116(B01401), 2011.
- [572] 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. *J. Geophys. Res.*, 102(B10):22,477–22,495, 1997.
- [573] P. van Thienen. Convective vigour and heat flow in chemically differentiated systems. *Geophys. J. Int.*, 169(2):747–766, 2007.
- [574] P. van Thienen, A.P. van den Berg, J.H. de Smet, J. van Hunen, and M.R. Drury. Interaction between small-scale mantle diapirs and a continental root. *Geochem. Geophys. Geosyst.*, 4(2):doi:10.1029/2002GC000338, 2003.
- [575] P. van Thienen, A.P. van den Berg, and N.J. Vlaar. On the formation of continental silicic melts in thermochemical mantle convection models: implications for early Earth. *Tectonophysics*, 394(1-2):111–124, 2004.
- [576] P. van Thienen, N.J. Vlaar, and A.P. van den Berg. Plate tectonics on the terrestrial planets. *Phys. Earth. Planet. Inter.*, 142:61–74, 2004.
- [577] P. van Thienen, N.J. Vlaar, and A.P. van den Berg. Assessment of the cooling capacity of plate tectonics and flood volcanism in the evolution of Earth, Mars and Venus. *Phys. Earth. Planet. Inter.*, 150:287–315, 2005.
- [578] J.W. van Wijk, W.S. Baldridge, J. van Hunen, S. Goes, R. Aster, D.D. Coblenz, S.P. Grand, and J. Ni. Small-scale convection at the edge of the Colorado Plateau: Implications for topography, magmatism, and evolution of Proterozoic lithosphere. *Geology*, 38:611–614, 2010.
- [579] O. Vanderhaeghe, S. Medvedev, P. Fullsack, C. Beaumont, and R.A. Jamieson. Evolution of orogenic wedges and continental plateaux: insights from crustal thermal-mechanical overlying subduction mantle lithosphere. *Geophys. J. Int.*, 153:27–51, 2003.
- [580] R.S. Varga. *Matrix Iterative Analysis*. Prentice-Hall, Inc., 1963.
- [581] J. Vatteville, P.E. van Keken, A. Limare, and A. Davaille. Starting laminar plumes: Comparison of laboratory and numerical modeling. *Geochem. Geophys. Geosyst.*, 10(12):doi:10.1029/2009GC002739, 2009.
- [582] Ph. Vernant and J. Chery. Mechanical modelling of oblique convergence in the Zagros, Iran. *Geophys. J. Int.*, 165:991–1002, 2006.
- [583] 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. *Earth Planet. Sci. Lett.*, 121:1–18, 1994.
- [584] K. Vogt and T. Gerya. Deep plate serpentinization triggers skinning of subducting slabs. *Geology*, page doi:10.1130/G35565.1, 2014.
- [585] M. von Tscharner and S. M. Schmalholz. A 3-D Lagrangian finite element algorithm with remeshing for simulating large-strain hydrodynamic instabilities in power law viscoelastic fluids. *Geochem. Geophys. Geosyst.*, 16:215–245, 2015.

- [586] L. Vynnytska, M.E. Rognes, and S.R. Clark. Benchmarking FEniCS for mantle convection simulations. *Computers & Geosciences*, 50:95–105, 2013.
- [587] 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. *Geochem. Geophys. Geosyst.*, 16:doi:10.1002/2015GC005824, 2015.
- [588] H. Wang, J. van Hunen, , and D.G. Pearson. The thinning of subcontinental lithosphere: The roles of plume impact and metasomatic weakening. *Geochem. Geophys. Geosyst.*, 16:1156–1171, 2015.
- [589] H. Wang, J. van Hunen, D.G. Pearson, and M.B. Allen. Craton stability and longevity: The roles of composition-dependent rheology and buoyancy. *Earth Planet. Sci. Lett.*, 391:224–233, 2014.
- [590] X. Wang, J. He, L. Ding, and R. Gao. A possible mechanism for the initiation of the Yinggehai Basin: A visco-elasto-plastic model. *Journal of Asian Earth Sciences*, 74:25–36, 2013.
- [591] Y. Wang, J. Huang, and S. Zhong. Episodic and multistaged gravitational instability of cratonic lithosphere and its implications for reactivation of the North China Craton. *Geochem. Geophys. Geosyst.*, 16:815–833, 2015.
- [592] 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. *Geochem. Geophys. Geosyst.*, 9, 2008.
- [593] 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. *Earth Planet. Sci. Lett.*, 267:129–145, 2008.
- [594] M.B. Weller and A. Lenardic. The energetics and convective vigor of mixed-mode heating: Velocity scalings and implications for the tectonics of exoplanets. *Geophys. Res. Lett.*, 43, 2016.
- [595] M.B. Weller, A. Lenardic, and W.B. Moore. Scaling relationships and physics for mixed heating convection in planetary interiors: Isoviscous spherical shells. *J. Geophys. Res.*, 121, 2016.
- [596] D.M. Whipp, C. Beaumont, and J. Braun. Feeding the aneurysm: Orogen-parallel mass transport into Nanga Parbat and the western Himalayan syntaxis. *J. Geophys. Res.*, 119:doi:10.1002/2013JB010929, 2014.
- [597] S.D. Willett. Dynamic and kinematic growth and change of a coulomb wedge. In K.R. McClay, editor, *Thrust Tectonics*, pages 19–31. Chapman and Hall, 1992.
- [598] S.D. Willett. Orogeny and orography: The effects of erosion on the structure of mountain belts. *J. Geophys. Res.*, 104(B12):28957, 1999.
- [599] S.D. Willett and C. Beaumont. Subduction of Asian lithosphere mantle beneath Tibet inferred from models of continental collision. *Nature*, 369:642–645, 1994.
- [600] Sean D. Willett and Daniel C. Pope. Thermo-mechanical models of convergent orogenesis: Thermal and rheologic dependence of crustal deformation. In *Rheology and deformation of the lithosphere at continental margins.*, pages 166–222. Columbia University Press, 2003.
- [601] M. Wolstencroft and J.H. Davies. Influence of the Ringwoodite-Perovskite transition on mantle convection in spherical geometry as a function of Clapeyron slope and Rayleigh number. *Solid Earth*, 2:315–326, 2011.
- [602] M. Wolstencroft, J.H. Davies, and D.R. Davies. NusseltRayleigh number scaling for spherical shell Earth mantle simulation up to a Rayleigh number of  $10^9$ . *Phys. Earth. Planet. Inter.*, 176:132–141, 2009.
- [603] G. Wu, L.L. Lavier, and E. Choi. Modes of continental extension in a crustal wedge. *Earth Planet. Sci. Lett.*, 421:89–97, 2015.

- [604] H. Xing, W. Yu, and J. Zhang. In *Advances in Geocomputing, Lecture Notes in Earth Sciences*. Springer-Verlag, Berlin Heidelberg, 2009.
- [605] P. Yamato, P. Agard, E. Burov, L. Le Pourhiet, L. Jolivet, and C. Tiberi. Burial and exhumation in a subduction wedge: Mutual constraints from thermomechanical modeling and natural P-T-t data (Schistes Lustres, western Alps). *J. Geophys. Res.*, 112(B07410):doi:10.1029/2006JB004441, 2007.
- [606] P. Yamato, E. Burov, P. Agard, L. Le Pourhiet, and L. Jolivet. HP-UHP exhumation during slow continental subduction: Self-consistent thermodynamically and thermomechanically coupled model with application to the Western Alps. *Earth Planet. Sci. Lett.*, 271:63–74, 2008.
- [607] P. Yamato, L. Husson, T.W. Becker, and K. Pedoja. Passive margins getting squeezed in the mantle convection vice. *Tectonics*, 2013.
- [608] P. Yamato, L. Husson, J. Braun, C. Loiselet, and C. Thieulot. Influence of surrounding plates on 3D subduction dynamics. *Geophys. Res. Lett.*, 36(L07303), 2009.
- [609] P. Yamato, R. Tartese, T. Duretz, and D.A. May. Numerical modelling of magma transport in dykes . *Tectonophysics*, 526-529:97–109, 2012.
- [610] Ph. Yamato, B.J.P. Kaus, F. Mouthereau, and S. Castellort. Dynamic constraints on the crustal-scale rheology of the Zagros fold belt, Iran. *Geology*, 39(9):815–818, 2011.
- [611] S. YANG and Y. SHI. Three-dimensional numerical simulation of glacial trough forming process. *Science China: Earth Sciences*, pages 10.1007/s11430-015-5120-8, 2015.
- [612] T. Yang, L. Moresi, M. Gurnis, S. Liu, D. Sandiford, S. Williams, and F.A. Capitanio. Contrasted East Asia and South America tectonics driven by deep mantle flow. *Earth Planet. Sci. Lett.*, 517:106–116, 2019.
- [613] C. Yao, F. Deschamps, J.P. Lowman, C. Sanchez-Valle, and P.J. Tackley. Stagnant lid convection in bottom-heated thin 3-D spherical shells: Influence of curvature and implications for dwarf planets and icy moons. *J. Geophys. Res.*, 119:1895–1913, 2014.
- [614] M. Yoshida. The role of harzburgite layers in the morphology of subducting plates and the behavior of oceanic crustal layers. *Geophys. Res. Lett.*, 40:5387–5392, 2013.
- [615] M. Yoshida, F. Tajima, S. Honda, and M. Morishige. The 3D numerical modeling of subduction dynamics: Plate stagnation and segmentation, and crustal advection in the wet mantle transition zone. *J. Geophys. Res.*, 117(B04104), 2012.
- [616] X. Yu and F. Tin-Loi. A simple mixed finite element for static limit analysis. *Computers and Structures*, 84:1906–1917, 2006.
- [617] Z. Xu Z. Li and T.V. Gerya. Numerical Geodynamic Modeling of Continental Convergent Margins, Earth Sciences. In Dr. Imran Ahmad Dar, editor, *Earth Sciences*. InTech, 2012.
- [618] N. Zhang and Z-X Li. Formation of mantle “lone plumes” in the global downwelling zone – a case for subduction-controlled plume generation beneath the south china sea. *Tectonophysics*, 2017.
- [619] N. Zhang, S. Zhong, W. Leng, , and Z.X. Li. A model for the evolution of the Earths mantle structure since the Early Paleozoic. *J. Geophys. Res.*, 115(B06401), 2010.
- [620] N. Zhang, S. Zhong, and A.K. McNamara. Supercontinent formation from stochastic collision and mantle convection models. *Gondwana Research*, 15:267–275, 2009.
- [621] S. Zhang and C. O’Neill. The early geodynamic evolution of mars-type planets. *Icarus*, 265:187–208, 2016.
- [622] Y. Zhao, J. de Vries, A.P. van den Berg, M.H.G. Jacobs, and W. van Westrenen. The participation of ilmenite-bearing cumulates in lunar mantle overturn. *Earth Planet. Sci. Lett.*, 511:1–11, 2019.

- [623] S. Zhong, M. Gurnis, and G. Hulbert. Accurate determination of surface normal stress in viscous flow from a consistent boundary flux method. *Phys. Earth. Planet. Inter.*, 78:1–8, 1993.
- [624] S. Zhong, M. Gurnis, and L. Moresi. Role of faults, nonlinear rheology, and viscosity structure in generating plates from instantaneous mantle flow models. *J. Geophys. Res.*, 103(B7):15,255–15,268, 1998.
- [625] 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. *Geochem. Geophys. Geosyst.*, 9(10), 2008.
- [626] S. Zhong, N. Zhang, Z.-X. Li, and J.H. Roberts. Supercontinent cycles, true polar wander, and very long-wavelength mantle convection. *Earth Planet. Sci. Lett.*, 261:551–564, 2007.
- [627] 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. *J. Geophys. Res.*, 105(B5):11,063–11,082, 2000.
- [628] Shijie Zhong. Analytic solutions for Stokes flow with lateral variations in viscosity. *Geophys. J. Int.*, 124(1):18–28, 1996.
- [629] S.J. Zhong, D.A. Yuen, L.N. Moresi, and M.G. Knepley. 7.05 - numerical methods for mantle convection. In Gerald Schubert, editor, *Treatise on Geophysics (Second Edition)*, pages 197 – 222. Elsevier, Oxford, second edition edition, 2015.
- [630] G. Zhu, T. Gerya, D.A. Yuen, S. Honda, T. Yoshida, and T. Connolly. Three-dimensional dynamics of hydrous thermal-chemical plumes in oceanic subduction zones. *Geochem. Geophys. Geosyst.*, 10:doi:10.1029/2009GC002625, 2009.
- [631] G. Zhu, T.V. Gerya, P.J. Tackley, and E. Kissling. Four-dimensional numerical modeling of crustal growth at active continental margins. *J. Geophys. Res.*, 118:4682–4698, 2013.
- [632] T. Zhu. Tomography-based mantle flow beneath Mongolia-Baikal area. *Phys. Earth. Planet. Inter.*, 237:40–50, 2014.
- [633] O. Zienkiewicz and S. Nakazawa. The penalty function method and its application to the numerical solution of boundary value problems. *The American Society of Mechanical Engineers*, 51, 1982.
- [634] O.C. Zienkiewicz, M. Huang, and M. Pastor. Localization problems in plasticity using finite elements with adaptive remeshing. *International Journal for Numerical and Analytical Methods in Geomechanics*, 19:127–148, 1995.
- [635] O.C. Zienkiewicz, C. Humpheson, and R.W. Lewis. Associated and non-associated visco-plasticity and plasticity in soil mechanics . *Géotechnique*, 25(4):671–689, 1975.
- [636] O.C. Zienkiewicz, J.P. Vilotte, and S. Toyoshima. Iterative method for constrained and mixed approximation. An inexpensive improvement of FEM performance. *Computer Methods in Applied Mechanics and Engineering*, 51:3–29, 1985.
- [637] S. Zlotnik, M. Fernandez, P. Diez, and J. Verges. Modelling gravitational instabilities: slab break-off and Rayleigh-Taylor diapirism. *Pure appl. geophys.*, 165:1491–1510, 2008.
- [638] F. Zwaan, G. Scheurs, J. Naliboff, and S.J.H. Buitier. Insights into the effects of oblique extension on continental rift interaction from 3D analogue and numerical models. *Tectonophysics*, 2016.

# Index

- $H^1$  norm, 128
- $H^1$  semi-norm, 128
- $H^1(\Omega)$  space, 128
- $L_1$  norm, 128
- $L_2$  norm, 128
- $P_1$ , 38, 44
- $P_1^+$ , 38, 45
- $P_2$ , 41
- $P_2^+$ , 42
- $P_3$ , 44
- $P_m \times P_n$ , 27
- $P_m \times P_{-n}$ , 28
- $Q_1 \times P_0$ , 62, 153, 159, 161, 167, 174, 178, 183, 185, 203, 210, 222, 224, 228, 243, 244
- $Q_1$ , 28, 33
- $Q_1^+ \times Q_1$ , 65
- $Q_2 \times P_0$ , 64
- $Q_2 \times Q_1$ , 28, 186, 191
- $Q_2$ , 29, 35, 47
- $Q_2^{(8)}$ , 36
- $Q_3 \times Q_2$ , 194
- $Q_3$ , 30, 37
- $Q_m \times P_{-n}$ , 28
- $Q_m \times Q_n$ , 28
- $Q_m \times Q_{-n}$ , 28
- , 293
- Advection-Diffusion, 55
- analytical solution, 153, 161, 178, 183, 191, 194, 197, 203, 210, 228, 244
- angular momentum, 144
- angular velocity, 144
- arithmetric mean, 175
- Augmented Lagrangian, 113
- Backward Euler, 58
- Barycentric Coordinates, 38
- basis functions, 33
- BDF-2, 59
- Biharmonic Operator, 150, 151
- Boussinesq, 15
- Bubble Function, 38
- bubble function, 28
- bulk modulus, 16, 205
- bulk viscosity, 14
- buoyancy-driven flow, 155
- CBF, 228
- CG, 109
- checkerboard, 103
- cohesion, 102
- compositional Field, 139
- Compressed Sparse Column, 92
- Compressed Sparse Row, 92
- compressibility, 205
- compressible flow, 210
- conforming element, 28
- conjugate gradient, 109
- connectivity array, 95
- convex polygon, 95
- Crank-Nicolson, 58
- Crouzeix-Raviart, 42
- CSC, 92
- CSR, 92
- Dirichlet boundary condition, 18
- divergence free, 61
- Domain Decomposition, 149
- Drucker-Prager, 102
- dynamic viscosity, 14
- ENO, 139
- Essential Boundary Conditions, 134
- Extrapolation, 129
- FANTOM, 293
- FLUIDITY, 293
- Forward Euler, 58
- Gauss quadrature, 23
- geometric mean, 175
- harmonic mean, 175
- Heun's emthod, 123
- hyperbolic PDE, 118
- IFISS, 293
- incompressible flow, 161, 167, 174, 178, 183, 185, 191, 194, 197, 203, 222, 224, 228, 243, 244
- isoparametric, 119
- isothermal, 153, 159, 161, 167, 174, 178, 183, 185, 191, 194, 197, 203, 222, 224, 228, 244
- isoviscous, 153, 159, 178, 183, 191, 194, 203, 210, 228, 243
- Lamé parameter, 16
- LBB, 61
- Legendre polynomial, 23
- Level-set Function, 139
- Level-set Method, 139
- LSF, 139
- LSM, 139
- MAC, 136
- Marker-and-Cell, 136
- method of manufactured solutions, 80
- midpoint method, 123

midpoint rule, 22  
 MINI element, 63  
 mixed formulation, 178, 183, 185, 191, 194, 197, 203, 210, 222, 224, 228, 243  
 MMS, 80  
 moment of inertia, 144  
  
 Natural Boundary Conditions, 134  
 Neumann boundary condition, 18  
 Newton's method, 126  
 Newton-Cotes, 23  
 non-conforming element, 28  
 non-isoviscous, 161, 167, 174, 185, 197, 244  
 nonconforming  $Q_1 \times P_0$ , 197  
 nonlinear, 147, 224  
 nonlinear rheology, 167  
 nullspace, 106  
 numerical benchmark, 222  
 Nusselt Number, 19  
  
 open boundary conditions, 243  
 optimal rate, 61, 129  
  
 Particle-in-Cell, 136  
 particle-in-cell, 174  
 penalty formulation, 65, 153, 155, 159, 161, 167, 174, 244  
 PIC, 136  
 Picard iterations, 147  
 piecewise, 27  
 Poiseuille flow, 84  
 Poisson ratio, 16  
 Prandtl Number, 19  
 preconditioned conjugate gradient, 112  
 pressure normalisation, 187  
 pressure nullspace, 86  
 pressure scaling, 105  
 pressure smoothing, 103, 178, 197, 203, 210, 228  
  
 quadrature, 23  
  
 Rayleigh Number, 19  
 rectangle rule, 22  
 relaxation, 147  
 Rens Elbertsen, 229, 250  
 Richardson iterations, 108  
 RK2, 123  
 RK3, 123  
 RK4, 123  
 RK45, 123  
 Runge-Kutta-Fehlberg method, 123  
  
 Schur complement, 108  
 Schur complement approach, 183, 185  
 second viscosity, 14  
 Shear Heating, 84  
 shear modulus, 16  
  
 solenoidal field, 61  
 SPD, 108  
 spin tensor, 12  
 Static Condensation, 63  
 static condensation, 141  
 Stokes sphere, 155, 185  
 strain rate, 12  
 Strain rate tensor, 19  
 Strain tensor, 19  
 strain tensor, 16  
 Stream Function, 150, 151  
 Stress Tensor, 19  
 strong form, 61  
 structured grid, 95  
 Symmetric Positive Definite, 295  
  
 Taylor-Hood, 61  
 tensor invariant, 101  
 thermal expansion, 205  
 trapezoidal rule, 22  
  
 unstructured grid, 95  
  
 velocity gradient, 12  
 Viscosity Rescaling Method, 102  
 von Mises, 102  
 VRM, 102  
  
 weak form, 61, 67  
 work against gravity, 208  
  
 Young's modulus, 16

## Notes

finish. not happy with definition. Look at literature . . . . .	20
derive formula for Earth size R1 and R2 . . . . .	20
check aspect manual The 2D cylindrical shell benchmarks by Davies et al. 5.4.12 . . . . .	20
insert here figure . . . . .	22
insert here figure . . . . .	22
citation needed . . . . .	35
verify those . . . . .	37
find reference . . . . .	42
VERIFY that when $\eta_1 = 1 - r - s$ , $\eta_2 = r$ and $\eta_3 = s$ we find the above $r, s$ shape functions . . . . .	43
verify those . . . . .	44
NOT happy with this statement!! . . . . .	48
add source term!! . . . . .	58
find literature . . . . .	61
Check Cuvelier book chapter 8 for modified element . . . . .	64
list codes which use this approach . . . . .	65
reduced integration . . . . .	68
write about 3D to 2D . . . . .	68
insert here link(s) to manual and literature . . . . .	68
sort out mess wrt Eq 26 of busa13 . . . . .	83
go through my papers and add relevant ones here . . . . .	88
produce drawing of node numbering . . . . .	96
insert here the rederivation 2.1.1 of spmw16 . . . . .	102
produce figure to explain this . . . . .	104
link to proto paper . . . . .	104
link to least square and nodal derivatives . . . . .	104
check this, and report page number . . . . .	109
add biblio . . . . .	110
how to compute $M$ for the Schur complement ? . . . . .	113
finish GMRES algo description. not sure what to do, hard to explain, not easy to code. . . . .	114
why does array contain only T?? . . . . .	116
finish . . . . .	116
write about case when element is rectangle/cuboid . . . . .	127
to finish . . . . .	131
to finish . . . . .	131
refs! . . . . .	135
tie to fieldstone 12 . . . . .	135
it would be nice to have a Q1 and Q2 drawing of a 1D element and show that indeed negative values arise . . . . .	138
write more about particle averaging and projection . . . . .	138
write about DG approach . . . . .	140
check . . . . .	141
implement and report . . . . .	142
missing picture . . . . .	143
explain why our eqs are nonlinear . . . . .	147
write about nonlinear residual . . . . .	147
FIND refs. check new version of Vol7 theoretical geophys . . . . .	151
build S and have python compute its smallest and largest eigenvalues as a function of resolution?183	183
finish structure of C matrix for q1q1 . . . . .	200
compare my rates with original paper! . . . . .	201
explain how Medge is arrived at! . . . . .	227
compare with ASPECT ??!. . . . .	227
gauss-lobatto integration? . . . . .	227
pressure average on surface instead of volume ? . . . . .	227
cross check with CIG database . . . . .	292