

## FINITE ELEMENT ERROR ANALYSIS OF A MANTLE CONVECTION MODEL

VOLKER JOHN, SONGUL KAYA, AND JULIA NOVO

(Communicated by )

*This paper is dedicated to William J. Layton on the occasion of his 60-th birthday.*

**Abstract.** A mantle convection model consisting of the stationary Stokes equations and a time-dependent convection-diffusion equation for the temperature is studied. The Stokes problem is discretized with a conforming inf-sup stable pair of finite element spaces and the temperature equation is stabilized with the SUPG method. Finite element error estimates are derived which show the dependency of the error of the solution of one problem on the error of the solution of the other equation. The dependency of the error bounds on the coefficients of the problem is monitored.

**Key words.** Mantel Convection, Stokes problem with variable viscosity, temperature problem with variable thermal convection, inf-sup stable finite elements, SUPG stabilization

### 1. Introduction

The process that occurs in the three-dimensional spherical shell between the crust and the core of the earth is called mantle convection. In this region, the magma moves very slowly. The movement is driven by the differences of the temperature at the hot core and the cool crust. Considering long time intervals, this movement is usually modeled with an incompressible viscous flow equation. Main features of this flow model are the high viscosity of order  $10^{20}$  Pa s [9], the small value of the thermal diffusivity (order  $\mathcal{O}(10^{-8} \text{ m}^2/\text{s})$  in [9]), and the dependency of the viscosity on other quantities, like the temperature. In turn, the temperature distribution is also driven by the movement of the magma, such that the modeling leads to a coupled problem. Simulations of mantle convection problems are quite challenging. One has to consider a three-dimensional problem in a very long time interval. With todays hardware and software capabilities, time intervals of almost  $10^9$  years are simulated [9], which results in performing many time steps. The resolution of important features, like plumes, requires to use adaptively refined grids. Massively parallel simulations with dynamic load balancing become necessary. The model (1) and (2) considered in this paper forms just the basic model. More advanced models use non-Newtonian fluids or they include a coupling to models for the behavior of the crust of the earth (solid material) to simulate the evolution of tectonic plates.

In this paper, the same model as in [27] is studied. Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , be bounded with polyhedral Lipschitz boundary  $\partial\Omega$ . Because of the large viscosity, the inertial term of the fundamental equations of fluid dynamics, the Navier–Stokes equations, can be neglected in mantle convection problems and thus, the equations reduce to the stationary incompressible Stokes equations. These equations with

variable kinematic viscosity  $\nu(\theta) > 0$  almost everywhere in  $\Omega$  are given by

$$(1) \quad \begin{aligned} -2\nabla \cdot (\nu(\theta) \mathbb{D}(\mathbf{u})) + \nabla p &= \mathbf{f} - \beta(\theta) \theta && \text{in } \Omega, \\ \nabla \cdot \mathbf{u} &= 0 && \text{in } \Omega, \\ \mathbf{u} &= \mathbf{0} && \text{on } \partial\Omega, \end{aligned}$$

where  $\mathbf{u}$  is the velocity field, the velocity deformation tensor  $\mathbb{D}(\mathbf{u}) = (\nabla \mathbf{u} + \nabla \mathbf{u}^T) / 2$  is the symmetric part of the gradient of  $\mathbf{u}$ ,  $p$  is the fluid pressure, and  $\mathbf{f}$  represents the body forces. Besides the dependency of the viscosity on the temperature  $\theta$ , a further impact of the temperature  $\theta$  is described by the function  $\beta$ .

The equation for the temperature is time-dependent. It is a convection-diffusion equation with a nonlinear diffusion term since the thermal diffusivity  $\kappa$  depends on the temperature

$$(2) \quad \begin{aligned} \partial_t \theta - \nabla \cdot (\kappa(\theta) \nabla \theta) + \mathbf{u} \cdot \nabla \theta &= g && \text{in } (0, T] \times \Omega, \\ \theta &= 0 && \text{in } (0, T] \times \partial\Omega, \\ \theta(0, \mathbf{x}) &= \theta_0(\mathbf{x}) && \text{in } \Omega. \end{aligned}$$

Altogether, (1) and (2) form a coupled system of equations. For the sake of easy implementation and efficiency, algorithms for the numerical solution of (1), (2) may decouple these problems and linearizations might be applied. Two algorithms in this spirit are as follows. Given a partition of the time interval into time steps  $0 = t_0 < t_1 < \dots < t_N = T$ :

**Algorithm 1.1.** *Nonlinear problem for the temperature.*

- (1) *the initial condition  $\theta_0$  is given*
- (2) *compute  $(\mathbf{u}_0, p_0)$  with  $\theta_0$*
- (3) *for  $i = 1, \dots, N$  do*
- (4) *compute  $\theta_i$  with  $\mathbf{u}_{i-1}$  or some other extrapolation, solving a nonlinear problem*
- (5) *compute  $(\mathbf{u}_i, p_i)$  with  $\theta_i$*
- (6) *end*

and

**Algorithm 1.2.** *Linear problem for the temperature.*

- (1) *the initial condition  $\theta_0$  is given*
- (2) *compute  $(\mathbf{u}_0, p_0)$  with  $\theta_0$*
- (3) *for  $i = 1, \dots, N$  do*
- (4) *compute  $\theta_i$  with  $\theta_{i-1}$  and  $\mathbf{u}_{i-1}$  or some other extrapolations solving a linear problem*
- (5) *compute  $(\mathbf{u}_i, p_i)$  with  $\theta_i$*
- (6) *end*

The finite element error analysis presented in this paper will focus on the individual equations which are solved in steps 4 and 5.

Finite element analysis of (1), (2) are already presented in [26, 27]. In [26], the case of constant viscosity ( $\nu = 1$ ) and thermal diffusivity is studied. In addition, the right-hand side of the Stokes equations depends linearly on the temperature in this paper. In both papers, the application of continuous piecewise linear ( $P_1$ ) finite elements for all quantities is considered. This approach requires the use of a stabilization of the discretization of the Stokes equations since the used pair of finite element spaces for velocity and pressure does not satisfy a discrete inf-sup condition. In [26, 27], the method of Brezzi and Pitkäranta [3] is applied. The convection-diffusion equation (2) is usually convection-dominated. Also this feature requires

the use of a stabilized method. The method of choice in [26, 27] was the SUPG method introduced in [4, 16]. Altogether, first order convergence in space and time was proved in [26, 27] for various norms of velocity, pressure, and temperature. Thermal convection problems with the stationary or evolutionary Navier–Stokes equations, instead of the Stokes equations, are studied in [22, 23, 28]. The papers [23, 28] consider inf-sup stable pairs of finite element spaces for the discretization of the Navier–Stokes equations and a Galerkin finite element discretization of both equations is analyzed. In [22], a divergence-conforming approximation of the velocity is studied. None of the papers mentioned above studies the dependency of the error bounds on the coefficients of the problem.

In the present paper, finite element pairs that satisfy a discrete inf-sup condition will be studied. Thus, higher than first order methods are included. Such methods are used in actual simulations [9]. Also for the temperature equation (2), higher order finite elements are considered. As in [26, 27], the SUPG stabilization is used. The dependency of the error bounds on the coefficients of the problem is tracked. As already mentioned, the finite element error analysis will focus on the individual problems (1) and (2) and it will study the impact of the error of the numerical solution of one problem on the error bound for the numerical solution of the other problem.

Standard notations for Lebesgue and Sobolev spaces and their norms will be used throughout the paper. In the analysis,  $C$  denotes a constant that is independent of the mesh width and the coefficients of (1), (2).

## 2. The Stokes problem with variable viscosity

Finite element methods for the Stokes equations with variable viscosity were analyzed in [18]. This section presents a slight generalization of the analysis which leads, however, to sharper error bounds provided the solution is sufficiently regular.

**2.1. The continuous problem.** Let the velocity space be denoted by  $V = (H_0^1(\Omega))^d$  and the pressure space by  $Q = L_0^2(\Omega)$ . A variational formulation of (1) reads as follows: Find  $(\mathbf{u}, p) \in (V, Q)$  satisfying

$$(3) \quad \begin{aligned} (2\nu\mathbb{D}(\mathbf{u}), \mathbb{D}(\mathbf{v})) - (\nabla \cdot \mathbf{v}, p) &= (\mathbf{f} - \beta(\theta)\theta, \mathbf{v}), \\ -(\nabla \cdot \mathbf{u}, q) &= 0 \end{aligned}$$

for all  $(\mathbf{v}, q) \in (V, Q)$ . It will be assumed that there is a positive constant  $\nu_{\min}$  such that

$$(4) \quad 0 < \nu_{\min} \leq \nu(\mathbf{x})$$

for almost all  $\mathbf{x} \in \Omega$ . With this assumption, it follows that  $\nu^{-1} \in L^\infty(\Omega)$ .

There holds the Poincaré inequality

$$(5) \quad \|\mathbf{v}\|_{L^2} \leq C\|\nabla\mathbf{v}\|_{L^2} \quad \forall \mathbf{v} \in V.$$

The space of weakly divergence-free functions is given by  $V_{\text{div}} = \{\mathbf{v} \in V : (\nabla \cdot \mathbf{v}, q) = 0, \forall q \in Q\}$ .

The following lemma proves that the weighted norm of the divergence is equivalent to the weighted norm of the deformation tensor.

**Lemma 2.1.** *Let  $\mathbf{v} \in H^1(\Omega)$ , then it holds*

$$(6) \quad \|\nu^{1/2}\nabla \cdot \mathbf{v}\|_{L^2} \leq \sqrt{d}\|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2}.$$

*Proof.* Let  $\mathbf{v} = (v_1, v_2, \dots, v_d)^T$ , then using the Cauchy-Schwarz inequality for sums gives

$$\begin{aligned}
\|\nu^{1/2}\nabla \cdot \mathbf{v}\|_{L^2}^2 &= \int_{\Omega} \nu \left( \sum_{i=1}^d \left( \frac{\partial v_i}{\partial x_i} \right) \right)^2 dx \leq \int_{\Omega} \nu \left( \sum_{i=1}^d 1 \right) \left( \sum_{i=1}^d \left( \frac{\partial v_i}{\partial x_i} \right)^2 \right) dx \\
&= d \int_{\Omega} \nu \left( \sum_{i=1}^d \left( \frac{\partial v_i}{\partial x_i} \right)^2 \right) dx \\
&\leq d \int_{\Omega} \nu \left( \sum_{i=1}^d \left( \frac{\partial v_i}{\partial x_i} \right)^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^d \frac{1}{4} \left( \frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i} \right)^2 \right) dx \\
&= \sqrt{d} \|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2}^2.
\end{aligned}$$

□

**Lemma 2.2** (Korn's inequality). *For all  $\mathbf{v} \in V$  it holds*

$$(7) \quad \frac{\nu_{\min}^{1/2}}{\sqrt{2}} \|\nabla \mathbf{v}\|_{L^2} \leq \|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2} \leq \|\nu^{1/2}\nabla \mathbf{v}\|_{L^2}.$$

*Proof.* The definition of the deformation tensor, triangle inequality, and the fact that  $\nu(\mathbf{x})$  is a scalar function yields

$$\|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2} \leq \frac{1}{2} \left( \|\nu^{1/2}\nabla \mathbf{v}\|_{L^2} + \|\nu^{1/2}\nabla \mathbf{v}^T\|_{L^2} \right) = \|\nu^{1/2}\nabla \mathbf{v}\|_{L^2},$$

which gives the right-hand side estimate of (7).

For the left-hand side estimate, Korn's inequality

$$(8) \quad \|\nabla \mathbf{v}\|_{L^2} \leq \sqrt{2} \|\mathbb{D}(\mathbf{v})\|_{L^2} \quad \forall \mathbf{v} \in V,$$

e.g., see [15], is used. With (8), one obtains

$$\frac{\nu_{\min}^{1/2}}{\sqrt{2}} \|\nabla \mathbf{v}\|_{L^2} \leq \nu_{\min}^{1/2} \|\mathbb{D}(\mathbf{v})\|_{L^2} \leq \|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2}.$$

□

**Remark 2.3.** The Korn inequality (7) is a slight improvement in comparison with the Korn inequality derived in [18] with respect to the right-hand side estimate. It is an open question whether the left-hand side estimate is improvable such that  $\|\nu^{1/2}\nabla \mathbf{v}\|_{L^2}$  appears. Pursuing the same approach that is used for constant viscosity, one finds on the one hand, using integration by parts and the symmetry of the deformation tensor

$$\begin{aligned}
-2(\nabla \cdot (\nu \mathbb{D}(\mathbf{v})), \mathbf{v}) &= 2(\nu \mathbb{D}(\mathbf{v}), \nabla \mathbf{v}) = (\nu \mathbb{D}(\mathbf{v}), \nabla \mathbf{v}) + \left( \nu (\mathbb{D}(\mathbf{v}))^T, \nabla \mathbf{v} \right) \\
&= (\nu \mathbb{D}(\mathbf{v}), \nabla \mathbf{v}) + (\nu \mathbb{D}(\mathbf{v}), \nabla \mathbf{v}^T) = 2\|\nu^{1/2}\mathbb{D}(\mathbf{v})\|_{L^2}^2
\end{aligned}$$

and on the other hand

$$-(\nabla \cdot (\nu \nabla \mathbf{v}), \mathbf{v}) = (\nu \nabla \mathbf{v}, \nabla \mathbf{v}) = \|\nu^{1/2}\nabla \mathbf{v}\|_{L^2}^2.$$

Since  $2(\nu \mathbb{D}(\mathbf{v}), \nabla \mathbf{v}) = (\nu \nabla \mathbf{v}, \nabla \mathbf{v}) + (\nu \nabla \mathbf{v}^T, \nabla \mathbf{v})$ , it is now sufficient to show that  $0 \leq (\nu \nabla \mathbf{v}^T, \nabla \mathbf{v})$ . For constant viscosity, one gets, using the identity  $\nabla \cdot (\nabla \mathbf{v}^T) = \nabla(\nabla \cdot \mathbf{v})$  and integration by parts, that  $(\nabla \mathbf{v}, \nabla \mathbf{v}) = (\nabla \cdot \mathbf{v}, \nabla \cdot \mathbf{v}) \geq 0$ . However, this approach cannot be applied for non-constant viscosity.

**2.2. Error analysis without temperature impact.** Consider a regular, non-degenerated family  $\{\mathcal{T}^h\}$  of triangulations of  $\bar{\Omega}$ . The mesh cells of a triangulation are denoted by  $K$ , their diameter by  $h_K$ , and the diameter of the largest ball inscribed in  $K$  by  $\rho_K$ . Let  $h = \max_{K \in \mathcal{T}^h} h_K$ . It is assumed that there exists a constant  $\sigma$ , independent of  $h$  and  $K$ , such that  $h_K/\rho_K \leq \sigma$  for all  $K \in \mathcal{T}^h$ .

Let  $V^h \subset V$  and  $Q^h \subset Q$  be conforming finite element spaces which fulfill the discrete inf-sup condition, i.e., there is a constant  $\beta$  independent of the mesh size parameter  $h$  such that

$$(9) \quad \inf_{q^h \in Q^h \setminus \{0\}} \sup_{\mathbf{v}^h \in V^h \setminus \{\mathbf{0}\}} \frac{(\nabla \cdot \mathbf{v}^h, q^h)}{\|\nabla \mathbf{v}^h\|_{L^2} \|q^h\|_{L^2}} \geq \beta > 0.$$

Since it will be assumed that the family of meshes is regular, the following inverse inequality holds

$$(10) \quad \|\mathbf{v}^h\|_{W^{m,p}(K)} \leq C_{\text{inv}} h_K^{n-m-d(\frac{1}{q}-\frac{1}{p})} \|\mathbf{v}^h\|_{W^{n,q}(K)},$$

for each  $\mathbf{v}^h \in V^h$ , where  $0 \leq n \leq m \leq 1$ ,  $1 \leq q \leq p \leq \infty$ , and  $h_K$  is the size (diameter) of the mesh cell  $K \in \mathcal{T}^h$ , see, e.g., [7, Thm. 3.2.6].

The finite element formulation of (3) with  $\beta = \mathbf{0}$  reads as follows: Find  $(\mathbf{u}^h, p^h) \in (V^h, Q^h)$  satisfying

$$(11) \quad \begin{aligned} 2(\nu \mathbb{D}(\mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) - (\nabla \cdot \mathbf{v}^h, p^h) &= (\mathbf{f}, \mathbf{v}^h), \\ -(\nabla \cdot \mathbf{u}^h, q^h) &= 0 \end{aligned}$$

for all  $(\mathbf{v}^h, q^h) \in (V^h, Q^h)$ .

Let  $V_{\text{div}}^h = \{\mathbf{v}^h \in V^h : (\nabla \cdot \mathbf{v}^h, q^h) = 0 \forall q^h \in Q^h\}$  be the space of discretely divergence-free functions. From the discrete inf-sup condition (9) it follows that this space is not empty. Then the finite element velocity from (11) is given by the problem: Find  $\mathbf{u}^h \in V_{\text{div}}^h$  such that

$$(12) \quad 2(\nu \mathbb{D}(\mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) = (\mathbf{f}, \mathbf{v}^h) \quad \forall \mathbf{v}^h \in V_{\text{div}}^h.$$

**Theorem 2.4.** *Let  $r, s \in [1, \infty]$  with  $r^{-1} + s^{-1} = 1$ ,  $\mathbf{u} \in (W^{1,2s}(\Omega))^d \cap V$ ,  $p \in L^{2s}(\Omega) \cap Q$ , and  $\nu \in L^r(\Omega)$  satisfying (4), then the following velocity error estimate is valid:*

$$(13) \quad \begin{aligned} &\|\nu^{1/2} \mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \\ &\leq 2\|\nu\|_{L^r}^{1/2} \inf_{\mathbf{v}^h \in V_{\text{div}}^h} \|\mathbb{D}(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}} + \frac{\sqrt{d}}{2} \|\nu^{-1}\|_{L^r}^{1/2} \inf_{q^h \in Q^h} \|p - q^h\|_{L^{2s}}. \end{aligned}$$

*Proof.* To obtain an error equation, consider the continuous formulation (3) for  $\mathbf{v}^h \in V_{\text{div}}^h$  and subtract the discrete equation (12) to get

$$(14) \quad 2(\nu \mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) - (\nabla \cdot \mathbf{v}^h, p) = 0 \quad \forall \mathbf{v}^h \in V_{\text{div}}^h.$$

Since  $(\nabla \cdot \mathbf{v}^h, q^h) = 0$  for all  $q^h \in Q^h$ , (14) can be written as

$$(15) \quad 2(\nu \mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) - (\nabla \cdot \mathbf{v}^h, p - q^h) = 0$$

for all  $(\mathbf{v}^h, q^h) \in (V_{\text{div}}^h, Q^h)$ . Then, the error is decomposed in two parts:  $\mathbf{u} - \mathbf{u}^h = \boldsymbol{\eta} - \boldsymbol{\phi}^h$ , where  $\boldsymbol{\eta} = \mathbf{u} - I^h(\mathbf{u})$  and  $\boldsymbol{\phi}^h = \mathbf{u}^h - I^h(\mathbf{u})$  and  $I^h : V \rightarrow V_{\text{div}}^h$  is some interpolant. Using this error decomposition in (15) and setting  $\mathbf{v}^h = \boldsymbol{\phi}^h$  gives

$$2\|\nu^{1/2} \mathbb{D}(\boldsymbol{\phi}^h)\|_{L^2}^2 = 2\left(\nu \mathbb{D}(\boldsymbol{\eta}), \mathbb{D}(\boldsymbol{\phi}^h)\right) - \left(\nabla \cdot \boldsymbol{\phi}^h, p - q^h\right).$$

Applying the Cauchy–Schwarz inequality on right-hand side of this estimate and using (6) yields

$$\begin{aligned} 2\|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2}^2 &\leq 2\|\nu^{1/2}\mathbb{D}(\boldsymbol{\eta})\|_{L^2}\|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2} \\ &\quad +\sqrt{d}\|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2}\|\nu^{-1/2}(p-q^h)\|_{L^2}. \end{aligned}$$

Dividing both sides by  $2\|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2}$  results in the estimate

$$(16) \quad \|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2} \leq \|\nu^{1/2}\mathbb{D}(\boldsymbol{\eta})\|_{L^2} + \frac{\sqrt{d}}{2}\|\nu^{-1/2}(p-q^h)\|_{L^2}.$$

Then, the application of the triangle inequality and (16) leads to

$$\begin{aligned} \|\nu^{1/2}\mathbb{D}(\mathbf{u}-\mathbf{u}^h)\|_{L^2} &\leq \|\nu^{1/2}\mathbb{D}(\boldsymbol{\eta})\|_{L^2} + \|\nu^{1/2}\mathbb{D}(\phi^h)\|_{L^2} \\ (17) \quad &\leq 2\|\nu^{1/2}\mathbb{D}(\boldsymbol{\eta})\|_{L^2} + \frac{\sqrt{d}}{2}\|\nu^{-1/2}(p-q^h)\|_{L^2}. \end{aligned}$$

Now, the application of the Hölder’s inequality to the first term on the right-hand side of (17) gives

$$(18) \quad \|\nu^{1/2}\mathbb{D}(\boldsymbol{\eta})\|_{L^2} = \|\nu\mathbb{D}(\boldsymbol{\eta}) : \mathbb{D}(\boldsymbol{\eta})\|_{L^1}^{1/2} \leq \|\nu\|_{L^r}^{1/2} \|\mathbb{D}(\boldsymbol{\eta}) : \mathbb{D}(\boldsymbol{\eta})\|_{L^s}^{1/2} = \|\nu\|_{L^r}^{1/2} \|\mathbb{D}(\boldsymbol{\eta})\|_{L^{2s}}.$$

In a similar way, the bound for the last term on the right-hand side of (17) is found to be

$$(19) \quad \|\nu^{-1/2}(p-q^h)\|_{L^2} \leq \|\nu^{-1}\|_{L^r}^{1/2} \|p-q^h\|_{L^{2s}}.$$

Inserting the bounds (18) and (19) in (17) and taking the infima gives (13).  $\square$

Of course, it would be possible to use different Lebesgue coefficients for the velocity and the pressure in the formulation of this theorem.

For many inf-sup stable pairs of finite element spaces there exists a linear interpolation operator  $I_{\text{div}}^h : V \rightarrow V^h$  with the properties

$$(20) \quad (\nabla \cdot (\mathbf{v} - I_{\text{div}}^h(\mathbf{v})), q^h) = 0 \quad \forall \mathbf{v} \in V, \forall q^h \in Q^h,$$

and

$$(21) \quad \|\nabla(\mathbf{v} - I_{\text{div}}^h(\mathbf{v}))\|_{L^s} \leq Ch^m \|\mathbf{v}\|_{W^{m+1,s}} \quad \forall \mathbf{v} \in W^{m+1,p}(\Omega), s \in [1, \infty],$$

see [13, Thm. 2.1]. The most notable case where the existence of such an operator could not be proved with the construction proposed in [13] is the Taylor–Hood pair of spaces  $P_2/P_1$  in three dimensions. Since for the solution of (3) it holds  $\mathbf{u} \in V_{\text{div}}$ , it follows from (20) that  $I_{\text{div}}^h(\mathbf{u}) \in V_{\text{div}}^h$ . Hence, the interpolation  $I_{\text{div}}^h(\mathbf{u})$  can be used to bound the best approximation error in (13) from above. Moreover, the Stokes projection defined in [12] can also be used to bound this error. This projection is discretely divergence-free and it has optimal approximations properties for any inf-sup stable pair of mixed finite elements.

In addition, to characterize the approximation properties for the space  $Q^h$ , let  $I^h(p)$  be the Lagrange interpolation of  $p$ . The following bound for the best approximation error can be found in [2, Thm. 4.4.4]

$$(22) \quad \|p - I^h(p)\|_{L^s} \leq Ch^{n+1} \|p\|_{W^{n+1,s}}, \quad \forall p \in W^{n+1,s}(\Omega), s \in [1, \infty],$$

with  $n+1 > d/s$  if  $1 < s \leq \infty$  and  $n+1 \geq d$  if  $s = 1$ .

**Corollary 2.5.** *Let  $r, s \in [1, \infty]$  with  $r^{-1} + s^{-1} = 1$ , let  $\mathbf{u} \in (W^{m+1,2s}(\Omega))^d \cap V$ ,  $p \in W^{n+1,2s}(\Omega) \cap Q$ , and  $\nu \in L^r(\Omega)$  with  $m, n \geq 0$ . Consider any pair of finite element spaces consisting of polynomials of degree  $m$  for the velocity and degree  $n$*

for the pressure that fits into the analysis presented in [13]. Then, the following velocity error estimate holds

$$(23) \quad \|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \leq C \left( h^m \|\nu\|_{L^r}^{1/2} \|\mathbf{u}\|_{W^{m+1,2s}} + h^{n+1} \|\nu^{-1}\|_{L^r}^{1/2} \|p\|_{W^{n+1,2s}} \right).$$

*Proof.* The velocity term on the right-hand side of (13) can be estimated by using  $\|\mathbb{D}(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}} \leq \|\nabla(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}}$ , setting  $\mathbf{v}^h = I_{\text{div}}^h(\mathbf{u})$ , and (21). For choosing  $q^h$  in the pressure term in the bound (13), one can use the Lagrangian interpolation and estimate (22).  $\square$

**Remark 2.6** (Comparison with an error bound from the literature). With different regularity assumptions, in particular on  $\nu(\mathbf{x})$ , it was proved in [18] that

$$(24) \quad \begin{aligned} & \|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \\ & \leq C \left[ h^m \nu_{\max}^{1/2} \left( 1 + \left( \frac{\nu_{\max}}{\nu_{\min}} \right)^{1/2} \right) \|\mathbf{u}\|_{H^{m+1}} + h^{n+1} \nu_{\min}^{-1/2} \|p\|_{H^{n+1}} \right], \end{aligned}$$

where  $\nu_{\min} = \min_{\mathbf{x} \in \bar{\Omega}} \nu(\mathbf{x})$ ,  $\nu_{\max} = \max_{\mathbf{x} \in \bar{\Omega}} \nu(\mathbf{x})$ .

Given the assumed regularity, the error bound (23) might be sharper than (24). This statement will be illustrated by a numerical example.

This example uses the same setup as the example presented in [18]. Let  $\Omega = (0, 1)^2$  and

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

then the prescribed velocity solution is defined by  $\mathbf{u} = (\partial_y \phi, -\partial_x \phi)^T$ . The pressure solution is given by

$$p(x, y) = \alpha_p \left( \pi^2 (xy^2 \cos(2\pi x^2 y) - x^2 y \sin(2\pi xy)) - \frac{1}{8} \right).$$

Simulations were performed with the two viscosity functions

$$\begin{aligned} \nu_1(x, y) &= 10^{-6} + (1 - 10^{-6}) \exp\left(-10^{13} \left((x - 0.5)^{10} + (y - 0.5)^{10}\right)\right), \\ \nu_2(x, y) &= 10^{-6} + (1 - 10^{-6}) \left(1 - \exp\left(-10^{13} \left((x - 0.5)^{10} + (y - 0.5)^{10}\right)\right)\right), \end{aligned}$$

see Figure 1. Whereas  $\nu_1$  is close to  $10^{-6}$  in the most part of the domain,  $\nu_2$  takes mostly values of around 1. But it is  $\nu_{1,\min} \approx \nu_{2,\min}$  and  $\nu_{1,\max} \approx \nu_{2,\max}$ . Consequently, the error bound (24) is almost the same for  $\nu_1$  and  $\nu_2$ . Consider the error bound (23) for  $r = 1$  and  $s = \infty$ . It is  $\|\nu_1\|_{L^1}^{1/2} = 0.09536616752$ ,  $\|\nu_1^{-1}\|_{L^1}^{1/2} = 991.6040647$ ,  $\|\nu_2\|_{L^1}^{1/2} = 0.9954427628$ ,  $\|\nu_2^{-1}\|_{L^1}^{1/2} = 25.86500587$  (the computation of the integrals was performed with Maple).

The Taylor–Hood pair of finite element spaces  $P_2/P_1$  was used on unstructured triangular grids, see Figure 2 for the coarsest grid, level 0. On the finest grid, level 8, there are 9 442 306 degrees of freedom for the velocity and 1 180 929 degrees of freedom for the pressure. The sparse direct solver umfpack [10] was used for solving the linear systems of equations. Since sometimes the bad condition number of the matrices resulted in notable round-off errors, a post-processing with an iterative solver was performed. This solver was the flexible GMRES(restart) method [25], with restart= 50, and with a coupled multigrid preconditioner, as described, e.g., in [17]. The iteration stopped if the Euclidean norm of the residual vector was less than  $10^{-12}$ . The simulations were performed with the code MoonMD [19].

Considering the left-hand side of (23) and (24), then there is a small scaling factor in the case of  $\nu_1$  and a larger scaling factor in the case of  $\nu_2$ . Two special

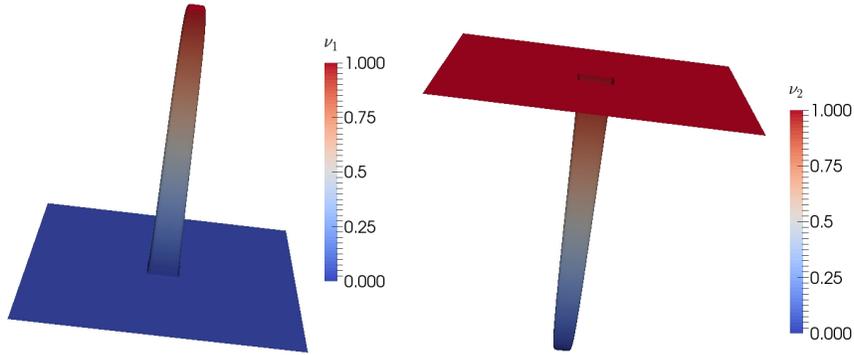
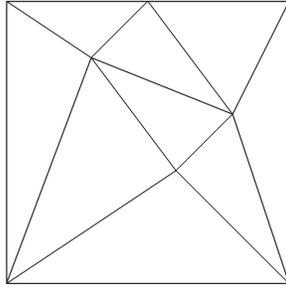
FIGURE 1. Viscosity functions  $\nu_1$  and  $\nu_2$ .

FIGURE 2. The coarsest grid, level 0.

situations will be studied, namely  $\alpha_v = 1, \alpha_p = 0$  and  $\alpha_v = 0, \alpha_p = 1$ . In the first situation, the pressure term drops from the right-hand side of (23) and this error bound is smaller for  $\nu_1$  than for  $\nu_2$ . The numerical results presented in Figure 3 show that the error itself is also smaller of one order of magnitude on finer grids. Note that  $\|\nu_1\|_{L^1}^{1/2}$  is smaller than  $\|\nu_2\|_{L^1}^{1/2}$  also by one order of magnitude. In the second situation, a so-called no-flow problem, the velocity term does not appear on the right-hand side of (23) and the error bound is larger for  $\nu_1$  than for  $\nu_2$ . Again, the numerical results behave in the same way. The ratio of the errors on the finest grid is around 400 whereas the ratio of  $\|\nu_1^{-1}\|_{L^1}^{1/2}$  and  $\|\nu_2^{-1}\|_{L^1}^{1/2}$  is only around 40. Thus, the error bound (23) still underpredicts the difference of the results for  $\nu_1$  and  $\nu_2$ , but it predicts qualitatively the correct behavior, in contrast to the error bound (24).

Finally, the general case with respect to  $\alpha_v$  and  $\alpha_p$  will be considered. The constants  $C$  in the error bounds (23) and (24) are essentially the result of applying estimates for the best approximation errors. If one assumes that these constants are of the same order for both error bounds and that the norms of the solution are of the same order, too, then they differ only in the factors with the viscosity function. In the considered example, both factors in (24) are of order  $\mathcal{O}(10^3)$  for both  $\nu_1$  and  $\nu_2$ . In contrast, the terms depending on the viscosity in (23) are of order  $\mathcal{O}(10^{-2})$  and  $\mathcal{O}(10^3)$  for  $\nu_1$ , i.e., only one factor is  $\mathcal{O}(10^3)$ , and of order  $\mathcal{O}(1)$  and  $\mathcal{O}(10)$  for  $\nu_2$ . Hence, under the assumptions from above, the error bound (23) can be expected to be smaller and therefore sharper.

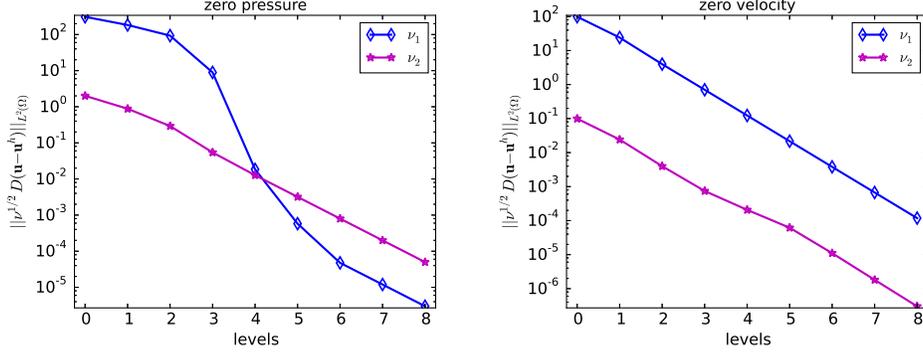


FIGURE 3. Errors  $\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}$  for the special situations  $\alpha_v = 1, \alpha_p = 0$  and  $\alpha_v = 0, \alpha_p = 1$ .

For completeness, the error analysis for the pressure will be presented briefly. It proceeds in the usual way, e.g., see [17].

Taking  $s = 1, r = \infty$  in (18) and using  $\|\mathbb{D}(\mathbf{v})\| \leq \|\nabla \mathbf{v}\|$  for all  $\mathbf{v} \in V$ , the inf-sup condition (9) can be written in the form

$$(25) \quad \|q^h\|_{L^2} \leq \sup_{\mathbf{v}^h \in V^h \setminus \{\mathbf{0}\}} \frac{(\nabla \cdot \mathbf{v}^h, q^h)}{\|\nabla \mathbf{v}^h\|_{L^2}} \leq \frac{\|\nu\|_{L^\infty}^{1/2}}{\beta} \sup_{\mathbf{v}^h \in V^h \setminus \{\mathbf{0}\}} \frac{(\nabla \cdot \mathbf{v}^h, q^h)}{\|\nu^{1/2}\mathbb{D}(\mathbf{v}^h)\|_{L^2}},$$

which is similar to the form used in [14].

**Theorem 2.7.** *Let all assumptions of Theorem 2.4 be satisfied and in addition assume that (9) holds, then*

$$(26) \quad \begin{aligned} \|p - p^h\|_{L^2} &\leq \left( C(s) + \frac{2\sqrt{d}}{\beta} \|\nu\|_{L^\infty}^{1/2} \|\nu^{-1}\|_{L^r}^{1/2} \right) \inf_{q^h \in Q^h} \|p - q^h\|_{L^{2s}} \\ &\quad + \frac{4}{\beta} \|\nu\|_{L^\infty}^{1/2} \|\nu\|_{L^r}^{1/2} \inf_{\mathbf{v}^h \in V^h} \|\mathbb{D}(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}}. \end{aligned}$$

*Proof.* Subtracting (11) from (3) and introducing an approximation  $\tilde{p}^h \in Q^h$  of the pressure yields the error equation

$$(27) \quad (\nabla \cdot \mathbf{v}^h, p^h - \tilde{p}^h) = (\nabla \cdot \mathbf{v}^h, p - \tilde{p}^h) - 2(\nu \mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\mathbf{v}^h))$$

for all  $\mathbf{v}^h \in V^h$ . The terms on the right-hand side of (27) can be bounded by using the Cauchy–Schwarz inequality and (6)

$$(28) \quad \begin{aligned} &|(\nabla \cdot \mathbf{v}^h, p^h - \tilde{p}^h)| \\ &\leq \left( \sqrt{d} \|\nu^{-1/2}(p - \tilde{p}^h)\|_{L^2} + 2\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \right) \|\nu^{1/2}\mathbb{D}(\mathbf{v}^h)\|_{L^2}. \end{aligned}$$

Dividing both sides by  $\|\nu^{1/2}\mathbb{D}(\mathbf{v}^h)\|_{L^2}$ , taking the supremum of (28) for  $\mathbf{v}^h \in V^h$ , and using (19) gives

$$(29) \quad \begin{aligned} \sup_{\mathbf{v}^h \in V^h \setminus \{\mathbf{0}\}} \frac{|(\nabla \cdot \mathbf{v}^h, p^h - \tilde{p}^h)|}{\|\nu^{1/2}\mathbb{D}(\mathbf{v}^h)\|_{L^2}} &\leq \sqrt{d} \|\nu^{-1/2}(p - \tilde{p}^h)\|_{L^2} + 2\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \\ &\leq \sqrt{d} \|\nu^{-1}\|_{L^r}^{1/2} \|p - \tilde{p}^h\|_{L^{2s}} + 2\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}. \end{aligned}$$

The statement of the theorem is now obtained by applying the triangle inequality

$$\|p - p^h\|_{L^2} \leq \|p - \tilde{p}^h\|_{L^2} + \|p^h - \tilde{p}^h\|_{L^2}$$

and using (25), (29), and the velocity error bound (13).  $\square$

**2.3. Error analysis with temperature impact.** In the next step, the impact of the temperature on the Stokes flow is taken into account. A finite element error analysis of a steady-state coupled Navier–Stokes and temperature system can be found in [23]. In this system, the temperature impact on the right-hand side is linear. Optimal order convergence in the  $H^1(\Omega)^d$  norm of the velocity and the  $L^2(\Omega)$  of the pressure was proved for inf-sup stable pairs of finite element spaces. The dependency of the constant in the error bound on the coefficients of the problem was not studied.

This section extends the analysis of the previous section to the situation that there is a special nonlinear temperature impact, as already given in (1), on the right-hand side of the Stokes equations. The dependency of the viscosity on the temperature will not be considered explicitly since it will be assumed that  $\nu(\theta)$  possesses the same regularity as in Section 2.2. Besides the errors  $\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}$  and  $\|p - p^h\|_{L^2}$ , also an estimate of the velocity error in  $L^2(\Omega)$  is provided. The latter error will appear in the error bound for the temperature.

Consider now the finite element problem: Find  $(\mathbf{u}^h, p^h) \in (V^h, Q^h)$  satisfying

$$(30) \quad \begin{aligned} 2(\nu\mathbb{D}(\mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) - (\nabla \cdot \mathbf{v}^h, p^h) &= (\mathbf{f} - \beta(\theta^h)\theta^h, \mathbf{v}^h), \\ -(\nabla \cdot \mathbf{u}^h, q^h) &= 0 \end{aligned}$$

for all  $(\mathbf{v}^h, q^h) \in (V^h, Q^h)$ , where  $\theta^h$  is some finite element approximation of the temperature field.

**Theorem 2.8.** *Let the assumptions of Theorem 2.4 be satisfied and assume that  $\theta, \theta^h \in H_0^1(\Omega)$ ,  $\beta \in L^6(\Omega)$  for all admissible temperature fields, and that  $\beta$  is Lipschitz continuous in this norm*

$$(31) \quad \|\beta(\theta_1) - \beta(\theta_2)\|_{L^6} \leq C\|\theta_1 - \theta_2\|_{L^6}$$

for all admissible temperature fields  $\theta_1, \theta_2$  and a constant that is independent of the temperature fields. Then, it holds

$$(32) \quad \begin{aligned} &\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \\ &\leq 2\|\nu\|_{L^r}^{1/2} \inf_{\mathbf{v}^h \in V_{\text{div}}^h} \|\mathbb{D}(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}} + \frac{\sqrt{d}}{2}\|\nu^{-1}\|_{L^r}^{1/2} \inf_{q^h \in Q^h} \|p - q^h\|_{L^{2s}} \\ &+ C\nu_{\min}^{-1/2} (\|\nabla\theta\|_{L^2} + \|\beta(\theta^h)\|_{L^6}) \|\nabla(\theta - \theta^h)\|_{L^2}. \end{aligned}$$

*Proof.* The proof starts in the same way as the proof of Theorem 2.4. Instead of (15), one arrives at the equation

$$2(\nu\mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\boldsymbol{\phi}^h)) - (\nabla \cdot \boldsymbol{\phi}^h, p - q^h) = (-\beta(\theta)\theta + \beta(\theta^h)\theta^h, \boldsymbol{\phi}^h).$$

The term on the right-hand side is split into

$$(33) \quad \left( (\beta(\theta^h) - \beta(\theta))\theta, \boldsymbol{\phi}^h \right) - \left( \beta(\theta^h)(\theta - \theta^h), \boldsymbol{\phi}^h \right).$$

The first term of (33) is estimated with Hölder's inequality, the Lipschitz continuity (31), the Sobolev embedding  $W^{1,1}(\Omega) \rightarrow L^{3/2}(\Omega)$ , the Sobolev embedding  $H^1(\Omega) \rightarrow L^6(\Omega)$ , Poincaré's inequality (5), the Cauchy–Schwarz inequality, and

Korn's inequality (7)

$$\begin{aligned}
\left( (\boldsymbol{\beta}(\theta^h) - \boldsymbol{\beta}(\theta)) \theta, \phi^h \right) &\leq \|\boldsymbol{\beta}(\theta^h) - \boldsymbol{\beta}(\theta)\|_{L^6} \|\theta\|_{L^6} \|\phi^h\|_{L^{3/2}} \\
&\leq C \|\theta - \theta^h\|_{L^6} \|\theta\|_{L^6} \|\phi^h\|_{W^{1,1}} \\
&\leq C \|\nabla(\theta - \theta^h)\|_{L^2} \|\nabla\theta\|_{L^2} \|\nabla\phi^h\|_{L^1} \\
&\leq C |\Omega|^{1/2} \|\nabla(\theta - \theta^h)\|_{L^2} \|\nabla\theta\|_{L^2} \|\nabla\phi^h\|_{L^2} \\
(34) \qquad \qquad \qquad &\leq C \nu_{\min}^{-1/2} \|\nabla(\theta - \theta^h)\|_{L^2} \|\nabla\theta\|_{L^2} \|\nu^{1/2} \mathbb{D}(\phi^h)\|_{L^2}.
\end{aligned}$$

The estimate of the second term of (33) is performed with the same tools, yielding

$$\begin{aligned}
\left( \boldsymbol{\beta}(\theta^h) (\theta - \theta^h), \phi^h \right) &\leq \|\boldsymbol{\beta}(\theta^h)\|_{L^6} \|\theta - \theta^h\|_{L^6} \|\phi^h\|_{L^{3/2}} \\
(35) \qquad \qquad \qquad &\leq C \nu_{\min}^{-1/2} \|\boldsymbol{\beta}(\theta^h)\|_{L^6} \|\nabla(\theta - \theta^h)\|_{L^2} \|\nu^{1/2} \mathbb{D}(\phi^h)\|_{L^2}.
\end{aligned}$$

Now, the proof continuous like the proof of Theorem 2.4, giving the statement of the theorem.  $\square$

**Theorem 2.9.** *Let the assumptions of Theorem 2.8 be satisfied. Then, it holds*

$$\begin{aligned}
\|p - p^h\|_{L^2} &\leq \left( C(s) + \frac{2\sqrt{d}}{\beta} \|\nu\|_{L^\infty}^{1/2} \|\nu^{-1}\|_{L^r}^{1/2} \right) \inf_{q^h \in Q^h} \|p - q^h\|_{L^{2s}} \\
&\quad + \frac{2}{\beta} \|\nu\|_{L^\infty}^{1/2} \|\nu\|_{L^r}^{1/2} \inf_{\mathbf{v}^h \in V^h} \|\mathbb{D}(\mathbf{u} - \mathbf{v}^h)\|_{L^{2s}} \\
(36) \qquad \qquad \qquad &\quad + \frac{C}{\beta} \|\nu\|_{L^\infty}^{1/2} \nu_{\min}^{-1/2} (\|\nabla\theta\|_{L^2} + \|\boldsymbol{\beta}(\theta^h)\|_{L^6}) \|\nabla(\theta - \theta^h)\|_{L^2}.
\end{aligned}$$

*Proof.* The proof proceeds along the lines of the proof of Theorem 2.7. By subtracting (3) and (30), one obtains for all  $\tilde{p}^h \in Q^h$  and all  $\mathbf{v}^h \in V^h$

$$\begin{aligned}
(\nabla \cdot \mathbf{v}^h, p^h - \tilde{p}^h) &= (\nabla \cdot \mathbf{v}^h, p - \tilde{p}^h) - 2(\nu \mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\mathbf{v}^h)) \\
&\quad - (\boldsymbol{\beta}(\theta) \theta - \boldsymbol{\beta}(\theta^h) \theta^h, \mathbf{v}^h).
\end{aligned}$$

The first two terms on the right-hand side are estimated in the same way as in the proof of Theorem 2.7 and the last term is bounded in the same way as in (33) – (35).  $\square$

Estimates (32) and (36) show that the order of convergence of the left-hand sides is bounded of the order of convergence of  $\|\nabla(\theta - \theta^h)\|_{L^2}$ , which is a usual term in the error bounds for scalar convection-diffusion equations (scaled with the square root of the diffusivity). The term  $\|\nabla\theta\|_{L^2}$  is usually bounded by a stability estimate and the term  $\|\boldsymbol{\beta}(\theta^h)\|_{L^6}$  is bounded by assumption (and can be even computed). The impact of the temperature error is scaled by  $\nu_{\min}^{-1/2}$ .

Finally, the error  $\|\mathbf{u} - \mathbf{u}^h\|_{L^2}$  will be studied. To this end, consider the dual Stokes problem Find  $(\phi_{\tilde{\mathbf{f}}}, \xi_{\tilde{\mathbf{f}}}) \in V \times Q$  such that for given  $\xi_{\tilde{\mathbf{f}}} \in L^2(\Omega)$

$$\begin{aligned}
(37) \qquad \qquad \qquad -2\nabla \cdot (\nu \mathbb{D}(\phi_{\tilde{\mathbf{f}}})) + \nabla \xi_{\tilde{\mathbf{f}}} &= \xi_{\tilde{\mathbf{f}}} \quad \text{in } \Omega, \\
&\quad \nabla \cdot \phi_{\tilde{\mathbf{f}}} = 0 \quad \text{in } \Omega
\end{aligned}$$

with homogeneous Dirichlet boundary conditions and its weak form

$$\begin{aligned}
(38) \qquad \qquad \qquad 2(\nu \mathbb{D}(\phi_{\tilde{\mathbf{f}}}), \mathbb{D}(\mathbf{v})) - (\nabla \cdot \mathbf{v}, \xi_{\tilde{\mathbf{f}}}) &= (\xi_{\tilde{\mathbf{f}}}, \mathbf{v}) \quad \forall \mathbf{v} \in V, \\
&\quad (\nabla \cdot \phi_{\tilde{\mathbf{f}}}, q) = 0 \quad \forall q \in Q.
\end{aligned}$$

It is assumed that the mapping

$$(\phi_{\tilde{\mathbf{f}}}, \xi_{\tilde{\mathbf{f}}}) \mapsto -2\nabla \cdot (\nu \mathbb{D}(\phi_{\tilde{\mathbf{f}}})) + \nabla \xi_{\tilde{\mathbf{f}}}$$

is an isomorphism from  $(H^2(\Omega) \cap V) \times (H^1(\Omega) \cap Q)$  to  $L^2(\Omega)$ .

**Theorem 2.10.** *Let the assumptions of Theorem 2.8 be satisfied and assume that  $\beta \in L^2(\Omega)$  and that  $\beta$  is Lipschitz continuous in with respect to the  $L^2(\Omega)$  norm*

$$(39) \quad \|\beta(\theta_1) - \beta(\theta_2)\|_{L^2} \leq C\|\theta_1 - \theta_2\|_{L^2}$$

for all admissible temperature fields  $\theta_1, \theta_2$  and a constant that is independent of the temperature fields. Then, it holds

$$\begin{aligned} & \|\mathbf{u} - \mathbf{u}^h\|_{L^2} \\ & \leq 2\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \sup_{\xi_{\hat{\mathbf{f}}} \in L^2(\Omega) \setminus \{0\}} \frac{1}{\|\xi_{\hat{\mathbf{f}}}\|_{L^2}} \left( \inf_{\phi^h \in V^h} \|\nu^{1/2}\nabla(\phi_{\hat{\mathbf{f}}} - \phi^h)\|_{L^2} \right) \\ & \quad + \sqrt{d}\|\nu^{-1}\|_{L^\infty}^{1/2}\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2} \sup_{\xi_{\hat{\mathbf{f}}} \in L^2(\Omega) \setminus \{0\}} \frac{1}{\|\xi_{\hat{\mathbf{f}}}\|_{L^2}} \left( \inf_{r^h \in Q^h} \|\xi_{\hat{\mathbf{f}}} - r^h\|_{L^2} \right) \\ & \quad + \sqrt{d} \inf_{q^h \in Q^h} \|p - q^h\|_{L^2} \sup_{\xi_{\hat{\mathbf{f}}} \in L^2(\Omega) \setminus \{0\}} \frac{1}{\|\xi_{\hat{\mathbf{f}}}\|_{L^2}} \left( \inf_{\phi^h \in V^h} \|\nabla(\phi_{\hat{\mathbf{f}}} - \phi^h)\|_{L^2} \right) \\ (40) \quad & + C(\|\theta\|_{L^2} + \|\beta(\theta^h)\|_{L^2})\|\theta - \theta^h\|_{L^2}. \end{aligned}$$

*Proof.* Choosing  $\mathbf{v} = \mathbf{u} - \mathbf{u}^h$  in (38) gives

$$(41) \quad (\xi_{\hat{\mathbf{f}}}, \mathbf{u} - \mathbf{u}^h) = 2\left(\nu\mathbb{D}(\phi_{\hat{\mathbf{f}}}), \mathbb{D}(\mathbf{u} - \mathbf{u}^h)\right) - \left(\nabla \cdot (\mathbf{u} - \mathbf{u}^h), \xi_{\hat{\mathbf{f}}}\right).$$

Using the weak form of the Stokes problem and the finite element problem (30), one gets for  $\phi^h \in V_{\text{div}}^h$  and  $q^h \in Q^h$  arbitrary

$$2\left(\nu\mathbb{D}(\mathbf{u} - \mathbf{u}^h), \mathbb{D}(\phi^h)\right) = (\nabla \cdot \phi^h, p - q^h) + (-\beta(\theta)\theta + \beta(\theta^h)\theta^h, \phi^h).$$

Inserting this identity in (41) and expanding it with some terms which are zero leads to

$$\begin{aligned} (\xi_{\hat{\mathbf{f}}}, \mathbf{u} - \mathbf{u}^h) & = 2(\nu\mathbb{D}(\phi_{\hat{\mathbf{f}}} - \phi^h), \mathbb{D}(\mathbf{u} - \mathbf{u}^h)) - (\nabla \cdot (\mathbf{u} - \mathbf{u}^h), \xi_{\hat{\mathbf{f}}} - r^h) \\ & \quad + (\nabla \cdot (\phi^h - \phi_{\hat{\mathbf{f}}}), p - q^h) + (-\beta(\theta)\theta + \beta(\theta^h)\theta^h, \phi^h) \end{aligned}$$

for arbitrary  $r^h \in Q^h$ . Then, it is straightforward to obtain (40) using the definition of the  $L^2(\Omega)$  norm, the decomposition (33), and the bound

$$\begin{aligned} (-\beta(\theta)\theta + \beta(\theta^h)\theta^h, \phi^h) & \leq \|\beta(\theta^h) - \beta(\theta)\|_{L^2}\|\theta\|_{L^2}\|\phi^h\|_{L^\infty} \\ & \quad + \|\beta(\theta^h)\|_{L^2}\|\theta - \theta^h\|_{L^2}\|\phi^h\|_{L^\infty} \\ & \leq (C\|\theta\|_{L^2} + \|\beta(\theta^h)\|_{L^2})\|\theta - \theta^h\|_{L^2}\|\phi^h\|_{L^\infty}. \end{aligned}$$

To conclude, one has to estimate  $\|\phi^h\|_{L^\infty}$  in terms of  $\|\phi_{\hat{\mathbf{f}}}\|_2$ , which in turn is bounded in terms of  $\|\hat{\mathbf{f}}\|$ . Denote by  $\mathbf{w}^h$  any approximation of  $\phi_{\hat{\mathbf{f}}}$  which is stable in the  $L^\infty(\Omega)$  norm. Then, using the inverse inequality (10), the isomorphism property, and a Sobolev embedding yields

$$\begin{aligned} \|\phi^h\|_{L^\infty} & \leq \|\phi^h - \mathbf{w}^h\|_{L^\infty} + \|\mathbf{w}^h\|_{L^\infty} \leq C_{\text{inv}}h^{-d/2}\|\phi^h - \mathbf{w}^h\|_{L^2} + \|\mathbf{w}^h\|_{L^\infty} \\ & \leq Ch^{-d/2}h^2\|\phi_{\hat{\mathbf{f}}}\|_{H^2} + C\|\phi_{\hat{\mathbf{f}}}\|_{L^\infty} \leq C\|\phi_{\hat{\mathbf{f}}}\|_{H^2} \leq C\|\xi_{\hat{\mathbf{f}}}\|_{L^2}. \end{aligned}$$

□

**Remark 2.11.** From estimate (32), it follows that the error bound for  $\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}$  contains

$$\frac{C(\theta, \theta^h)}{\nu_{\min}^{1/2}}\|\nabla(\theta - \theta^h)\|_{L^2}$$

and from (36), one can see that the error bound for  $\|p - p^h\|_{L^2}$  contains

$$\frac{C(\theta, \theta^h)}{\nu_{\min}^{1/2}} \|\nu\|_{L^\infty}^{1/2} \|\nabla(\theta - \theta^h)\|_{L^2}.$$

In applications, e.g., [9],  $\|\nu\|_{L^\infty}$  is larger by several orders of magnitude than  $\nu_{\min}$ . Hence, the error of approximating the temperature will have a notably larger impact on the pressure error  $\|p - p^h\|_{L^2}$  than on the scaled velocity error  $\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}$ .

The situation for  $\|\mathbf{u} - \mathbf{u}^h\|_{L^2}$  is more complicated since the dual problem is defined with  $\nu$  and thus  $\xi_{\mathcal{F}}$  will depend on the viscosity. Therefore it is hard to describe the dependency of this error on the error of the temperature since not only the last term in the error bound (40) depends on  $\|\theta - \theta^h\|_{L^2}$  but also the first two terms depend, via  $\|\nu^{1/2}\mathbb{D}(\mathbf{u} - \mathbf{u}^h)\|_{L^2}$ , on  $\|\nabla(\theta - \theta^h)\|_{L^2}$ . Altogether, we could not obtain a clear description of the impact of the temperature error on  $\|\mathbf{u} - \mathbf{u}^h\|_{L^2}$ .

### 3. The temperature equation with nonlinear diffusivity

The equation for the temperature is a time-dependent convection-diffusion equation with a nonlinear diffusion term. In practice, this is convection-dominated. It is well known that in this situation, the standard finite element method produces spurious oscillations and the use of a stabilized method is necessary. There are several methods that can be applied, see [24], see also [5, 8, 21, 1]. In this paper, as in [27], the most popular method, the streamline-upwind Petrov–Galerkin (SUPG) method is studied. An error analysis of this method for transient convection-reaction-diffusion equations can be found in [20]. As temporal discretization, the backward Euler scheme is considered.

Compared with [20], there are a few different aspects. First, the reaction field is missing such that the usual assumption, that reaction minus one half of the divergence of the convection is bounded from below by a positive constant, does not hold (see the comments at the end of the section). Second, it cannot be assumed that the convection field is divergence-free since it is the finite element solution of a Stokes problem. And third, for technical reasons (application of an inverse inequality) the situation has to be considered that the thermal diffusivity  $\kappa$  is approximated by a function  $\kappa^h$  that belongs to a finite-dimensional space. These differences give rise to several new technical issues.

**3.1. The continuous problem and its discretization.** Let  $V_\theta = H_0^1(\Omega)$ . A variational form of the equation for the temperature reads as follows: find  $\theta : (0, T] \rightarrow V_\theta$  such that

$$(42) \quad (\partial_t \theta, \phi) + (\kappa(\theta) \nabla \theta, \nabla \phi) + (\mathbf{u} \cdot \nabla \theta, \phi) = (g, \phi) \quad \forall \phi \in V_\theta.$$

A Poincaré estimate of form (5) holds for functions from  $V_\theta$ .

Let the time interval be decomposed in equidistant time steps  $0 = t_0 < t_1 < \dots < t_N = T$  and denote the length of the time step by  $\Delta t = t_n - t_{n-1}$ ,  $n = 1, \dots, N$ . Let  $V_\theta^h \subset V_\theta$  be a conforming finite element space defined on a triangulation  $\mathcal{T}^h$  of  $\Omega$ . Because the triangulations are assumed to be regular, an inverse estimate of type (10) holds for functions from  $V_\theta^h$ . Let  $\kappa^h(\theta)$ , defined on  $V_\theta$ , be a piecewise polynomial approximation of  $\kappa(\theta)$  on  $\mathcal{T}^h$  and let  $\mathbf{u}^h = \mathbf{u}_{n-1}^h \in V^h$  be given. Then, the backward Euler/SUPG method reads as follows: For  $n = 1, 2, \dots, N$  find

$\theta_n^h \in V_\theta^h$  such that

$$\begin{aligned}
& \left( \frac{\theta_n^h - \theta_{n-1}^h}{\Delta t}, \phi^h \right) + \left( \kappa^h(\hat{\theta}^h) \nabla \theta_n^h, \nabla \phi^h \right) + \frac{1}{2} [(\mathbf{u}^h \cdot \nabla \theta_n^h, \phi^h) - (\mathbf{u}^h \cdot \nabla \phi^h, \theta_n^h)] \\
& = (g_n, \phi^h) \\
(43) \quad & + \sum_{K \in \mathcal{T}^h} \delta_K \left( g_n - \frac{\theta_n^h - \theta_{n-1}^h}{\Delta t} + \nabla \cdot \left( \kappa^h(\hat{\theta}^h) \nabla \theta_n^h \right) - \mathbf{u}^h \cdot \nabla \theta_n^h, \mathbf{u}^h \cdot \nabla \phi^h \right)_K
\end{aligned}$$

for all  $\phi^h \in V_\theta^h$ , where either  $\hat{\theta}^h = \theta_{n-1}^h$  or  $\hat{\theta}^h = \theta_n^h$  and  $\{\delta_K\}$  is the local stabilization parameters. The skew-symmetric form of the convective term is used since  $\mathbf{u}^h$  is in general not weakly divergence-free. For the approximation  $\kappa^h$  of  $\kappa$ , it will be assumed that there are constants independent of  $h$  with

$$(44) \quad 0 < \kappa_{\min} \leq \kappa(t, \mathbf{x}, \phi^h), \kappa^h(t, \mathbf{x}, \phi^h) \leq \kappa_{\max} < \infty$$

for all  $\phi^h \in V_\theta^h$ ,  $t \in [0, T]$ ,  $\mathbf{x} \in \Omega$ .

Method (43) is written in short form

$$\begin{aligned}
& (\theta_n^h - \theta_{n-1}^h, \phi^h) + \Delta t a_{\text{SUPG}}(\hat{\theta}^h; \theta_n^h, \phi^h) = \Delta t (g_n, \phi^h) \\
(45) \quad & + \Delta t \sum_{K \in \mathcal{T}^h} \delta_K (g_n, \mathbf{u}^h \cdot \nabla \phi^h)_K - \Delta t \sum_{K \in \mathcal{T}^h} \delta_K (\theta_n^h - \theta_{n-1}^h, \mathbf{u}^h \cdot \nabla \phi^h)_K,
\end{aligned}$$

where

$$\begin{aligned}
a_{\text{SUPG}}(\hat{\theta}^h; \theta_n^h, \phi^h) & = \left( \kappa^h(\hat{\theta}^h) \nabla \theta_n^h, \nabla \phi^h \right) + \frac{1}{2} [(\mathbf{u}^h \cdot \nabla \theta_n^h, \phi^h) - (\mathbf{u}^h \cdot \nabla \phi^h, \theta_n^h)] \\
& \quad - \sum_{K \in \mathcal{T}^h} \delta_K \left( \nabla \cdot \left( \kappa^h(\hat{\theta}^h) \nabla \theta_n^h \right) - \mathbf{u}^h \cdot \nabla \theta_n^h, \mathbf{u}^h \cdot \nabla \phi^h \right)_K.
\end{aligned}$$

### 3.2. Error analysis.

**Lemma 3.1** (Coercivity of the SUPG form). *Let*

$$(46) \quad \delta_K \leq \frac{h_K^2}{C_{\text{inv}}^2 \kappa_{\max}}$$

and assume (44), then it holds

$$\begin{aligned}
a_{\text{SUPG}}(\psi^h; \phi^h, \phi^h) & \geq \frac{1}{2} \left( \|\kappa^h(\psi^h)^{1/2} \nabla \phi^h\|_{L^2}^2 + \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right) \\
(47) \quad & =: \frac{1}{2} \|\phi^h\|_{\text{SUPG}(\psi^h)}^2
\end{aligned}$$

for all  $\phi^h, \psi^h \in V_\theta^h$ .

*Proof.* Choosing the second and the third argument of the SUPG form to be the same, one finds that the convective term vanishes, due to its skew-symmetry, and one gets

$$\begin{aligned}
a_{\text{SUPG}}(\psi^h; \phi^h, \phi^h) & = \|\kappa^h(\psi^h)^{1/2} \nabla \phi^h\|_{L^2}^2 + \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \\
(48) \quad & - \sum_{K \in \mathcal{T}^h} \delta_K (\nabla \cdot (\kappa^h(\psi^h) \nabla \phi^h), \mathbf{u}^h \cdot \nabla \phi^h)_K.
\end{aligned}$$

The usual technique for estimating the last term from above includes the application of the inverse estimate (10). However, the use of this estimates requires that the

term  $\nabla \cdot (\kappa^h (\psi^h) \nabla \phi^h)$  belongs to a finite-dimensional space. For this reason, the piecewise polynomial approximation  $\kappa^h$  was introduced. One obtains, using the Cauchy–Schwarz inequality, (44), the inverse estimate (10), Young’s inequality, and the condition (46) on the stabilization parameter

$$\begin{aligned} & \sum_{K \in \mathcal{T}^h} \delta_K (\nabla \cdot (\kappa^h (\psi^h) \nabla \phi^h), \mathbf{u}^h \cdot \nabla \phi^h)_K \\ & \leq \sum_{K \in \mathcal{T}^h} \delta_K \kappa_{\max}^{1/2} \|\nabla \cdot (\kappa^h (\psi^h)^{1/2} \nabla \phi^h)\|_{L^2(K)} \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)} \\ & \leq \frac{1}{2} \sum_{K \in \mathcal{T}^h} \left[ \frac{C_{\text{inv}}^2 \kappa_{\max}}{h_K^2} \delta_K \|\kappa^h (\psi^h)^{1/2} \nabla \phi^h\|_{L^2(K)}^2 + \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right] \\ & \leq \frac{1}{2} \left( \|\kappa^h (\psi^h)^{1/2} \nabla \phi^h\|_{L^2}^2 + \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right). \end{aligned}$$

Inserting this upper bound in the last term of (48) gives the statement of the lemma.  $\square$

**Lemma 3.2** (Stability of the solution). *Let the assumptions of Lemma 3.1 be satisfied and let in addition*

$$(49) \quad \delta_K \leq \frac{\Delta t}{4}.$$

Then, it is for all discrete times  $t_n = n\Delta t$ ,  $n = 0, 1, \dots, N$

$$\|\theta_n^h\|_{L^2}^2 + \frac{\Delta t}{2} \sum_{j=1}^n \|\theta_j^h\|_{\text{SUPG}(\hat{\theta}^h)}^2 \leq \|\theta_0^h\|_{L^2}^2 + \Delta t \left( \frac{2}{\kappa_{\min}} + \Delta t \right) \sum_{j=1}^n \|g_j\|_{L^2}^2,$$

where in the sum it is either  $\hat{\theta}^h = \theta_j^h$  or  $\hat{\theta}^h = \theta_{j-1}^h$ .

*Proof.* The proof follows the lines of the proof of Theorem 3.1 in [20], using in particular the relation

$$(50) \quad (\theta_n^h - \theta_{n-1}^h, \theta_n^h) = \frac{1}{2} (\|\theta_n^h\|_{L^2}^2 + \|\theta_n^h - \theta_{n-1}^h\|_{L^2}^2 - \|\theta_{n-1}^h\|_{L^2}^2).$$

The only difference is the estimate of the first term on the right-hand side of (45), which has to take into account the absence of a reactive term. Thus, this term is bounded by using Poincaré’s inequality (5), the Cauchy–Schwarz inequality, and Young’s inequality in the following way

$$(51) \quad \Delta t (g_n, \phi^h) \leq \frac{\Delta t}{4} \|\kappa^h (\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2}^2 + \Delta t \|\kappa^h (\hat{\theta}^h)^{-1/2} g_n\|_{L^2}^2.$$

The first term is then absorbed in the SUPG norm.  $\square$

It follows from (49) that the stabilization parameter depends on the length of the time step. This issue is discussed comprehensively in [20]. In this paper, error estimates for stabilization parameters independently of the time step could be derived under the assumption that the convection field is stationary. This assumption cannot be made in the context of the application in mind.

In the sequel,  $\pi^h \theta \in V_\theta^h$  will denote the following elliptic projection

$$(52) \quad (\kappa(\theta(t)) \nabla (\pi^h \theta(t) - \theta(t)), \nabla \phi^h) = 0 \quad \forall \phi^h \in V_\theta^h.$$

Following [29, Lemma 13.1], there are optimal error bounds for both  $\pi^h\theta$  and  $\partial_t(\pi^h\theta) = \pi^h(\partial_t\theta)$ . Moreover, the following bound holds for all  $t$  [29, Lemma 13.3]

$$(53) \quad \|\nabla\pi_h(\theta(t))\|_{L^\infty} \leq C(\theta).$$

**Theorem 3.3** (Finite element error estimate.). *Let the assumptions of Lemmas 3.1 and 3.2 be satisfied and let the regularities appearing on the right-hand side of the following estimate be assumed, then the temperature error satisfies*

$$(54) \quad \begin{aligned} & \|\theta_n - \theta_n^h\|_{L^2}^2 + \Delta t \sum_{j=1}^n \|\theta_j - \theta_j^h\|_{\text{SUPG}(\hat{\theta}_j^h)}^2 \\ & \leq 2\|\theta_n - \pi^h\theta_n\|_{L^2}^2 + 2\Delta t \sum_{j=1}^n \|\theta_j - \pi^h\theta_j\|_{\text{SUPG}(\hat{\theta}_j^h)}^2 + \|\theta_0^h - \pi^h\theta_0\|_{L^2}^2 \\ & \quad + C\Delta t \sum_{j=1}^n \left[ \left( \frac{C(\theta_j)^2}{\kappa_{\min}} + \frac{h}{\kappa_{\max}} \right) \|\kappa(\theta_j) - \kappa^h(\hat{\theta}_j^h)\|_{L^2}^2 \right. \\ & \quad + \left( C(\theta_j)^2 + \frac{\|\nabla \cdot \mathbf{u}^h\|_{L^{2s}}^2}{\kappa_{\min}} + \frac{1}{\delta_{\min}} + \frac{\|\mathbf{u}^h\|_{L^\infty}^2}{\kappa_{\max}} \right) \|\theta_j - \pi^h\theta_j\|_{L^2}^2 \\ & \quad + \left. \left( \frac{\|\nabla\theta_j\|_{L^{2s}}^2}{\kappa_{\min}} + \frac{\|\theta_j\|_{L^\infty}^2}{\kappa_{\min}} + \max_{K \in \mathcal{T}^h} \{\delta_K\} \|\nabla\theta_j\|_{L^\infty}^2 \right) \|\mathbf{u}_j - \mathbf{u}^h\|_{L^2}^2 \right. \\ & \quad \left. + \|(I - \pi^h)\partial_t\theta_j\|_{L^2} + \Delta t \int_{t_{j-1}}^{t_j} \|\partial_{tt}\theta\|_{H^1}^2 d\tau \right], \end{aligned}$$

where  $s > 1$  if  $d = 2$  and  $s \geq 3/2$  if  $d = 3$ .

*Proof.* Let the error at time  $t_n$  be decomposed as follows  $\theta_n - \theta_n^h = (\theta_n - \pi^h\theta_n) - (\theta_n^h - \pi^h\theta_n) = \eta_n - \epsilon_n^h$ . An error equation is obtained by subtracting (42) from (43). A straightforward calculation, noting that a diffusive term vanishes because the elliptic projection (52) is used, yields

$$(55) \quad \begin{aligned} & (\epsilon_n^h - \epsilon_{n-1}^h, \phi^h) + \Delta t a_{\text{SUPG}}(\hat{\theta}^h; \epsilon_n^h, \phi^h) + \Delta t \left( (\kappa^h(\hat{\theta}^h) - \kappa(\theta_n)) \nabla\pi^h\theta_n, \nabla\phi^h \right) \\ & \quad + \Delta t (\mathbf{u}_n \cdot \nabla\theta_n, \phi^h) - \frac{\Delta t}{2} \left[ (\mathbf{u}^h \cdot \nabla\pi^h\theta_n, \phi^h) - (\mathbf{u}^h \cdot \nabla\phi^h, \pi^h\theta_n) \right] \\ & = \Delta t (T_n^h, \phi^h) - \sum_{K \in \mathcal{T}^h} \delta_K (\epsilon_n^h - \epsilon_{n-1}^h, \mathbf{u}^h \cdot \nabla\phi^h)_K + \sum_{K \in \mathcal{T}^h} \delta_K (T_n^h, \mathbf{u}^h \cdot \nabla\phi^h)_K \\ & \quad + \Delta t \sum_{K \in \mathcal{T}^h} \delta_K \left( \nabla \cdot (\kappa^h(\hat{\theta}^h) \nabla\pi^h\theta_n - \kappa(\theta_n) \nabla\theta_n), \mathbf{u}^h \cdot \nabla\phi^h \right)_K \\ & \quad + \Delta t \sum_{K \in \mathcal{T}^h} \delta_K (\mathbf{u}_n \cdot \nabla\theta_n - \mathbf{u}^h \cdot \nabla\pi^h\theta_n, \mathbf{u}^h \cdot \nabla\phi^h)_K \end{aligned}$$

with the truncation error

$$T_n^h = (\partial_t\theta_n - \pi^h(\partial_t\theta_n)) + \left( \pi^h(\partial_t\theta_n) - \frac{\pi^h\theta_n - \pi^h\theta_{n-1}}{\Delta t} \right).$$

Now, the arguments proceed along the lines of the proof of [20, Theorem 4.1]. Using Hölder's inequality and (53) gives

$$\begin{aligned}
& \left( (\kappa(\theta_n) - \kappa^h(\hat{\theta}^h)) \nabla \pi^h \theta_n, \nabla \phi^h \right) \\
& \leq \|\nabla \pi^h \theta_n\|_{L^\infty} \|\kappa(\theta_n) - \kappa^h(\hat{\theta}^h)\|_{L^2} \kappa_{\min}^{-1/2} \|\kappa^h(\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2} \\
& \leq \frac{C(\theta_n)}{\kappa_{\min}^{1/2}} \|\kappa(\theta_n) - \kappa^h(\hat{\theta}^h)\|_{L^2} \|\kappa^h(\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2}.
\end{aligned}$$

Next, a bound for the convective term will be derived. Using integration by parts and that  $\mathbf{u}_n$  is weakly divergence-free, this term can be split into

$$\begin{aligned}
& (\mathbf{u}_n \cdot \nabla \theta_n, \phi^h) - \frac{1}{2} \left[ (\mathbf{u}^h \cdot \nabla \pi^h \theta_n, \phi^h) - (\mathbf{u}^h \cdot \nabla \phi^h, \pi^h \theta_n) \right] \\
& = \frac{1}{2} \left( (\mathbf{u}_n \cdot \nabla \theta_n, \phi^h) - (\mathbf{u}^h \cdot \nabla \pi^h \theta_n, \phi^h) \right) \\
(56) \quad & - \frac{1}{2} \left( (\mathbf{u}_n \cdot \nabla \phi^h, \theta_n) - (\mathbf{u}^h \cdot \nabla \phi^h, \pi^h \theta_n) \right).
\end{aligned}$$

An estimate of the first term on the right-hand side of (56) is derived by using integration by parts, Hölder's inequality, the Sobolev imbedding  $H^1(\Omega) \rightarrow L^{2s/(s-1)}(\Omega)$  which holds for  $s > 1$  if  $d = 2$  and  $s \geq 3/2$  if  $d = 3$ , and the Cauchy-Schwarz inequality for sums

$$\begin{aligned}
& \left( (\mathbf{u}_n \cdot \nabla \theta_n, \phi^h) - (\mathbf{u}^h \cdot \nabla \pi^h \theta_n, \phi^h) \right) \\
& = \left( (\mathbf{u}_n - \mathbf{u}^h) \cdot \nabla \theta_n, \phi^h \right) + (\mathbf{u}^h \cdot \nabla (\theta_n - \pi^h \theta_n), \phi^h) \\
& = \left( (\mathbf{u}_n - \mathbf{u}^h) \cdot \nabla \theta_n, \phi^h \right) - (\nabla \cdot \mathbf{u}^h, (\theta_n - \pi^h \theta_n) \phi^h) - (\theta_n - \pi^h \theta_n, \mathbf{u}^h \cdot \nabla \phi^h) \\
& \leq \|\mathbf{u}_n - \mathbf{u}^h\|_{L^2} \|\nabla \theta_n\|_{L^{2s}} \|\phi^h\|_{L^{2s/(s-1)}} + \|\nabla \cdot \mathbf{u}^h\|_{L^{2s}} \|\theta_n - \pi^h \theta_n\|_{L^2} \|\phi^h\|_{L^{2s/(s-1)}} \\
& \quad - \sum_{K \in \mathcal{T}^h} \delta_K \left( \frac{\theta_n - \pi^h \theta_n}{\delta_K}, \mathbf{u}^h \cdot \nabla \phi^h \right) \\
& \leq \frac{1}{\kappa_{\min}^{1/2}} \left( \|\nabla \theta_n\|_{L^{2s}} \|\mathbf{u}_n - \mathbf{u}^h\|_{L^2} + \|\nabla \cdot \mathbf{u}^h\|_{L^{2s}} \|\theta_n - \pi^h \theta_n\|_{L^2} \right) \\
& \quad \times \|\kappa^h(\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2} + \frac{1}{\delta_{\min}^{1/2}} \|\theta_n - \pi^h \theta_n\|_{L^2} \left( \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right)^{1/2}.
\end{aligned}$$

A bound of the second term of (56) is obtained in a similar way

$$\begin{aligned}
& (\mathbf{u}_n \cdot \nabla \phi^h, \theta_n) - (\mathbf{u}^h \cdot \nabla \phi^h, \pi^h \theta_n) \\
& = \left( (\mathbf{u}_n - \mathbf{u}^h) \cdot \nabla \phi^h, \theta_n \right) + (\mathbf{u}^h \cdot \nabla \phi^h, \theta_n - \pi^h \theta_n) \\
& \leq \frac{\|\theta_n\|_{L^\infty}}{\kappa_{\min}^{1/2}} \|\mathbf{u}_n - \mathbf{u}^h\|_{L^2} \|\kappa^h(\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2} \\
& \quad + \frac{1}{\delta_{\min}^{1/2}} \|\theta_n - \pi^h \theta_n\|_{L^2} \left( \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right)^{1/2}.
\end{aligned}$$

The estimate of the diffusion term in (55), which comes from the stabilization, starts with Hölder's inequality and the inverse estimate (10)

$$(57) \quad \sum_{K \in \mathcal{T}^h} \delta_K \left( \nabla \cdot \left( \kappa^h(\hat{\theta}^h) \nabla \pi^h \theta_n - \kappa(\theta_n) \nabla \theta_n \right), \mathbf{u}^h \cdot \nabla \phi^h \right)_K \\ \leq \sum_{K \in \mathcal{T}^h} \delta_K C_{\text{inv}} h_K^{-1} \|\kappa^h(\hat{\theta}^h) \nabla \pi^h \theta_n - \kappa(\theta_n) \nabla \theta_n\|_{L^2(K)} \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}.$$

Using (53) and the inverse inequality (10) yields

$$\begin{aligned} & \|\kappa^h(\hat{\theta}^h) \nabla \pi^h \theta_n - \kappa(\theta_n) \nabla \theta_n\|_{L^2(K)} \\ & \leq \|\nabla \pi^h \theta_n\|_{L^\infty} \|\kappa^h(\hat{\theta}^h) - \kappa(\theta_n)\|_{L^2(K)} + \|\kappa(\theta_n)\|_{L^\infty(K)} \|\nabla(\theta_n - \pi^h \theta_n)\|_{L^2(K)} \\ & \leq C(\theta_n) \|\kappa^h(\hat{\theta}^h) - \kappa(\theta_n)\|_{L^2(K)} + C \kappa_{\max} h_K^{-1} \|\theta_n - \pi^h \theta_n\|_{L^2(K)}. \end{aligned}$$

Inserting this estimate in (57) and using the bound (46) of the stabilization parameter gives

$$\begin{aligned} & \sum_{K \in \mathcal{T}^h} \delta_K \left( \nabla \cdot \left( \kappa^h(\hat{\theta}^h) \nabla \pi^h \theta_n - \kappa(\theta_n) \nabla \theta_n \right), \mathbf{u}^h \cdot \nabla \phi^h \right)_K \\ & \leq C(\theta_n) \left[ \left( \frac{C}{\kappa_{\max}} \sum_{K \in \mathcal{T}^h} h_K \|\kappa^h(\hat{\theta}^h) - \kappa(\theta_n)\|_{L^2(K)}^2 \right)^{1/2} + \|\theta_n - \pi^h \theta_n\|_{L^2} \right] \\ & \quad \times \left( \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right)^{1/2}. \end{aligned}$$

For bounding the contribution of the convective term to the stabilization in (55), the same tools are used as for the previous terms, which leads to

$$\begin{aligned} & \sum_{K \in \mathcal{T}^h} \delta_K \left( \mathbf{u}_n \cdot \nabla \theta_n - \mathbf{u}^h \cdot \nabla \pi^h \theta_n, \mathbf{u}^h \cdot \nabla \phi^h \right)_K \\ & \leq \left[ \max_{K \in \mathcal{T}^h} \{\delta_K^{1/2}\} \|\nabla \theta_n\|_{L^\infty} \|\mathbf{u}_n - \mathbf{u}^h\|_{L^2} + \frac{\|\mathbf{u}^h\|_{L^\infty}}{\kappa_{\max}^{1/2}} \|\theta_n - \pi^h \theta_n\|_{L^2} \right] \\ & \quad \times \left( \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{u}^h \cdot \nabla \phi^h\|_{L^2(K)}^2 \right)^{1/2}. \end{aligned}$$

The bound of the first term with the truncation error in (55) starts with the Cauchy-Schwarz inequality

$$(T_n^h, \phi^h) \leq \kappa_{\min}^{-1/2} \|T_n^h\|_{L^2} \|\kappa^h(\hat{\theta}^h)^{1/2} \nabla \phi^h\|_{L^2}.$$

Then, the truncation error can be bounded using the same techniques as in [20, Eq. (4.3)] to give

$$\|T_n^h\|_{L^2} \leq \|(I - \pi^h) \partial_t \theta_n\|_{L^2} + C \left( \Delta t \int_{t_{n-1}}^{t_n} \|\partial_{tt} \theta\|_{H^1}^2 d\tau \right)^{1/2}.$$

With similar arguments and using (49), one obtains for the second term with the truncation error

$$\begin{aligned} & \sum_{K \in \mathcal{T}^h} \delta_K (T_n^h, \mathbf{u}^h \cdot \nabla \phi^h)_K \\ & \leq C \left( \Delta t \| (I - \pi^h) \partial_t \theta_n \|_{L^2}^2 + \Delta t^2 \int_{t_{n-1}}^{t_n} \| \partial_{tt} \theta \|_{H^1}^2 d\tau \right)^{1/2} \\ & \quad \times \left( \sum_{K \in \mathcal{T}^h} \delta_K \| \mathbf{u}^h \cdot \nabla \phi^h \|_{L^2(K)}^2 \right)^{1/2}. \end{aligned}$$

Last, note that the second term on the right-hand side of (55) can be absorbed in the left-hand side by using the first term on the left-hand side and a relation like (50).

The final steps consist in using a relation of the form (50) and the coercivity (47) to estimate the first two terms on the left-hand side of (55), by applying Young's inequality to all estimates such that the factors that belong to the SUPG norm are absorbed from the left-hand side, by summing over all discrete times, and then by applying the triangle inequality with respect to the decomposition of the error.  $\square$

**Remark 3.4.** Estimate (54) provides the information that  $\| \theta_n - \theta_n^h \|_{L^2}$  and a discrete analog of  $\| (\kappa^h)^{1/2} \nabla (\theta - \theta^h) \|_{L^2(L^2)}$  depend on a discrete version of

$$\frac{C(\theta)}{\kappa_{\min}^{1/2}} \| \mathbf{u} - \mathbf{u}^h \|_{L^2(L^2)}.$$

The term  $\kappa_{\min}^{1/2}$  can be expected to be very small in applications such that  $\| \mathbf{u} - \mathbf{u}^h \|_{L^2(L^2)}$  is scaled with a large factor. However, the term that describes the approximation of the thermal diffusivity and the interpolation error on the right-hand side of (54) are scaled with the same factor. From this point of view, all terms are of equal importance and  $\| \mathbf{u} - \mathbf{u}^h \|_{L^2(L^2)}$  should be preferably of the same order as  $\| \theta - \pi^h \theta \|_{L^2(L^2)}$ .

Error estimate (54) requires some interpretations.

The last term on the right-hand side of (54) gives the first order of convergence with respect to time, as it is expected from the backward Euler scheme.

The error bound is not uniform with respect to the thermal diffusivity. From the analytical point of view, the reason is that the usual assumption (reaction minus one half of the divergence of the convection) made for convection-diffusion problems, see the description at the beginning of Section 3, cannot be made. If this assumption holds, then the terms whose estimation leads to inverse powers of the thermal diffusivity in the proof of Theorem 3.3 could be bounded by using inverse powers of the positive constant that appears in the assumption. A change of variable as suggested in [11, Remark 1] could be applied to transform the equation into one of the same type satisfying the assumption. However, since the thermal diffusivity depends on the temperature, the thermal diffusivity of the new equation after the change of variables includes an exponential factor of type  $e^{\alpha t}$  for  $\alpha$  being a constant and the error bounds have to be carefully revised. For this reason, this will be subject of future research.

The interpolation errors  $\theta_j - \pi^h \theta_j$  and  $(I - \pi^h) \partial_t \theta_j$  appear at many occasions in the  $L^2(\Omega)$  norm such that a higher order of convergence can be expected for the respective terms, provided that the solution is sufficiently smooth. However, in the

second term on the right-hand side, the scaled  $H^1(\Omega)$  norm of the interpolation error appears. This term gives the order of convergence with respect to the mesh width.

For an optimal order of convergence, the error  $\|\mathbf{u}_j - \mathbf{u}^h\|_{L^2}$  has to be sufficiently small, i.e., the extrapolation that is used as  $\mathbf{u}^h$  has to be sufficiently accurate. It is known that appropriate IMEX schemes can be constructed, e.g., see [6].

Finally, the approximation error of the thermal diffusivity has to be considered. As mentioned above,  $\kappa^h$  was introduced for technical reasons and in practice, it is possible to evaluate  $\kappa(\hat{\theta}^h)$  at the quadrature points. For the term that appears in the error bound, the triangle inequality yields

$$(58) \quad \|\kappa(\theta_j) - \kappa^h(\hat{\theta}^h)\|_{L^2} \leq \|\kappa(\theta_j) - \kappa^h(\theta_j)\|_{L^2} + \|\kappa^h(\theta_j) - \kappa^h(\hat{\theta}^h)\|_{L^2}.$$

The first term on the right-hand side of (58) measures how good  $\kappa^h$  approximates  $\kappa$ . The second term should become small if the arguments of  $\kappa^h$  are in some sense close. This situation is given if  $\kappa^h$  is Lipschitz continuous with respect to the  $L^2(\Omega)$  norm

$$\|\kappa^h(\bar{\theta}) - \kappa^h(\tilde{\theta})\|_{L^2} \leq C\|\bar{\theta} - \tilde{\theta}\|_{L^2}$$

for all  $\bar{\theta}, \tilde{\theta} \in V_\theta$  and a constant  $C > 0$ . Considering concretely  $\|\theta_j - \hat{\theta}^h\|_{L^2}$ , it can be expected that this error is smaller if  $\hat{\theta}^h = \theta_j^h$ , i.e., Algorithm 1.1 is used, than if  $\hat{\theta}^h$  is some extrapolation from previous discrete times.

The appearance of  $\delta_{\min}^{-1}$  is already discussed in [20].

#### 4. Summary

This paper studied a model for mantle convection consisting of a coupled problem of the Stokes equations and a time-dependent convection-diffusion equation. A finite element analysis was performed for the individual equations, thereby tracking the dependency of the error bounds on the coefficients of the problem and on the finite element error coming from the other equation. In the following, realistic magnitudes of the coefficients will be assumed. Then, it was found that the temperature error possesses a large impact on the pressure error. The concrete dependency of the  $L^2(\Omega)$  error of the velocity on the temperature error is an open question. On the other hand, the velocity error in  $L^2(0, T; L^2(\Omega))$  has a large impact on the temperature error.

Considering just the order of convergence, then one can derive from the error estimates (32), (36), and (40) for the Stokes problem optimal orders for the temperature error. Considering the typical situation that the degree of the velocity finite element space is larger by one than the degree of the pressure finite element space, which is given, e.g., for Taylor–Hood pairs of finite element spaces, then one should choose the degree of the temperature finite element space equal to the degree of the velocity finite element space.

#### References

- [1] Gabriel R. Barrenechea, Volker John, and Petr Knobloch. A local projection stabilization finite element method with nonlinear crosswind diffusion for convection-diffusion-reaction equations. *ESAIM Math. Model. Numer. Anal.*, 47(5):1335–1366, 2013.
- [2] Susanne C. Brenner and L. Ridgway Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer, New York, third edition, 2008.
- [3] F. Brezzi and J. Pitkäranta. On the stabilization of finite element approximations of the Stokes equations. In *Efficient solutions of elliptic systems (Kiel, 1984)*, volume 10 of *Notes Numer. Fluid Mech.*, pages 11–19. Friedr. Vieweg, Braunschweig, 1984.

- [4] Alexander N. Brooks and Thomas J. R. Hughes. Streamline upwind/Petrov-Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier-Stokes equations. *Comput. Methods Appl. Mech. Engrg.*, 32(1-3):199–259, 1982. FENOMECH '81, Part I (Stuttgart, 1981).
- [5] Erik Burman and Miguel A. Fernández. Finite element methods with symmetric stabilization for the transient convection-diffusion-reaction equation. *Comput. Methods Appl. Mech. Engrg.*, 198(33-36):2508–2519, 2009.
- [6] M. P. Calvo, J. de Frutos, and J. Novo. Linearly implicit Runge-Kutta methods for advection-reaction-diffusion equations. *Appl. Numer. Math.*, 37(4):535–549, 2001.
- [7] Philippe G. Ciarlet. *The finite element method for elliptic problems*. North-Holland Publishing Co., Amsterdam, 1978. Studies in Mathematics and its Applications, Vol. 4.
- [8] Ramon Codina and Jordi Blasco. Analysis of a stabilized finite element approximation of the transient convection-diffusion-reaction equation using orthogonal subscales. *Comput. Vis. Sci.*, 4(3):167–174, 2002.
- [9] Juliane Dannberg and Timo Heister. Compressible magma/mantle dynamics: 3-d, adaptive simulations in ASPECT. *Geophysical Journal International*, 207(3):1343–1366, 2016.
- [10] Timothy A. Davis. Algorithm 832: UMFPACK V4.3—an unsymmetric-pattern multifrontal method. *ACM Trans. Math. Software*, 30(2):196–199, 2004.
- [11] Javier de Frutos, Bosco García Archilla, and Julia Novo. Stabilization of Galerkin finite element approximations to transient convection-diffusion problems. *SIAM J. Numer. Anal.*, 48(3):953–979, 2010.
- [12] Javier de Frutos, Bosco García-Archilla, Volker John, and Julia Novo. Grad-div stabilization for the evolutionary Oseen problem with inf-sup stable finite elements. *J. Sci. Comput.*, 66(3):991–1024, 2016.
- [13] V. Girault and L. R. Scott. A quasi-local interpolation operator preserving the discrete divergence. *Calcolo*, 40(1):1–19, 2003.
- [14] Piotr P. Grinevich and Maxim A. Olshanskii. An iterative method for the Stokes-type problem with variable viscosity. *SIAM J. Sci. Comput.*, 31(5):3959–3978, 2009.
- [15] C. O. Horgan. Korn's inequalities and their applications in continuum mechanics. *SIAM Rev.*, 37(4):491–511, 1995.
- [16] T. J. R. Hughes and A. Brooks. A multidimensional upwind scheme with no crosswind diffusion. In *Finite element methods for convection dominated flows (Papers, Winter Ann. Meeting Amer. Soc. Mech. Engrs., New York, 1979)*, volume 34 of *AMD*, pages 19–35. Amer. Soc. Mech. Engrs. (ASME), New York, 1979.
- [17] Volker John. *Finite Element Methods for Incompressible Flow Problems*, volume 51 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2016.
- [18] Volker John, Kristine Kaiser, and Julia Novo. Finite element methods for the incompressible Stokes equations with variable viscosity. *ZAMM Z. Angew. Math. Mech.*, 96(2):205–216, 2016.
- [19] Volker John and Gunar Matthies. MoonNMD—a program package based on mapped finite element methods. *Comput. Vis. Sci.*, 6(2-3):163–169, 2004.
- [20] Volker John and Julia Novo. Error analysis of the SUPG finite element discretization of evolutionary convection-diffusion-reaction equations. *SIAM J. Numer. Anal.*, 49(3):1149–1176, 2011.
- [21] Volker John and Ellen Schmeier. Finite element methods for time-dependent convection-diffusion-reaction equations with small diffusion. *Comput. Methods Appl. Mech. Engrg.*, 198(3-4):475–494, 2008.
- [22] Ricardo Oyarzúa, Tong Qin, and Dominik Schötzau. An exactly divergence-free finite element method for a generalized Boussinesq problem. *IMA J. Numer. Anal.*, 34(3):1104–1135, 2014.
- [23] Carlos E. Pérez, Jean-Marie Thomas, Serge Blancher, and René Creff. The steady Navier-Stokes/energy system with temperature-dependent viscosity. II. The discrete problem and numerical experiments. *Internat. J. Numer. Methods Fluids*, 56(1):91–114, 2008.
- [24] Hans-Görg Roos, Martin Stynes, and Lutz Tobiska. *Robust numerical methods for singularly perturbed differential equations*, volume 24 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, second edition, 2008. Convection-diffusion-reaction and flow problems.
- [25] Youcef Saad. A flexible inner-outer preconditioned GMRES algorithm. *SIAM J. Sci. Comput.*, 14(2):461–469, 1993.
- [26] M. Tabata and A. Suzuki. A stabilized finite element method for the Rayleigh-Bénard equations with infinite Prandtl number in a spherical shell. *Comput. Methods Appl. Mech. Engrg.*, 190(3-4):387–402, 2000.

- [27] Masahisa Tabata. Finite element approximation to infinite Prandtl number Boussinesq equations with temperature-dependent coefficients – thermal convection problems in a spherical shell. *Future Generation Computer Systems*, 22(4):521 – 531, 2006.
- [28] Masahisa Tabata and Daisuke Tagami. Error estimates of finite element methods for nonstationary thermal convection problems with temperature-dependent coefficients. *Numer. Math.*, 100(2):351–372, 2005.
- [29] Vidar Thomée. *Galerkin finite element methods for parabolic problems*, volume 25 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, second edition, 2006.

Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstr. 39, 10117 Berlin, Germany and Freie Universität Berlin, Department of Mathematics and Computer Science, Arnimallee 6, 14195 Berlin, Germany

*E-mail:* john@wias-berlin.de

*URL:* <http://www.wias-berlin.de/people/john/>

Department of Mathematics, Middle East Technical University, 06800 Ankara, Turkey

*E-mail:* smerdan@metu.edu.tr

*URL:* <http://www.metu.edu.tr/~smerdan>

Departamento de Matemáticas, Universidad Autónoma de Madrid, Madrid, Spain. Research supported by Spanish MINECO under grants MTM2013-42538-P (MINECO, ES) and MTM2016-78995-P (AEI/FEDER, UE)

*E-mail:* julia.novo@uam.es

*URL:* <http://verso.mat.uam.es/web/index.php/es/directorio/24-pdi/61-novo-martin-julia>