

A LOCAL PROJECTION STABILIZATION FINITE ELEMENT METHOD WITH NONLINEAR CROSSWIND DIFFUSION FOR CONVECTION-DIFFUSION-REACTION EQUATIONS

GABRIEL R. BARRENECHEA, VOLKER JOHN, AND PETR KNOBLOCH

ABSTRACT. An extension of the local projection stabilization (LPS) finite element method for convection-diffusion-reaction equations is presented and analyzed, both in the steady-state and the transient setting. In addition to the standard LPS method, a nonlinear crosswind diffusion term is introduced that accounts for the reduction of spurious oscillations. The existence of a solution can be proved and, depending on the choice of the stabilization parameter, also its uniqueness. Error estimates are derived which are supported by numerical studies. These studies demonstrate also the reduction of the spurious oscillations.

1. INTRODUCTION

The solution of convection-dominated convection-diffusion-reaction equations with finite element methods constitutes a very challenging (and open) problem. Over the last three decades, the amount of work devoted to this problem is impressive. The usual way of treating dominating convection, at least in the context of finite element methods, consists in adding extra terms to the standard Galerkin formulation, aimed at enhancing the stability of the discrete solution by means of introducing artificial diffusion. These new terms vary according to the method, and can be residual-based, as in the SUPG/GLS/SDFEM family (see [6, 16, 13, 14, 29]), or edge based, such as the CIP method (see [9, 7]). For an up-to-date and thorough review of these and other techniques, see [31]. It is striking to notice that, despite the impressive amount of work that has been devoted to this topic, up to now there is not a method that 'ticks all the boxes', i.e., a method that produces sharp layers while avoiding oscillations, see [1] for a recent review and a numerical assessment.

Among the various stabilized finite element methods, the local projection stabilization (LPS) method has received some attention over the last decade. Originally proposed for the Stokes problem in [2], and extended to the Oseen equations in [4] (see also [5, 30]), the LPS method has been also used recently to treat convection-diffusion equations (see

Date: March 5, 2013.

Key words and phrases. finite element method; local projection stabilization; crosswind diffusion; convection-diffusion-reaction equation; well posedness; time dependent problem; stability; error estimates.

[26, 15, 24, 25]). The basic idea of this method consists in restricting the direct application of the stabilization to so-called fluctuations or resolved small scales, which are defined by local projections. It has several attractive features, such as adding symmetric terms to the formulation and avoiding the computation of second derivatives of the basis functions (thus using only information that is needed for the assembly of the matrices from the standard Galerkin method). Unfortunately, the solutions obtained with the LPS method possess the same deficiency like solutions computed, e.g., with the SUPG method: non-negligible spurious oscillations are often present in a vicinity of layers.

Motivated by the wish of recovering the monotonicity properties of the continuous problem, which might be crucial in applications, a number of so-called Spurious Oscillations at Layers Diminishing (SOLD) methods were proposed. SOLD methods add an extra term to the already stabilized formulation, which usually depends on the discrete solution in a nonlinear way, vanishes for small residuals (thus acting mostly at layers), and adds some extra, but different, diffusivity to the formulation. In particular, methods that add crosswind diffusion, like the one proposed in [11], have been proved to belong to the best SOLD methods in comprehensive studies [17, 18]. Although these methods diminish oscillations considerably, no single method succeeds to fully eliminate them [17, 18, 23]. Also, from a purely mathematical point of view, it is unknown if these methods lead to well-posed problems. In fact, existence of solutions is usually possible to prove, but, to our best knowledge, there is no nonlinear SOLD method that is known to produce a unique solution, see [27] and [7] for a discussion of this topic.

Based on the previous considerations, this paper has three major objectives, namely:

- to improve the quality of the LPS solution (especially in the vicinity of layers);
- to explore the applicability of SOLD-type strategies within a LPS context; and
- to contribute to the mathematical understanding of nonlinear stabilization techniques for the convection-diffusion equation.

Hence, in this work we propose a LPS method with nonlinear crosswind diffusion for convection-diffusion-reaction equations. Two ways for choosing the parameter in the crosswind diffusion term will be studied. The first choice uses global information obtained from the data of the problem, whereas the second proposal is completely local, employing information of the computed solution instead of the data. For the first approach, which is the simpler one, the existence and the uniqueness of the solution can be proved for the steady-state and time-dependent equations, where the latter is discretized in time with an implicit one-step

θ -scheme. To our best knowledge, this is the first nonlinear discretization for convection-diffusion-reaction equations for which both, existence and uniqueness of a solution can be shown. The form of the crosswind term resembles the Smagorinsky Large Eddy Simulation (LES) model which was analyzed in [28]. It involves fluctuations of a term mimicking a p -Laplacian. The crucial analytical property for proving the uniqueness of the solution is the strong monotonicity of the corresponding operator. For the more complicated local definition of the parameter, the analysis will show the existence of a solution and its uniqueness for the time-dependent discretization in the case of sufficiently small time steps.

The analysis is performed for the model problems of linear steady-state and time-dependent convection-diffusion-reaction equations. Applying a nonlinear discretization scheme to a linear problem leads certainly to a considerable complication of the solution process and to an additional numerical cost. This latter aspect can be overcome in the transient regime by using a semi-implicit (linearized) approach that computes the stabilization parameter with the solution from the previous discrete time. With respect to the former aspect, it has to be mentioned that the most important motivation for studying discretizations that reduce spurious oscillations comes from the need to address applications that lead to nonlinear coupled systems of convection-diffusion-reaction equations as in [21]. It was demonstrated in [21] that the locally large spurious oscillations of the SUPG method might lead to a fast blow-up of the simulations, and hence the reduction of the spurious oscillations is essential to perform simulations at all. Thus, the reduction of the oscillations at layers becomes a priority, even over computational cost. It should be noted that in many applications, like in [21], only interior or characteristic layers are present, such that a method for reducing the oscillations has to work properly in particular for these types of layers. Finally, it is worth mentioning that our final aim is to address applications that lead to such coupled problems. Since these problems are nonlinear, the use of a nonlinear stabilization usually does not result in a notable complication of the solution procedure.

The plan of the paper is as follows. In the remaining part of this introduction, the problems of interest are stated and some basic notations are given. Section 2 will summarize the main abstract hypothesis imposed on the different partitions of the domain and the finite element spaces considered. Section 3 presents the method for the steady-state case, for which well-posedness is analyzed in Section 3.1 and error estimates are proved in Section 3.2. In Section 4, the method for the time-dependent problem is presented. Well-posedness and stability are proved in Section 4.1 and error estimates in Section 4.2. Since the analysis is based on the abstract framework from Section 2, Section 5 presents some concrete examples that

fit into this framework. Finally, numerical illustrations that support the analytical results and which demonstrate the reduction of spurious oscillations are presented in Section 6.

Throughout the paper, standard notations are used for Sobolev spaces and corresponding norms, see, e.g., [10]. In particular, given a measurable set $D \subset \mathbb{R}^d$, the inner product in $L^2(D)$ or $L^2(D)^d$ is denoted by $(\cdot, \cdot)_D$ and the notation (\cdot, \cdot) is used instead of $(\cdot, \cdot)_\Omega$. The norm (seminorm) in $W^{m,p}(D)$ will be denoted by $\|\cdot\|_{m,p,D}$ ($|\cdot|_{m,p,D}$), with the convention $\|\cdot\|_{m,D} = \|\cdot\|_{m,2,D}$, and the same notation is used for scalar and vector-valued functions.

1.1. The problems of interest. Let $\Omega \subset \mathbb{R}^d$, $d \in \{2, 3\}$, be a bounded polygonal (polyhedral) domain with a Lipschitz-continuous boundary $\partial\Omega$ and let us consider the steady-state convection-diffusion-reaction equation

$$(1) \quad -\varepsilon \Delta u + \mathbf{b} \cdot \nabla u + c u = f \quad \text{in } \Omega, \quad u = u_b \quad \text{on } \partial\Omega.$$

It is assumed that ε is a positive constant and $\mathbf{b} \in W^{1,\infty}(\Omega)^d$, $c \in L^\infty(\Omega)$, $f \in L^2(\Omega)$, and $u_b \in H^{1/2}(\partial\Omega)$ are given functions satisfying

$$(2) \quad \sigma := c - \frac{1}{2} \nabla \cdot \mathbf{b} \geq \sigma_0 > 0 \quad \text{in } \Omega,$$

where σ_0 is a constant. Then the boundary value problem (1) has a unique solution in $H^1(\Omega)$.

The condition $\sigma_0 > 0$ is often used in the analysis of stabilized finite element methods for the numerical solution of (1), see, e.g., [31], but it limits the applications of the theory since many problems of interest involve solenoidal convective velocities and no zero-order terms, which leads to $\sigma_0 = 0$. Unfortunately, it is not known how to prove optimal convergence results even for the underlying linear local projection stabilization without assuming $\sigma_0 > 0$, although numerical results do not indicate any deterioration of the convergence rates when $\sigma_0 = 0$. The analysis of the nonlinear term introduced in this paper does not require this assumption.

Besides the steady-state case, also the time-dependent convection-diffusion-reaction equation

$$(3) \quad \begin{cases} u_t - \varepsilon \Delta u + \mathbf{b} \cdot \nabla u + c u = f & \text{in } (0, T] \times \Omega, \\ u = u_b & \text{in } [0, T] \times \partial\Omega, \\ u(0, \cdot) = u_0 & \text{in } \Omega, \end{cases}$$

will be considered. In (3), $[0, T]$ is a finite time interval, ε is assumed to be a positive constant, $\mathbf{b} \in L^\infty(0, T; W^{1,\infty}(\Omega)^d)$, $c \in L^\infty(0, T; L^\infty(\Omega))$, $f \in L^2(0, T; L^2(\Omega))$, $u_b \in L^2(0, T; H^{1/2}(\partial\Omega))$, and $u_0 \in H^1(\Omega)$ denotes the initial condition. The function σ is defined analogously to (2)

and the inequality (2) is assumed to hold for all $t \in [0, T]$. In this case, the condition $\sigma_0 > 0$ can be circumvented by considering instead of (3) an equivalent problem for $v = u e^{-\alpha t}$ which satisfies $\sigma_0 > 0$ for sufficiently large α .

2. ASSUMPTIONS ON APPROXIMATION SPACES AND THE SET \mathcal{M}_h

From now on, C , \tilde{C} or \bar{C} denote generic constants which may take different values at different occurrences but are always independent of the data ε , \mathbf{b} , c , f , and u_b , the constant σ_0 , and the discretization parameters (h and δt in the following).

Given $h > 0$, let $W_h \subset W^{1,\infty}(\Omega)$ be a finite-dimensional space approximating the space $H^1(\Omega)$ and set $V_h = W_h \cap H_0^1(\Omega)$. Next, let \mathcal{M}_h be a set consisting of a finite number of open subsets M of Ω such that $\bar{\Omega} = \cup_{M \in \mathcal{M}_h} \bar{M}$. It will be supposed that, for any $M \in \mathcal{M}_h$,

$$(4) \quad \text{card}\{M' \in \mathcal{M}_h; M \cap M' \neq \emptyset\} \leq C,$$

$$(5) \quad h_M := \text{diam}(M) \leq C h,$$

$$(6) \quad h_M \leq C h_{M'} \quad \forall M' \in \mathcal{M}_h, M \cap M' \neq \emptyset,$$

$$(7) \quad h_M^d \leq C \text{meas}_d(M).$$

The space W_h is assumed to satisfy the local inverse inequality

$$(8) \quad |v_h|_{1,M} \leq C h_M^{-1} \|v_h\|_{0,M} \quad \forall v_h \in W_h, M \in \mathcal{M}_h.$$

For any $M \in \mathcal{M}_h$, a finite-dimensional space $D_M \subset L^\infty(M)$ is introduced. It is assumed that there exists a positive constant β_{LP} independent of h such that

$$(9) \quad \sup_{v \in V_M} \frac{(v, q)_M}{\|v\|_{0,M}} \geq \beta_{LP} \|q\|_{0,M} \quad \forall q \in D_M, M \in \mathcal{M}_h,$$

where $V_M = \{v_h \in V_h; v_h = 0 \text{ in } \Omega \setminus M\}$. This hypothesis will be needed in what follows for the construction of a special interpolation operator (see Lemma 6 below). Concrete examples of spaces W_h and D_M satisfying the assumptions formulated here will be presented in Section 5.

Furthermore, for any $M \in \mathcal{M}_h$, a finite-dimensional space $G_M \subset L^\infty(M)$ with $G_M \supset D_M$ is introduced such that

$$\left. \frac{\partial v_h}{\partial x_i} \right|_M \in G_M \quad \forall v_h \in W_h, i = 1, \dots, d,$$

and it is assumed that, for any $p \in [1, \infty]$, there is a constant C such that

$$(10) \quad \|q\|_{0,p,M} \leq C h_M^{\frac{d}{p} - \frac{d}{2}} \|q\|_{0,M} \quad \forall q \in G_M, M \in \mathcal{M}_h.$$

To characterize the approximation properties of the spaces W_h and D_M , it is assumed that there exist interpolation operators $i_h \in \mathcal{L}(C(\overline{\Omega}), W_h) \cap \mathcal{L}(C(\overline{\Omega}) \cap H_0^1(\Omega), V_h)$ and $j_M \in \mathcal{L}(H^1(M), D_M)$, $M \in \mathcal{M}_h$, such that, for some constants $l \in \mathbb{N}$ and $C > 0$ and for any set $M \in \mathcal{M}_h$, it holds

$$(11) \quad |v - i_h v|_{1,M} + h_M^{-1} \|v - i_h v\|_{0,M} \leq C h_M^k |v|_{k+1,M} \quad \forall v \in H^{k+1}(M), \quad k = 1, \dots, l,$$

$$(12) \quad \|q - j_M q\|_{0,M} \leq C h_M^k |q|_{k,M} \quad \forall q \in H^k(M), \quad k = 1, \dots, l.$$

In addition, it is assumed that, for any $p \in [1, 6]$,

$$(13) \quad |v - i_h v|_{1,p,M} \leq C h_M^{k+\frac{d}{p}-\frac{d}{2}} |v|_{k+1,M} \quad \forall v \in H^{k+1}(M), \quad k = 1, \dots, l.$$

3. A LOCAL PROJECTION DISCRETIZATION OF THE STEADY-STATE PROBLEM

The weak form of problem (1) is: Find $u \in H^1(\Omega)$ such that $u = u_b$ on $\partial\Omega$ and

$$(14) \quad a(u, v) = (f, v) \quad \forall v \in H_0^1(\Omega),$$

where the bilinear form a is given by

$$a(u, v) := \varepsilon (\nabla u, \nabla v) + (\mathbf{b} \cdot \nabla u, v) + (c u, v).$$

As it was mentioned in the introduction, the most often used approach to cure the instabilities of the Galerkin method consists in adding extra terms to the formulation. To build these additional terms for the method studied here, for any $M \in \mathcal{M}_h$, a continuous linear projection operator π_M is introduced which maps the space $L^2(M)$ onto the space D_M . It is assumed that

$$(15) \quad \|\pi_M\|_{\mathcal{L}(L^2(M), L^2(M))} \leq C \quad \forall M \in \mathcal{M}_h.$$

E.g., if π_M is the orthogonal L^2 projection, then $C = 1$. Using this operator, the fluctuation operator $\kappa_M := id - \pi_M$ is defined, where id is the identity operator on $L^2(M)$. Then, clearly

$$(16) \quad \|\kappa_M\|_{\mathcal{L}(L^2(M), L^2(M))} \leq C \quad \forall M \in \mathcal{M}_h.$$

Since κ_M vanishes on D_M , it follows from (16) and (12) that

$$(17) \quad \|\kappa_M q\|_{0,M} \leq C h_M^k |q|_{k,M} \quad \forall q \in H^k(M), \quad M \in \mathcal{M}_h, \quad k = 0, \dots, l.$$

An application of κ_M to a vector-valued function means that κ_M is applied component-wise.

For any $M \in \mathcal{M}_h$, a constant $\mathbf{b}_M \in \mathbb{R}^d$ is chosen such that

$$(18) \quad |\mathbf{b}_M| \leq \|\mathbf{b}\|_{0,\infty,M}, \quad \|\mathbf{b} - \mathbf{b}_M\|_{0,\infty,M} \leq C h_M |\mathbf{b}|_{1,\infty,M},$$

where $|\cdot|$ denotes the Euclidean norm in \mathbb{R}^d . A typical choice for \mathbf{b}_M is the value of \mathbf{b} at one point of M , or the integral mean value of \mathbf{b} over M . In addition, a function $\tilde{u}_{bh} \in W_h$ is introduced such that its trace approximates the boundary condition u_b .

We are now ready to present the finite element method to be studied: Find $u_h \in W_h$ such that $u_h - \tilde{u}_{bh} \in V_h$ and

$$(19) \quad a(u_h, v_h) + s_h(u_h, v_h) + d_h(u_h; u_h, v_h) = (f, v_h) \quad \forall v_h \in V_h,$$

where

$$\begin{aligned} s_h(u, v) &= \sum_{M \in \mathcal{M}_h} \tau_M (\kappa_M(\mathbf{b}_M \cdot \nabla u), \kappa_M(\mathbf{b}_M \cdot \nabla v))_M, \\ d_h(w; u, v) &= \sum_{M \in \mathcal{M}_h} (\tau_M^{\text{sold}}(w) \kappa_M(P_M \nabla u), \kappa_M(P_M \nabla v))_M, \end{aligned}$$

and $P_M : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the projection onto the line (plane) orthogonal (crosswind) to the vector \mathbf{b}_M defined by

$$P_M = \begin{cases} I - \frac{\mathbf{b}_M \otimes \mathbf{b}_M}{|\mathbf{b}_M|^2} & \text{if } \mathbf{b}_M \neq \mathbf{0}, \\ 0 & \text{if } \mathbf{b}_M = \mathbf{0}, \end{cases}$$

I being the identity tensor. The stabilization parameters are given by

$$(20) \quad \begin{aligned} \tau_M &= \tau_0 \min \left\{ \frac{h_M}{\|\mathbf{b}\|_{0,\infty,M}}, \frac{h_M^2}{\varepsilon} \right\}, \\ \tau_M^{\text{sold}}(u_h) &= \tilde{\tau}_M(u_h) |\kappa_M(P_M \nabla u_h)|, \end{aligned}$$

where τ_0 is a positive constant and $\tilde{\tau}_M$ is a non-negative function of u_h and the data of (1). Note that the crosswind stabilization term is of p -Laplacian type with $p = 3$.

It remains to specify the function $\tilde{\tau}_M$. First, inspired by the definition of s_h , where each term in the sum is bounded by $\tau_0 h_M |\mathbf{b}_M| \|\kappa_M \nabla u\|_{0,M} \|\kappa_M \nabla v\|_{0,M}$, we set $\tilde{\tau}_M(u_h) = \gamma_M(u_h) h_M |\mathbf{b}_M|$ with a function γ_M still depending on u_h and/or the data of (1). Second, the function γ_M has to be chosen in such a way that the discrete problem preserves the following scaling properties of the problem (1):

- if the data ε , \mathbf{b} , c , and f are replaced by $\alpha \varepsilon$, $\alpha \mathbf{b}$, αc , and αf , respectively, with some constant $\alpha \neq 0$, then the solution of (1) does not change;
- if f and u_b are replaced by αf and αu_b , respectively, then u changes to αu ;
- if Ω is transformed to $F^{-1}(\Omega)$ with $F(x) = x/\alpha$, then $u \circ F$ solves an analog of (1) in $F^{-1}(\Omega)$ with the data $\alpha^2 \varepsilon$, $\alpha \mathbf{b} \circ F$, $c \circ F$, $f \circ F$, and $u_b \circ F$.

Note that the discrete problem (19) without the nonlinear term d_h preserves these properties. To preserve the properties also when using the nonlinear term, the function γ_M has to satisfy

$$\begin{aligned}\gamma_M(\varepsilon, \mathbf{b}, c, f, u_b, \Omega, u_h) &= \gamma_M(\alpha \varepsilon, \alpha \mathbf{b}, \alpha c, \alpha f, u_b, \Omega, u_h) \\ &= \alpha \gamma_M(\varepsilon, \mathbf{b}, c, \alpha f, \alpha u_b, \Omega, \alpha u_h) \\ &= \alpha^{-1} \gamma_{F^{-1}(M)}(\alpha^2 \varepsilon, \alpha \mathbf{b} \circ F, c \circ F, f \circ F, u_b \circ F, F^{-1}(\Omega), u_h \circ F)\end{aligned}$$

for any admissible data, $\alpha \neq 0$, and $u_h \in W_h$. We shall consider two choices of the scaling function γ_M : a global one independent of u_h and a local one depending on u_h . In the former case, one may set

$$(21) \quad \gamma_M = \gamma_0 \text{diam}(\Omega)^{d/2} \left(\frac{\|f\|_{0,\Omega} \text{diam}(\Omega)}{\varepsilon + \|\mathbf{b}\|_{0,\infty,\Omega} \text{diam}(\Omega) + \|c\|_{0,\infty,\Omega} \text{diam}(\Omega)^2} + \frac{\|u_b\|_{0,\partial\Omega}}{\text{diam}(\Omega)^{1/2}} \right)^{-1}$$

with a positive constant γ_0 . The local scaling can be defined by setting $\gamma_M = \beta h_M^{d/2} / |u_h|_{1,M}$ with a positive constant β if $|u_h|_{1,M} \neq 0$. Thus, we arrive at the following two formulas for the function $\tilde{\tau}_M$:

$$(22) \quad \tilde{\tau}_M = \beta h_M |\mathbf{b}_M|,$$

and

$$(23) \quad \tilde{\tau}_M(u_h) = \begin{cases} \frac{\beta h_M^{1+d/2} |\mathbf{b}_M|}{|u_h|_{1,M}} & \text{if } |u_h|_{1,M} \neq 0, \\ 0 & \text{if } |u_h|_{1,M} = 0, \end{cases}$$

where β is a positive real number independent of u_h and h . The parameter β depends on the data of (1) in case of (22) (e.g., like γ_M in (21)), but it is independent of the data of (1) in case of (23). For these two choices of $\tilde{\tau}_M$, we shall investigate the properties of the discrete problem (19). Although the local scaling is likely to lead to better numerical results than the global one, we consider both variants since the choice (22) turns out to be more appealing for the analysis.

Remark.

- If $d = 2$ and $\mathbf{b}_M \neq \mathbf{0}$, one has $P_M = \mathbf{b}_M^\perp \otimes \mathbf{b}_M^\perp$ where \mathbf{b}_M^\perp is a vector satisfying $\mathbf{b}_M^\perp \cdot \mathbf{b}_M = 0$ and $|\mathbf{b}_M^\perp| = 1$. Thus, in this case, the nonlinear stabilization term can be written in the form

$$d_h(w; u, v) = \sum_{M \in \mathcal{M}_h} (\tau_M^{\text{sold}}(w) \kappa_M(\mathbf{b}_M^\perp \cdot \nabla u), \kappa_M(\mathbf{b}_M^\perp \cdot \nabla v))_M.$$

- It is useful for the analysis of the discrete problem to note that $\kappa_M(\mathbf{b}_M \cdot \nabla u) = \mathbf{b}_M \cdot \kappa_M \nabla u$ and $\kappa_M(P_M \nabla u) = P_M \kappa_M \nabla u$. Note also that $\|P_M\|_2 = 1$.
- Finally, if $\tilde{\tau}_M$ is defined by (23), then, using the stability of κ_M and \mathbf{b}_M (18) and (16), respectively, and $\|P_M\|_2 = 1$, one obtains

$$(24) \quad \|\tau_M^{\text{sold}}(v)\|_{0,M} \leq C h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \quad \forall v \in H^1(\Omega), M \in \mathcal{M}_h.$$

In the analysis, the error will be measured using the following mesh-dependent norm

$$\|v\|_{\text{LPS}} := (\varepsilon |v|_{1,\Omega}^2 + \|\sigma^{1/2} v\|_{0,\Omega}^2 + s_h(v, v))^{1/2},$$

and a term involving the crosswind derivative of the error. Note that integrating by parts gives

$$(25) \quad a(v, v) + s_h(v, v) = \|v\|_{\text{LPS}}^2 \quad \forall v \in H_0^1(\Omega).$$

3.1. Well-posedness of the nonlinear discrete problem. This section studies the existence and uniqueness of solutions for the nonlinear discrete problem (19). The results of this section are valid also for $\sigma_0 = 0$.

Let us define the nonlinear operator $T_h : V_h \rightarrow V_h$ by

$$(26) \quad (T_h z_h, v_h) = a(z_h + \tilde{u}_{bh}, v_h) + s_h(z_h + \tilde{u}_{bh}, v_h) + d_h(z_h + \tilde{u}_{bh}; z_h + \tilde{u}_{bh}, v_h) - (f, v_h)$$

for any $z_h, v_h \in V_h$. Then $u_h \in W_h$ is a solution of (19) if and only if $u_h|_{\partial\Omega} = \tilde{u}_{bh}|_{\partial\Omega}$ and

$$T_h(u_h - \tilde{u}_{bh}) = 0,$$

or, equivalently, $u_h = \tilde{u}_h + \tilde{u}_{bh} \in W_h$ is a solution of (19) if $\tilde{u}_h \in V_h$ and $T_h(\tilde{u}_h) = 0$. Thus, our aim is to prove that the operator T_h has a zero in V_h . To this end, the properties of the form d_h shall be investigated first. As these properties are different with respect to the definition of $\tilde{\tau}_M$, we start supposing that $\tilde{\tau}_M$ is given by (22).

Lemma 1. *Let $\tilde{\tau}_M$ be defined by (22). Consider any $u, v, z \in W^{1,3}(\Omega)$ and set $w := u - v$. Then*

$$(27) \quad d_h(u; u, w) - d_h(v; v, w) \geq \frac{1}{7} \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla w)\|_{0,3,M}^3 = \frac{1}{7} d_h(w; w, w),$$

$$(28) \quad |d_h(u; u, z) - d_h(v; v, z)| \leq \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M (\|\kappa_M(P_M \nabla u)\|_{0,3,M} + \|\kappa_M(P_M \nabla v)\|_{0,3,M}) \\ \times \|\kappa_M(P_M \nabla w)\|_{0,3,M} \|\kappa_M(P_M \nabla z)\|_{0,3,M}.$$

Proof. Let us denote

$$(29) \quad d_h(u; u, z) - d_h(v; v, z) = \sum_{M \in \mathcal{M}_h} N_M(u, v, z),$$

where

$$N_M(u, v, z) := \left(\tau_M^{\text{sold}}(u) \kappa_M(P_M \nabla u) - \tau_M^{\text{sold}}(v) \kappa_M(P_M \nabla v), \kappa_M(P_M \nabla z) \right)_M.$$

For $t \in [0, 1]$, let $u^t := tu + (1-t)v$ and set

$$g(t) := \tilde{\tau}_M |\kappa_M(P_M \nabla u^t)| \kappa_M(P_M \nabla u^t), \quad t \in [0, 1].$$

Then

$$N_M(u, v, z) = \left(g(1) - g(0), \kappa_M(P_M \nabla z) \right)_M = \left(\int_0^1 g'(t) dt, \kappa_M(P_M \nabla z) \right)_M.$$

Since

$$(30) \quad g'(t) = \tilde{\tau}_M \frac{\kappa_M(P_M \nabla u^t)}{|\kappa_M(P_M \nabla u^t)|} \kappa_M(P_M \nabla u^t) \cdot \kappa_M(P_M \nabla w) + \tilde{\tau}_M |\kappa_M(P_M \nabla u^t)| \kappa_M(P_M \nabla w),$$

one has

$$\begin{aligned} |g'(t)| &\leq 2 \tilde{\tau}_M |\kappa_M(P_M \nabla u^t)| |\kappa_M(P_M \nabla w)| \\ &\leq 2 \tilde{\tau}_M (t |\kappa_M(P_M \nabla u)| + (1-t) |\kappa_M(P_M \nabla v)|) |\kappa_M(P_M \nabla w)|, \end{aligned}$$

which implies (28). On the other hand, since multiplication of the first term on the right-hand side of (30) by $\kappa_M(P_M \nabla w)$ gives a non-negative expression, one obtains

$$(31) \quad N_M(u, v, w) \geq \left(\tilde{\tau}_M \int_0^1 |\kappa_M(P_M \nabla u^t)| dt \kappa_M(P_M \nabla w), \kappa_M(P_M \nabla w) \right)_M.$$

Next, clearly

$$\int_0^1 |\kappa_M(P_M \nabla u^t)| dt \geq \max_{i=1, \dots, d} \int_0^1 |t \kappa_M(P_M \nabla u)_i + (1-t) \kappa_M(P_M \nabla v)_i| dt.$$

Denoting

$$I(a, b) = \int_0^1 |ta + (1-t)b| dt, \quad a, b \in \mathbb{R},$$

a direct computation gives

$$I(a, b) = \frac{|a| + |b|}{2} \quad \text{if } ab \geq 0, \quad I(a, b) = \frac{1}{2} \frac{a^2 + b^2}{|a| + |b|} \quad \text{if } ab < 0.$$

Thus, for any $a, b \in \mathbb{R}$, it follows

$$I(a, b) \geq \frac{|a| + |b|}{4} \geq \frac{|a - b|}{4}.$$

Consequently,

$$\int_0^1 |\kappa_M(P_M \nabla u^t)| dt \geq \frac{1}{4} \max_{i=1, \dots, d} |\kappa_M(P_M \nabla w)_i| \geq \frac{1}{4\sqrt{d}} |\kappa_M(P_M \nabla w)| \geq \frac{1}{7} |\kappa_M(P_M \nabla w)|.$$

Combining this estimate with (31) and using (29) gives (27). \square

Next, the properties of d_h are explored for the case that $\tilde{\tau}_M$ is defined by (23).

Lemma 2. *Let $\tilde{\tau}_M$ be defined by (23). Consider any $u, v, z \in W^{1,4}(\Omega)$. Then*

$$(32) \quad |d_h(u; v, z)| \leq C \sum_{M \in \mathcal{M}_h} h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \|\kappa_M(P_M \nabla v)\|_{0,4,M} \|\kappa_M(P_M \nabla z)\|_{0,4,M},$$

$$(33) \quad |d_h(u; u, z) - d_h(v; v, z)| \leq C \sum_{M \in \mathcal{M}_h} h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \zeta_M(u, v) \times \\ \times (\|\kappa_M(P_M \nabla u)\|_{0,4,M} + \|\kappa_M(P_M \nabla v)\|_{0,4,M}) \|\kappa_M(P_M \nabla z)\|_{0,4,M},$$

where

$$\zeta_M(u, v) = \begin{cases} \frac{|u - v|_{1,M}}{|u|_{1,M} + |v|_{1,M}} & \text{if } |u|_{1,M} \neq 0 \text{ or } |v|_{1,M} \neq 0, \\ 0 & \text{if } |u|_{1,M} = |v|_{1,M} = 0. \end{cases}$$

Proof. Denoting

$$d_M(u; v, z) = (\tau_M^{\text{sold}}(u) \kappa_M(P_M \nabla v), \kappa_M(P_M \nabla z))_M,$$

it is easy to realize that

$$d_h(u; v, z) = \sum_{M \in \mathcal{M}_h} d_M(u; v, z).$$

Applying Hölder's inequality yields

$$|d_M(u; v, z)| \leq \|\tau_M^{\text{sold}}(u)\|_{0,M} \|\kappa_M(P_M \nabla v)\|_{0,4,M} \|\kappa_M(P_M \nabla z)\|_{0,4,M},$$

which, using (24), gives

$$(34) \quad |d_M(u; v, z)| \leq C h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \|\kappa_M(P_M \nabla v)\|_{0,4,M} \|\kappa_M(P_M \nabla z)\|_{0,4,M},$$

thus proving (32). Now it will be shown that

$$(35) \quad |d_M(u; u, z) - d_M(v; v, z)| \leq C h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \zeta_M(u, v) \\ \times (\|\kappa_M(P_M \nabla u)\|_{0,4,M} + \|\kappa_M(P_M \nabla v)\|_{0,4,M}) \|\kappa_M(P_M \nabla z)\|_{0,4,M}.$$

If $|u|_{1,M} = 0$ or $|v|_{1,M} = 0$, then (35) is a particular case of (34). Thus, it suffices to consider the case $|u|_{1,M} \neq 0$, $|v|_{1,M} \neq 0$. Denoting $\xi(x) = |x|x$, one obtains

$$(36) \quad \begin{aligned} d_M(u; u, z) - d_M(v; v, z) &= \frac{\beta h_M^{1+d/2} |\mathbf{b}_M|}{|u|_{1,M}} (\xi(\kappa_M(P_M \nabla u)) - \xi(\kappa_M(P_M \nabla v)), \kappa_M(P_M \nabla z))_M \\ &+ \beta h_M^{1+d/2} |\mathbf{b}_M| \left(\frac{1}{|u|_{1,M}} - \frac{1}{|v|_{1,M}} \right) (\xi(\kappa_M(P_M \nabla v)), \kappa_M(P_M \nabla z))_M. \end{aligned}$$

The integral terms on M possess the same structure as the term $N_M(u, v, z)$ in the proof of Lemma 1 (the second term corresponds to $N_M(0, v, z)$). They are estimated using the same technique, only with a different Hölder inequality. Then, (16) is applied to $\|\kappa_M(P_M \nabla(u-v))\|_{0,M}$ resp. $\|\kappa_M(P_M \nabla v)\|_{0,M}$. Furthermore, the first inequality from (18) is employed. To finish the estimate of the second term in (36), the triangle inequality is used. One obtains

$$\begin{aligned} |d_M(u; u, z) - d_M(v; v, z)| &\leq C h_M^{1+d/2} \|\mathbf{b}\|_{0,\infty,M} \frac{|u-v|_{1,M}}{|u|_{1,M}} \\ &\quad \times (\|\kappa_M(P_M \nabla u)\|_{0,4,M} + \|\kappa_M(P_M \nabla v)\|_{0,4,M}) \|\kappa_M(P_M \nabla z)\|_{0,4,M}. \end{aligned}$$

The same type of inequality follows by interchanging u and v . Then, using the sharper of these two estimates and $\min\{|u|_{1,M}^{-1}, |v|_{1,M}^{-1}\} \leq 2/(|u|_{1,M} + |v|_{1,M})$ gives (35). \square

The properties of the operator T_h , namely its monotonicity and local Lipschitz continuity, follow now by the results of the two previous lemmas and the representation of the LPS norm (25).

Lemma 3. *If $\tilde{\tau}_M$ is defined by (22), then the operator T_h defined in (26) is locally Lipschitz-continuous and strongly monotone, i.e., it satisfies*

$$(37) \quad (T_h w_h - T_h z_h, w_h - z_h) \geq \|w_h - z_h\|_{\text{LPS}}^2 + \frac{1}{7} \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla(w_h - z_h))\|_{0,3,M}^3$$

for all $w_h, z_h \in V_h$. If $\tilde{\tau}_M$ is defined by (23), then the operator T_h is Lipschitz-continuous and it satisfies

$$(38) \quad (T_h z_h, z_h) \geq \frac{\varepsilon}{2} |z_h|_{1,\Omega}^2 - C_0 (\|\tilde{u}_{bh}\|_{1,\Omega}^2 + \|f\|_{0,\Omega}^2)$$

for all $z_h \in V_h$, where $C_0 > 0$ depends on ε , \mathbf{b} , and c , but not on z_h , h , and σ_0

Proof. Let us define the operators $A_h, N_h : V_h \rightarrow V_h$ by

$$\begin{aligned} (A_h z_h, v_h) &= a(z_h, v_h) + s_h(z_h, v_h) \quad \forall z_h, v_h \in V_h, \\ (N_h z_h, v_h) &= d_h(z_h + \tilde{u}_{bh}; z_h + \tilde{u}_{bh}, v_h) \quad \forall z_h, v_h \in V_h. \end{aligned}$$

Then, for any $w_h, z_h \in V_h$, there holds

$$T_h w_h - T_h z_h = A_h(w_h - z_h) + N_h w_h - N_h z_h.$$

The operator A_h is linear on a finite-dimensional space and hence it is Lipschitz continuous. Thus, the (local) Lipschitz-continuity of T_h follows from (28), (33), and the equivalence of norms on finite-dimensional spaces. The strong monotonicity (37) follows from (25) and (27). Finally, let $\tilde{\tau}_M$ be defined by (23). In view of (25), it holds

$$(39) \quad (T_h z_h, z_h) = \|z_h\|_{\text{LPS}}^2 + d_h(z_h + \tilde{u}_{bh}; z_h, z_h) \\ + a(\tilde{u}_{bh}, z_h) + s_h(\tilde{u}_{bh}, z_h) + d_h(z_h + \tilde{u}_{bh}; \tilde{u}_{bh}, z_h) - (f, z_h).$$

Applying (32), (10), (16), (18), (4), and (5), one obtains

$$|d_h(z_h + \tilde{u}_{bh}; \tilde{u}_{bh}, z_h)| \leq C h \|\mathbf{b}\|_{0,\infty,\Omega} |\tilde{u}_{bh}|_{1,\Omega} |z_h|_{1,\Omega}.$$

The same estimate also holds for $s_h(\tilde{u}_{bh}, z_h)$. Using the fact that $d_h(z_h + \tilde{u}_{bh}; z_h, z_h) \geq 0$ and applying the Cauchy–Schwarz inequality to the third and last term on the right-hand side of (39), one derives

$$(T_h z_h, z_h) \geq \varepsilon |z_h|_{1,\Omega}^2 - (\varepsilon + C \|\mathbf{b}\|_{0,\infty,\Omega} + \|c\|_{0,\infty,\Omega}) \|\tilde{u}_{bh}\|_{1,\Omega} \|z_h\|_{1,\Omega} - \|f\|_{0,\Omega} \|z_h\|_{0,\Omega}.$$

Now, employing the Poincaré and Young inequalities, one obtains (38). \square

To prove that the discrete problem (19) has at least one solution, we shall use the following simple consequence of Brouwer’s fixed-point theorem, whose proof can be found in [32, p. 164, Lemma 1.4].

Lemma 4. *Let X be a finite-dimensional Hilbert space with inner product (\cdot, \cdot) and norm $\|\cdot\|$. Let $P : X \rightarrow X$ be a continuous mapping and $K > 0$ a real number such that $(Px, x) > 0$ for any $x \in X$ with $\|x\| = K$. Then there exists $x \in X$ such that $\|x\| \leq K$ and $Px = 0$.*

Collecting the previous results, the main result of this section can be stated now, namely, the well-posedness of the problem (19).

Theorem 5. *If $\tilde{\tau}_M$ is defined by (22) or (23), then the problem (19) has a solution. If $\tilde{\tau}_M$ is defined by (22), the solution of (19) is unique.*

Proof. If $\tilde{\tau}_M$ is defined by (22), then it follows from the strong monotonicity (37) that, for any $z_h \in V_h$,

$$(T_h z_h, z_h) \geq \|z_h\|_{\text{LPS}}^2 + (T_h 0, z_h) \geq \varepsilon |z_h|_{1,\Omega}^2 - \|T_h 0\|_{0,\Omega} \|z_h\|_{0,\Omega}.$$

Thus, using Young's inequality and the equivalence of norms in the space V_h one gets

$$(T_h z_h, z_h) \geq C_1 \|z_h\|_{0,\Omega}^2 - C_2,$$

where C_1, C_2 are positive constants that depend on h and the data of (1), but not on z_h and σ_0 . According to (38), the same inequality holds if $\tilde{\tau}_M$ is defined by (23). Thus, in view of Lemma 4 with any $K > \sqrt{C_2/C_1}$, the operator T_h has a zero and hence the problem (19) has a solution. The uniqueness in the case that $\tilde{\tau}_M$ is defined by (22) follows from the strong monotonicity (37). \square

3.2. Error estimates. For the analysis of the methods introduced in Section 3, we will need an appropriate interpolation operator. An important tool for the construction of such an operator is provided by the following result, whose proof can be found in [25, Lemma 1].

Lemma 6. *Let us suppose the inf-sup condition (9) to be satisfied. Then, there exists an operator $\varrho_h : L^2(\Omega) \rightarrow V_h$ such that, for any $v, w \in L^2(\Omega)$, the estimates*

$$(40) \quad |(v - \varrho_h v, w)| \leq C \sum_{M \in \mathcal{M}_h} \|v\|_{0,M} \|\kappa_M w\|_{0,M},$$

$$(41) \quad |\varrho_h v|_{1,M}^2 + h_M^{-2} \|\varrho_h v\|_{0,M}^2 \leq C \sum_{\substack{M' \in \mathcal{M}_h, \\ M \cap M' \neq \emptyset}} h_{M'}^{-2} \|v\|_{0,M'}^2 \quad \forall M \in \mathcal{M}_h$$

are valid. Consequently, for any $\alpha \in \mathbb{R}$, it holds

$$(42) \quad \sum_{M \in \mathcal{M}_h} h_M^\alpha (|\varrho_h v|_{1,M}^2 + h_M^{-2} \|\varrho_h v\|_{0,M}^2) \leq C \sum_{M \in \mathcal{M}_h} h_M^{\alpha-2} \|v\|_{0,M}^2,$$

where the constant C is independent of v and h but can depend on α .

With the operators i_h and ϱ_h , an operator $r_h \in \mathcal{L}(H^2(\Omega), W_h) \cap \mathcal{L}(H^2(\Omega) \cap H_0^1(\Omega), V_h)$ is defined by

$$(43) \quad r_h v := i_h v + \varrho_h(v - i_h v).$$

To formulate the interpolation properties of r_h , it is convenient to introduce the mesh dependent norm

$$\|v\|_{1,h} = \left(\sum_{M \in \mathcal{M}_h} \{|v|_{1,M}^2 + h_M^{-2} \|v\|_{0,M}^2\} \right)^{1/2}.$$

Then, using (41), the geometrical hypotheses (4) and (5), and the approximation property of i_h (11), one obtains

$$(44) \quad \|v - r_h v\|_{1,h} \leq C \|v - i_h v\|_{1,h} \leq \tilde{C} h^k |v|_{k+1,\Omega} \quad \forall v \in H^{k+1}(\Omega), \quad k = 1, \dots, l,$$

and consequently

$$(45) \quad |v - r_h v|_{1,\Omega} + h^{-1} \|v - r_h v\|_{0,\Omega} \leq C h^k |v|_{k+1,\Omega} \quad \forall v \in H^{k+1}(\Omega), \quad k = 1, \dots, l.$$

The derivation of the error estimates will be based on the following two lemmas. The first one states an interpolation error estimate and the second one states a bound on the nonlinear form d_h .

Lemma 7. *Let $u \in H^{k+1}(\Omega)$ for some $k \in \{1, \dots, l\}$, and let $\eta := u - r_h u$. Then, for any $v_h \in V_h \setminus \{0\}$, the following estimate holds*

$$(46) \quad \|\eta\|_{\text{LPS}} + \frac{a(\eta, v_h) + s_h(\eta, v_h) - s_h(u, v_h)}{\|v_h\|_{\text{LPS}}} \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 \|\sigma\|_{0,\infty,\Omega} + h^2 |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1})^{1/2} h^k |u|_{k+1,\Omega}.$$

Proof. Since, in view of (5), (16), (18), and the definition of τ_M (20)

$$\|v\|_{\text{LPS}} \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 \|\sigma\|_{0,\infty,\Omega})^{1/2} \|v\|_{1,h} \quad \forall v \in H^1(\Omega),$$

it follows from (44) that

$$\|\eta\|_{\text{LPS}} \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 \|\sigma\|_{0,\infty,\Omega})^{1/2} h^k |u|_{k+1,\Omega}.$$

Next, for any $v_h \in V_h \setminus \{0\}$, integration by parts gives

$$(\mathbf{b} \cdot \nabla \eta, v_h) = -(\eta, \mathbf{b} \cdot \nabla v_h) - ((\nabla \cdot \mathbf{b}) \eta, v_h).$$

Thus, applying the Cauchy–Schwarz inequality and (45), it follows that

$$a(\eta, v_h) + s_h(\eta, v_h) \leq \left(\|\eta\|_{\text{LPS}} + C |\mathbf{b}|_{1,\infty,\Omega} \sigma_0^{-1/2} h^{k+1} |u|_{k+1,\Omega} \right) \|v_h\|_{\text{LPS}} - (\eta, \mathbf{b} \cdot \nabla v_h).$$

The use of (40), the approximation property of i_h (11), (4), and (5) leads to

$$\begin{aligned} (\eta, \mathbf{b} \cdot \nabla v_h) &\leq C \sum_{M \in \mathcal{M}_h} \|u - i_h u\|_{0,M} \|\kappa_M(\mathbf{b} \cdot \nabla v_h)\|_{0,M} \\ &\leq C h^k |u|_{k+1,\Omega} \left(\sum_{M \in \mathcal{M}_h} h_M^2 \|\kappa_M(\mathbf{b} \cdot \nabla v_h)\|_{0,M}^2 \right)^{1/2}. \end{aligned}$$

Applying (16), (18), (20), and the inverse inequality (8), one derives

$$\begin{aligned} \|\kappa_M(\mathbf{b} \cdot \nabla v_h)\|_{0,M} &\leq \|\kappa_M((\mathbf{b} - \mathbf{b}_M) \cdot \nabla v_h)\|_{0,M} + \|\kappa_M(\mathbf{b}_M \cdot \nabla v_h)\|_{0,M} \\ &\leq C |\mathbf{b}|_{1,\infty,M} \|v_h\|_{0,M} + \tau_0^{-1/2} (\varepsilon + h_M \|\mathbf{b}\|_{0,\infty,M})^{1/2} h_M^{-1} \tau_M^{1/2} \|\kappa_M(\mathbf{b}_M \cdot \nabla v_h)\|_{0,M}, \end{aligned}$$

which leads to the estimate

$$(\eta, \mathbf{b} \cdot \nabla v_h) \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1})^{1/2} h^k |u|_{k+1,\Omega} \|v_h\|_{\text{LPS}}.$$

Finally, using (17), (18), (20), and the geometrical hypotheses (4) and (5), one obtains

$$s_h(u, u) \leq \sum_{M \in \mathcal{M}_h} \tau_M |\mathbf{b}_M|^2 \|\kappa_M \nabla u\|_{0,M}^2 \leq C \|\mathbf{b}\|_{0,\infty,\Omega} h^{2k+1} |u|_{k+1,\Omega}^2,$$

and hence

$$s_h(u, v_h) \leq \sqrt{s_h(u, u)} \sqrt{s_h(v_h, v_h)} \leq C \|\mathbf{b}\|_{0,\infty,\Omega}^{1/2} h^{k+1/2} |u|_{k+1,\Omega} \|v_h\|_{\text{LPS}},$$

which completes the proof. \square

Lemma 8. *For any $w_h \in W_h$ and $u, v \in H^{k+1}(\Omega)$ with $k \in \{1, \dots, l\}$, it holds*

$$(47) \quad d_h(w_h; r_h u, r_h v) \leq C h^{2k-d/2} \left(\max_{M \in \mathcal{M}_h} \|\tau_M^{\text{sold}}(w_h)\|_{0,M} \right) |u|_{k+1,\Omega} |v|_{k+1,\Omega}.$$

Proof. The application of Hölder's inequality and (10) leads to

$$(48) \quad \begin{aligned} d_h(w_h; r_h u, r_h v) &\leq \sum_{M \in \mathcal{M}_h} \|\tau_M^{\text{sold}}(w_h)\|_{0,M} \|\kappa_M(P_M \nabla(r_h u))\|_{0,4,M} \|\kappa_M(P_M \nabla(r_h v))\|_{0,4,M} \\ &\leq C \sum_{M \in \mathcal{M}_h} \|\tau_M^{\text{sold}}(w_h)\|_{0,M} h_M^{-d/2} \|\kappa_M(P_M \nabla(r_h u))\|_{0,M} \|\kappa_M(P_M \nabla(r_h v))\|_{0,M} \\ &\leq C \left(\max_{M \in \mathcal{M}_h} \|\tau_M^{\text{sold}}(w_h)\|_{0,M} \right) \left(\sum_{M \in \mathcal{M}_h} h_M^{-d/2} \|\kappa_M(P_M \nabla(r_h u))\|_{0,M}^2 \right)^{1/2} \\ &\quad \times \left(\sum_{M \in \mathcal{M}_h} h_M^{-d/2} \|\kappa_M(P_M \nabla(r_h v))\|_{0,M}^2 \right)^{1/2}. \end{aligned}$$

Let us estimate the term with u ; the term with v can be treated analogously. Using (16) and (17), for $u \in H^{k+1}(\Omega)$ with $k \in \{1, \dots, l\}$ there holds

$$(49) \quad \begin{aligned} \|\kappa_M(P_M \nabla(r_h u))\|_{0,M} &\leq \|\kappa_M(P_M \nabla u)\|_{0,M} + \|\kappa_M(P_M \nabla(u - r_h u))\|_{0,M} \\ &\leq C h_M^k |u|_{k+1,M} + C |u - r_h u|_{1,M}. \end{aligned}$$

According to (42), one has for any $\alpha \in \mathbb{R}$

$$\begin{aligned} \sum_{M \in \mathcal{M}_h} h_M^\alpha |u - r_h u|_{1,M}^2 &\leq 2 \sum_{M \in \mathcal{M}_h} h_M^\alpha |u - i_h u|_{1,M}^2 + 2 \sum_{M \in \mathcal{M}_h} h_M^\alpha |\varrho_h(u - i_h u)|_{1,M}^2 \\ &\leq C \sum_{M \in \mathcal{M}_h} h_M^\alpha (|u - i_h u|_{1,M}^2 + h_M^{-2} \|u - i_h u\|_{0,M}^2), \end{aligned}$$

and hence it follows from the approximation property of i_h (11), (4), and (5) that, for $\alpha \geq -2$,

$$(50) \quad \sum_{M \in \mathcal{M}_h} h_M^\alpha \|\kappa_M(P_M \nabla(r_h u))\|_{0,M}^2 \leq C h^{2k+\alpha} |u|_{k+1,\Omega}^2.$$

Inserting (50) with $\alpha = -d/2$ into (48), the statement of the lemma is proved. \square

We are now in position to prove the first error estimate. The following theorem states the error estimate in the case $\tilde{\tau}_M$ is given by (22).

Theorem 9. *Let $\tilde{\tau}_M$ be defined by (22). Let the weak solution of (1) satisfy $u \in H^{k+1}(\Omega)$ for some $k \in \{1, \dots, l\}$. Let $\tilde{u}_b \in H^2(\Omega)$ be an extension of u_b and let $\tilde{u}_{bh} = i_h \tilde{u}_b$. Then the solution u_h of the local projection discretization (19) satisfies the error estimate*

$$\begin{aligned} & \|u - u_h\|_{\text{LPS}} + \left(\sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla(u - u_h))\|_{0,3,M}^3 \right)^{1/2} \\ & \leq C \left\{ \varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} (1 + \beta h^{k-d/2} |u|_{k+1,\Omega}) + h^2 (\|\sigma\|_{0,\infty,\Omega} + |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1}) \right\}^{1/2} h^k |u|_{k+1,\Omega}. \end{aligned}$$

If $u \in W^{k+1,\infty}(\Omega)$ with $k \in \{1, \dots, l\}$, then

$$\begin{aligned} & \|u - u_h\|_{\text{LPS}} + \left(\sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla(u - u_h))\|_{0,3,M}^3 \right)^{1/2} \\ & \leq C \left\{ \varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} (1 + \beta h^k |u|_{k+1,\infty,\Omega}) + h^2 (\|\sigma\|_{0,\infty,\Omega} + |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1}) \right\}^{1/2} h^k |u|_{k+1,\Omega}. \end{aligned}$$

Proof. The error $u - u_h$ is split into the interpolation error $\eta := u - r_h u$ and the discrete error $e_h := u_h - r_h u$. Then $e_h \in V_h$ and also $r_h u - \tilde{u}_{bh} \in V_h$. From the monotonicity (37) it follows with the discrete problem (19) and the continuous problem (14) that

$$\begin{aligned} & \|e_h\|_{\text{LPS}}^2 + \frac{1}{7} \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla e_h)\|_{0,3,M}^3 \leq (T_h(u_h - \tilde{u}_{bh}) - T_h(r_h u - \tilde{u}_{bh}), e_h) \\ & = a(u_h, e_h) + s_h(u_h, e_h) + d_h(u_h; u_h, e_h) - (T_h(r_h u - \tilde{u}_{bh}), e_h) \\ & = (f, e_h) - (T_h(r_h u - \tilde{u}_{bh}), e_h) \\ & = a(u, e_h) - a(r_h u, e_h) - s_h(r_h u, e_h) - d_h(r_h u; r_h u, e_h) \\ & = a(\eta, e_h) + s_h(\eta, e_h) - s_h(u, e_h) - d_h(r_h u; r_h u, e_h). \end{aligned}$$

The first three terms on the right-hand side can be estimated using (46). To bound the nonlinear term, Hölder's and Young's inequalities are applied to conclude

$$\begin{aligned} (51) \quad d_h(r_h u; r_h u, e_h) & \leq \{d_h(r_h u; r_h u, r_h u)\}^{\frac{2}{3}} \{d_h(e_h; e_h, e_h)\}^{\frac{1}{3}} \\ & \leq 2 d_h(r_h u; r_h u, r_h u) + \frac{3}{70} d_h(e_h; e_h, e_h). \end{aligned}$$

Then (47), (49), the bound of h_M (5), (18), and (45) yield

$$(52) \quad d_h(r_h u; r_h u, r_h u) \leq C \beta \|\mathbf{b}\|_{0,\infty,\Omega} h^{3k+1-d/2} |u|_{k+1,\Omega}^3.$$

Therefore,

$$(53) \quad \|e_h\|_{\text{LPS}}^2 + \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla e_h)\|_{0,3,M}^3 \\ \leq C \left\{ \varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} (1 + \beta h^{k-d/2} |u|_{k+1,\Omega}) + h^2 \|\sigma\|_{0,\infty,\Omega} + h^2 \|\mathbf{b}\|_{1,\infty,\Omega}^2 \sigma_0^{-1} \right\} h^{2k} |u|_{k+1,\Omega}^2.$$

Next, to estimate the interpolation error, for any $p \in [1, 6]$, it follows from the commutation property of κ_M and P_M , the estimate of the $L^p(M)$ norm by the $L^2(M)$ norm (10), (15), and (13) that

$$(54) \quad \|\kappa_M(P_M \nabla \eta)\|_{0,p,M} \leq \|\nabla \eta - \pi_M \nabla \eta\|_{0,p,M} \\ \leq \|\nabla(u - i_h u)\|_{0,p,M} + \|\nabla(i_h u - r_h u) - \pi_M \nabla \eta\|_{0,p,M} \\ \leq |u - i_h u|_{1,p,M} + C h_M^{\frac{d}{p} - \frac{d}{2}} \|\nabla(i_h u - r_h u) - \pi_M \nabla \eta\|_{0,M} \\ \leq |u - i_h u|_{1,p,M} + \tilde{C} h_M^{\frac{d}{p} - \frac{d}{2}} (|\varrho_h(u - i_h u)|_{1,M} + |u - i_h u|_{1,M}) \\ \leq \bar{C} h_M^{k + \frac{d}{p} - \frac{d}{2}} |u|_{k+1,M} + \tilde{C} h_M^{\frac{d}{p} - \frac{d}{2}} |\varrho_h(u - i_h u)|_{1,M}.$$

Then, applying (54), (22), (5), (18), (41), (11), (4), and (6), one derives

$$(55) \quad \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla \eta)\|_{0,3,M}^3 \leq C \beta h \|\mathbf{b}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} h_M^{3k-d/2} |u|_{k+1,M}^3.$$

Thus, combining (53), (55), and (46), the first estimate of the theorem follows.

If $u \in W^{k+1,\infty}(\Omega)$ with $k \in \{1, \dots, l\}$, then local norms of Sobolev spaces with $p = 2$ can be estimated with norms of Sobolev spaces with $p = \infty$, thereby gaining powers of h from the smallness of the local domain: $|u|_{k+1,M} \leq C h_M^{d/2} |u|_{k+1,\infty,M}$ for any $M \in \mathcal{M}_h$. Hence, it follows from (55) and the geometrical hypotheses (4) and (5) that

$$\sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla \eta)\|_{0,3,M}^3 \leq C \beta \|\mathbf{b}\|_{0,\infty,\Omega} h^{3k+1} |u|_{k+1,\infty,\Omega} |u|_{k+1,\Omega}^2.$$

Furthermore, using (41), (11), and (4), one gets

$$|u - r_h u|_{1,M} \leq C \sum_{\substack{M' \in \mathcal{M}_h, \\ M \cap M' \neq \emptyset}} h_{M'}^k |u|_{k+1,M'} \leq \tilde{C} h^{k+d/2} |u|_{k+1,\infty,\Omega} \quad \forall M \in \mathcal{M}_h.$$

Therefore, according to (47) and (49),

$$(56) \quad d_h(r_h u; r_h u, r_h u) \leq C \beta \|\mathbf{b}\|_{0,\infty,\Omega} h^{3k+1} |u|_{k+1,\infty,\Omega} |u|_{k+1,\Omega}^2,$$

which implies the second estimate of the theorem. \square

Remark. Theorem 9 implies, in particular, the following convergence estimates in the convection-dominated case $\varepsilon < h$: If $u \in H^2(\Omega)$, then

$$\|u - u_h\|_{\text{LPS}} \leq C_0 h^{2-d/4} (h^{(d-2)/4} + |u|_{2,\Omega}^{1/2}) |u|_{2,\Omega},$$

where C_0 depends on the data of the problem. If $u \in W^{2,\infty}(\Omega)$, then

$$\|u - u_h\|_{\text{LPS}} \leq C_0 h^{3/2} (1 + h^{1/2} |u|_{2,\infty,\Omega}^{1/2}) |u|_{2,\Omega}.$$

If $u \in H^{k+1}(\Omega)$ with $k \in \{2, \dots, l\}$, then

$$\|u - u_h\|_{\text{LPS}} \leq C_0 h^{k+1/2} (1 + h^{(2k-d)/4} |u|_{k+1,\Omega}^{1/2}) |u|_{k+1,\Omega}.$$

Remark. A situation of practical interest is that the convective field \mathbf{b} arises from a finite element approximation of the Navier–Stokes equations. In this case, a necessary condition for a uniform convergence of $\|\mathbf{b}\|_{1,\infty,\Omega}$ with respect to h is that the exact velocity is sufficiently regular. This condition might not be fulfilled, e.g., if the domain possesses re-entrant corners, and therefore estimates involving weaker norms of \mathbf{b} are also of interest. Changing the arguments in the proof of Lemma 7 slightly, one obtains, e.g., the following result

$$\begin{aligned} & \|u - u_h\|_{\text{LPS}} + \left(\sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M \nabla(u - u_h))\|_{0,3,M}^3 \right)^{1/2} \\ (57) \quad & \leq C \left\{ \varepsilon + \|\mathbf{b}\|_{0,\infty,\Omega}^2 \sigma_0^{-1} + h \|\mathbf{b}\|_{0,\infty,\Omega} (1 + \beta h^{k-d/2} |u|_{k+1,\Omega}) \right. \\ & \left. + h^{2-\frac{d}{2}} \max_{M \in \mathcal{M}_h} \|\nabla \cdot \mathbf{b}\|_{0,4,M}^2 \sigma_0^{-1} + h^2 \|\sigma\|_{0,\infty,\Omega} \right\}^{1/2} h^k |u|_{k+1,\Omega}. \end{aligned}$$

If the norms of \mathbf{b} in (57) are still too strong, one can use the discrete character of a computed convection field \mathbf{b} and apply inverse inequalities to derive estimates involving the weaker norms $\|\mathbf{b}\|_{1,\Omega}$ and $\|\nabla \cdot \mathbf{b}\|_{0,\Omega}$. However, the relaxation of the regularity assumption on \mathbf{b} in the error bounds is accompanied with a reduction of the order of convergence, e.g., the order of convergence of (57) is reduced by 1/2 compared with the orders given in the previous remark.

Remark. The right-hand sides of the estimates in Theorem 9 can be stated in terms of local (semi)norms of the data and of the solution on macro-elements multiplied by diameters of the macro-elements. However, due to the use of the interpolation operator r_h , such estimates are more complicated than usually. For example, a counterpart of (52) using local quantities

has the form

$$d_h(r_h u; r_h u, r_h u) \leq C \beta \sum_{M \in \mathcal{M}_h} \|\mathbf{b}\|_{0,\infty,M} h_M^{1-d/2} \left(\sum_{\substack{M' \in \mathcal{M}_h, \\ M \cap M' \neq \emptyset}} h_{M'}^{2k} |u|_{k+1,M'}^2 \right)^{3/2}.$$

Therefore, for clarity, we decided to state the estimates in terms of global quantities.

We end this section by presenting the error estimate in the case $\tilde{\tau}_M$ is defined by (23).

Theorem 10. *Let $\tilde{\tau}_M$ be defined by (23). Let the weak solution of (1) satisfy $u \in H^{k+1}(\Omega)$ for some $k \in \{1, \dots, l\}$. Let $\tilde{u}_b \in H^2(\Omega)$ be an extension of u_b and let $\tilde{u}_{bh} = i_h \tilde{u}_b$. Then the solution u_h of the local projection discretization (19) satisfies the error estimate*

$$\begin{aligned} \|u - u_h\|_{\text{LPS}} + (d_h(u_h; u - u_h, u - u_h))^{1/2} \\ \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 \|\sigma\|_{0,\infty,\Omega} + h^2 |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1})^{1/2} h^k |u|_{k+1,\Omega}. \end{aligned}$$

Proof. Set again $\eta := u - r_h u$ and $e_h := u_h - r_h u$. From (19) and (14), it follows that

$$\begin{aligned} a(e_h, e_h) + s_h(e_h, e_h) + d_h(u_h; u_h, e_h) \\ = a(u_h, e_h) + s_h(u_h, e_h) + d_h(u_h; u_h, e_h) - a(r_h u, e_h) - s_h(r_h u, e_h) \\ = a(\eta, e_h) + s_h(\eta, e_h) - s_h(u, e_h). \end{aligned}$$

Thus, in view of the representation of the LPS norm (25), one gets

$$\|e_h\|_{\text{LPS}}^2 + d_h(u_h; e_h, e_h) = a(\eta, e_h) + s_h(\eta, e_h) - s_h(u, e_h) - d_h(u_h; r_h u, e_h).$$

The first three terms on the right-hand side can be estimated using (46). To bound the nonlinear term, Hölder's and Young's inequalities are again applied

$$(58) \quad d_h(u_h; r_h u, e_h) \leq \sqrt{d_h(u_h; r_h u, r_h u)} \sqrt{d_h(u_h; e_h, e_h)} \leq d_h(u_h; r_h u, r_h u) + \frac{1}{4} d_h(u_h; e_h, e_h).$$

Using (47), (24), and (5), one obtains

$$(59) \quad d_h(u_h; r_h u, r_h u) \leq C \|\mathbf{b}\|_{0,\infty,\Omega} h^{2k+1} |u|_{k+1,\Omega}^2.$$

Therefore,

$$\|e_h\|_{\text{LPS}}^2 + d_h(u_h; e_h, e_h) \leq C (\varepsilon + h \|\mathbf{b}\|_{0,\infty,\Omega} + h^2 \|\sigma\|_{0,\infty,\Omega} + h^2 |\mathbf{b}|_{1,\infty,\Omega}^2 \sigma_0^{-1}) h^{2k} |u|_{k+1,\Omega}^2.$$

Note that an application of the triangle inequality gives

$$(60) \quad d_h(u_h; u - u_h, u - u_h) \leq 2 d_h(u_h; \eta, \eta) + 2 d_h(u_h; e_h, e_h).$$

It follows from Hölder's inequality, (24), (54), (42) with $\alpha = 0$, (11), (4), and (5), that

$$(61) \quad d_h(u_h; \eta, \eta) \leq \sum_{M \in \mathcal{M}_h} \|\tau_M^{\text{sold}}(u_h)\|_{0,M} \|\kappa_M(P_M \nabla \eta)\|_{0,4,M}^2 \leq C \|\mathbf{b}\|_{0,\infty,\Omega} h^{2k+1} |u|_{k+1,\Omega}^2.$$

Finally, using the triangle inequality and the estimate (46), the statement of the theorem follows. \square

Remark. Theorems 9 and 10 prove the convergence of the method in the LPS norm plus an extra term involving the crosswind derivative of the error. Hence, these estimates give, essentially, an extra control of the whole gradient of the error.

4. THE TIME-DEPENDENT PROBLEM

We now move on to the study of the time-dependent problem (3). A weak form of problem (3) reads as follows: Find $u \in L^2(0, T; H^1(\Omega)) \cap H^1(0, T; L^2(\Omega))$ such that $u = u_b$ on $[0, T] \times \partial\Omega$, $u(0, \cdot) = u_0$ and

$$(62) \quad (u_t, v) + a(u, v) = (f, v) \quad \forall v \in H_0^1(\Omega), \quad \text{for almost every } t \in (0, T].$$

To avoid technicalities in the analysis, it is assumed that the boundary condition does not depend on time, $u_b(t, \cdot) = u_b$. The initial condition u_0 is assumed to satisfy $u_0|_{\partial\Omega} = u_b$ and it is approximated by a function $u_h^0 \in W_h$ such that $u_h^0 - \tilde{u}_{bh} \in V_h$.

To perform the discretization of the time derivative, the time interval $[0, T]$ is divided into N_T equidistant strips of length $\delta t = T/N_T$. The constant time step is used only for simplicity of presentation; for variable time steps the same techniques can be applied leading to essentially the same results. The nodes are denoted by $t^n = n \delta t$ for $n = 0, 1, \dots, N_T$ and the abbreviations $u^n := u(t^n, \cdot)$, $f^n := f(t^n, \cdot)$, etc. are used. Since this section studies the LPS method with nonlinear crosswind diffusion in combination with a one-step θ -scheme as temporal discretization, from now on, the superscript $n + \theta$ denotes for all functions which are defined in $[0, T]$ the values at time $t^{n+\theta} := \theta t^{n+1} + (1 - \theta) t^n$ with any $n \in \{0, \dots, N_T - 1\}$ and $\theta \in [0, 1]$, e.g. $\mathbf{b}^{n+\theta} = \mathbf{b}(t^{n+\theta}, \cdot)$. For functions, which are defined only at the discrete times t^n and t^{n+1} , it denotes the linear interpolation, e.g. $u_h^{n+\theta} = \theta u_h^{n+1} + (1 - \theta) u_h^n$. Finally, it is convenient to introduce the interpolation operator $\tilde{r}_h^{n+\theta}$ satisfying

$$(63) \quad \tilde{r}_h^{n+\theta} u = \theta r_h u^{n+1} + (1 - \theta) r_h u^n$$

with r_h from (43). Thus, writing α instead of $n + \theta$, functions u^α , u_h^α , $\tilde{r}_h^\alpha u$, etc. are defined for any $\alpha \in [0, N_T]$.

Then, given $\theta \in (0, 1]$, the fully discrete problem reads as follows: For $n = 0, 1, \dots, N_T - 1$, find $u_h^{n+1} \in W_h$ such that $u_h^{n+1} - \tilde{u}_{bh} \in V_h$ and

$$(64) \quad \left(\frac{u_h^{n+1} - u_h^n}{\delta t}, v_h \right) + a^{n+\theta}(u_h^{n+\theta}, v_h) + s_h^{n+\theta}(u_h^{n+\theta}, v_h) + d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, v_h) \\ = (f^{n+\theta}, v_h) \quad \forall v_h \in V_h.$$

For $\theta = 1/2$, the Crank–Nicolson scheme is recovered and for $\theta = 1$, the implicit Euler scheme is obtained.

Remark. To simplify the notation, we will not explicitly indicate at which time instant the functions \mathbf{b} and σ in the definition of the norm $\|\cdot\|_{\text{LPS}}$ are evaluated. This will be implicitly determined from the context or by the argument of the norm. Thus, if we write, e.g., $\|u_h^{n+\theta}\|_{\text{LPS}}$, the norm $\|\cdot\|_{\text{LPS}}$ is defined using $\mathbf{b}^{n+\theta}$ and $\sigma^{n+\theta}$.

4.1. Well-posedness and stability. The well-posedness of (64) can be traced back to the well-posedness of the LPS scheme with crosswind diffusion for the steady-state problem. The discretization of the temporal derivative can be written in the form

$$\left(\frac{u_h^{n+1} - u_h^n}{\delta t}, v_h \right) = \frac{1}{\theta} \left(\frac{u_h^{n+\theta} - u_h^n}{\delta t}, v_h \right).$$

The first part of this term has the form of a reaction term for $u_h^{n+\theta}$. Thus, given u_h^n , the equation at the discrete time t^{n+1} is an equation for $u_h^{n+\theta}$ which has the same form as (19) with the data of the problem at $t^{n+\theta}$ and with a reaction coefficient which has a contribution from the temporal derivative. Thus, defining the operator $\tilde{T}_h^{n+\theta} : V_h \rightarrow V_h$ by

$$(\tilde{T}_h^{n+\theta} z_h, v_h) = (T_h^{n+\theta} z_h, v_h) + \frac{1}{\theta \delta t} (z_h + \tilde{u}_{bh}, v_h) - \frac{1}{\theta \delta t} (u_h^n, v_h) \quad \forall z_h, v_h \in V_h,$$

it follows that $\tilde{T}_h^{n+\theta}(u_h^{n+\theta} - \tilde{u}_{bh}) = 0$. Therefore, the existence and uniqueness of a solution $u_h^{n+\theta}$ can be proved in the same way as in the steady-state case, see Section 3.1. This fact is stated in the next result.

Corollary 11. *Let $n \in \{0, 1, \dots, N_T - 1\}$ and $u_h^n \in W_h$ with $u_h^n|_{\partial\Omega} = \tilde{u}_{bh}$ be given. If $\tilde{\tau}_M$ is defined by (22) or (23), then the problem (64) possesses a solution u_h^{n+1} . In the case that $\tilde{\tau}_M$ is defined by (22), the solution of (64) is unique. Furthermore, there is a constant $C > 0$ such that the solution of the scheme (64) with $\tilde{\tau}_M$ given by (23) is unique if $\delta t \|\mathbf{b}^{n+\theta}\|_{0,\infty,M} \leq C h_M$ for any $M \in \mathcal{M}_h$.*

Proof. The only point remaining to prove is the uniqueness in the case $\tilde{\tau}_M$ is given by (23). For this, let $v_h, w_h \in W_h$ and $z_h := v_h - w_h$. Then, applying (33), the estimate of the $L^p(M)$

norm by the $L^2(M)$ norm (10), (16), $\|P_M^{n+\theta}\|_2 = 1$, and the inverse inequality (8), one arrives at

$$|d_h^{n+\theta}(v_h; v_h, z_h) - d_h^{n+\theta}(w_h; w_h, z_h)| \leq C \sum_{M \in \mathcal{M}_h} h_M^{-1} \|\mathbf{b}^{n+\theta}\|_{0,\infty,M} \|z_h\|_{0,M}^2.$$

Thus, if $v_h, w_h \in V_h$, one obtains

$$(\tilde{T}_h^{n+\theta} v_h - \tilde{T}_h^{n+\theta} w_h, z_h) \geq \sum_{M \in \mathcal{M}_h} \left(\frac{\tilde{C}}{\theta \delta t} - \frac{C \|\mathbf{b}^{n+\theta}\|_{0,\infty,M}}{h_M} \right) \|z_h\|_{0,M}^2 + \|z_h\|_{\text{LPS}}^2.$$

Consequently, for δt small enough, the operator $\tilde{T}_h^{n+\theta}$ is strongly monotone and hence the solution to the discrete problem (64) is unique. \square

The next result states the stability of the method.

Lemma 12. *Let $\theta \in [1/2, 1]$ be given. Let $\tilde{u}_h^\alpha := u_h^\alpha - \tilde{u}_{bh}$ for any $\alpha \in [0, N_T]$. Then any solution of (64) satisfies the following stability estimate for all $N = 1, 2, \dots, N_T$:*

$$(65) \quad \|\tilde{u}_h^N\|_{0,\Omega}^2 + (2\theta - 1) \sum_{n=0}^{N-1} \|\tilde{u}_h^{n+1} - \tilde{u}_h^n\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2 \\ + \delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(\bar{u}_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) \leq \|\tilde{u}_h^0\|_{0,\Omega}^2 + C \delta t \sum_{n=0}^{N-1} \left\{ \sigma_0^{-1} \|f^{n+\theta}\|_{0,\Omega}^2 \right. \\ \left. + \left[\varepsilon + \sigma_0^{-1} (\|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega}^2 + \|c^{n+\theta}\|_{0,\infty,\Omega}^2) + h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \right] \|\tilde{u}_{bh}\|_{1,\Omega}^2 + \mu_h \right\},$$

where

$$(66) \quad \bar{u}_h^{n+\theta} = \tilde{u}_h^{n+\theta}, \quad \mu_h = \beta h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} |\tilde{u}_{bh}|_{1,3,\Omega}^3 \quad \text{if } \tilde{\tau}_M \text{ is given by (22),}$$

$$(67) \quad \bar{u}_h^{n+\theta} = u_h^{n+\theta}, \quad \mu_h = 0 \quad \text{if } \tilde{\tau}_M \text{ is given by (23).}$$

Proof. The proof starts in the usual way by setting $v_h = \tilde{u}_h^{n+\theta} \in V_h$ in (64) and using that $u_h^{n+1} - u_h^n = \tilde{u}_h^{n+1} - \tilde{u}_h^n$, which leads to

$$(68) \quad (\tilde{u}_h^{n+1} - \tilde{u}_h^n, \tilde{u}_h^{n+\theta}) + \delta t \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2 + \delta t d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, \tilde{u}_h^{n+\theta}) \\ = \delta t (f^{n+\theta}, \tilde{u}_h^{n+\theta}) - \delta t a^{n+\theta}(\tilde{u}_{bh}, \tilde{u}_h^{n+\theta}) - \delta t s_h^{n+\theta}(\tilde{u}_{bh}, \tilde{u}_h^{n+\theta}).$$

A straightforward computation gives

$$(69) \quad (\tilde{u}_h^{n+1} - \tilde{u}_h^n, \tilde{u}_h^{n+\theta}) = \frac{1}{2} (\|\tilde{u}_h^{n+1}\|_{0,\Omega}^2 - \|\tilde{u}_h^n\|_{0,\Omega}^2) + \frac{2\theta - 1}{2} \|\tilde{u}_h^{n+1} - \tilde{u}_h^n\|_{0,\Omega}^2.$$

Next, the application of the Cauchy–Schwarz inequality, the Young inequality, (16), (18), the definition of τ_M (20), and the geometrical hypotheses (4) and (5) yield

$$\begin{aligned} (f^{n+\theta}, \tilde{u}_h^{n+\theta}) &\leq \frac{1}{\sigma_0} \|f^{n+\theta}\|_{0,\Omega}^2 + \frac{1}{4} \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2, \\ a^{n+\theta}(\tilde{u}_{bh}, \tilde{u}_h^{n+\theta}) &\leq 6 [\varepsilon + \sigma_0^{-1} (\|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega}^2 + \|c^{n+\theta}\|_{0,\infty,\Omega}^2)] \|\tilde{u}_{bh}\|_{1,\Omega}^2 + \frac{1}{8} \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2, \\ s_h^{n+\theta}(\tilde{u}_{bh}, \tilde{u}_h^{n+\theta}) &\leq C h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} |\tilde{u}_{bh}|_{1,\Omega}^2 + \frac{1}{8} \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2. \end{aligned}$$

If $\tilde{\tau}_M$ is given by (22), then, from (27) and an analog of (51), one obtains

$$\begin{aligned} d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, \tilde{u}_h^{n+\theta}) &\geq \frac{1}{7} d_h^{n+\theta}(\tilde{u}_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) + d_h^{n+\theta}(\tilde{u}_{bh}; \tilde{u}_{bh}, \tilde{u}_h^{n+\theta}) \\ &\geq \frac{1}{10} d_h^{n+\theta}(\tilde{u}_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) - 2 d_h^{n+\theta}(\tilde{u}_{bh}; \tilde{u}_{bh}, \tilde{u}_{bh}). \end{aligned}$$

Furthermore, the use of (10), (16), (18), $\|P_M^{n+\theta}\|_2 = 1$, (4), and (5) leads to

$$d_h^{n+\theta}(\tilde{u}_{bh}; \tilde{u}_{bh}, \tilde{u}_{bh}) \leq C \beta \sum_{M \in \mathcal{M}_h} h_M^{1-d/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,M} |\tilde{u}_{bh}|_{1,M}^3 \leq \tilde{C} \beta h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} |\tilde{u}_{bh}|_{1,3,\Omega}^3.$$

If $\tilde{\tau}_M$ is given by (23), then, using an inequality like (58), one gets

$$\begin{aligned} d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, \tilde{u}_h^{n+\theta}) &= d_h^{n+\theta}(u_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) + d_h^{n+\theta}(u_h^{n+\theta}; \tilde{u}_{bh}, \tilde{u}_h^{n+\theta}) \\ &\geq \frac{1}{2} d_h^{n+\theta}(u_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) - \frac{1}{2} d_h^{n+\theta}(u_h^{n+\theta}; \tilde{u}_{bh}, \tilde{u}_{bh}). \end{aligned}$$

Applying the Hölder inequality, (24), the estimate of the $L^p(M)$ norm by the $L^2(M)$ norm (10), (16), $\|P_M^{n+\theta}\|_2 = 1$, (4), and (5), one deduces that

$$\begin{aligned} d_h^{n+\theta}(u_h^{n+\theta}; \tilde{u}_{bh}, \tilde{u}_{bh}) &\leq C \sum_{M \in \mathcal{M}_h} h_M^{1+d/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,M} \|\kappa_M(P_M^{n+\theta} \nabla \tilde{u}_{bh})\|_{0,4,M}^2 \\ &\leq \tilde{C} h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} |\tilde{u}_{bh}|_{1,\Omega}^2. \end{aligned}$$

Now, inserting the above relations into (68) and using the notation (66) and (67), one obtains

$$\begin{aligned} &\frac{1}{2} (\|\tilde{u}_h^{n+1}\|_{0,\Omega}^2 - \|\tilde{u}_h^n\|_{0,\Omega}^2) + \frac{2\theta-1}{2} \|\tilde{u}_h^{n+1} - \tilde{u}_h^n\|_{0,\Omega}^2 + \frac{\delta t}{2} \|\tilde{u}_h^{n+\theta}\|_{\text{LPS}}^2 + \frac{\delta t}{6} d_h^{n+\theta}(\bar{u}_h^{n+\theta}; \tilde{u}_h^{n+\theta}, \tilde{u}_h^{n+\theta}) \\ &\leq \delta t \sigma_0^{-1} \|f^{n+\theta}\|_{0,\Omega}^2 + C \delta t \{ \varepsilon + \sigma_0^{-1} (\|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega}^2 + \|c^{n+\theta}\|_{0,\infty,\Omega}^2) + h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \} \|\tilde{u}_{bh}\|_{1,\Omega}^2 \\ &\quad + C \delta t \mu_h, \end{aligned}$$

and (65) follows by summing up from $n = 0$ to $N - 1$. \square

Remark. The inequality (65) is a proper stability result provided that $\|u_h^0\|_{0,\Omega}$, $\|\tilde{u}_{bh}\|_{1,\Omega}$ and, if $\tilde{\tau}_M$ is given by (22), also $|\tilde{u}_{bh}|_{1,3,\Omega}$ are bounded when $h \rightarrow 0$. One may set $u_h^0 = I_h u_0$ and $\tilde{u}_{bh} = I_h \tilde{u}_b$, where $I_h : H^1(\Omega) \rightarrow W_h$ is the Scott–Zhang interpolation operator (cf., e.g., [12])

and $\tilde{u}_b \in H^1(\Omega)$ is an extension of u_b . Then $\|u_h^0\|_{0,\Omega} \leq C \|u_0\|_{1,\Omega}$ and $\|\tilde{u}_{bh}\|_{1,\Omega} \leq C \|\tilde{u}_b\|_{1,\Omega}$. If $\tilde{u}_b \in W^{1,3}(\Omega)$ (requiring the stronger assumption $u_b \in W^{2/3,3}(\partial\Omega)$), then also $|\tilde{u}_{bh}|_{1,3,\Omega} \leq C \|\tilde{u}_b\|_{1,3,\Omega}$. It is important that I_h preserves homogeneous boundary conditions since one has to assure that u_h^0 and \tilde{u}_{bh} coincide on the boundary of Ω . If $u_0 \in H^2(\Omega)$ and $u_b \in H^{3/2}(\partial\Omega)$, which are the minimal regularity assumptions for deriving the error estimates in the next section, one may use the operator i_h from Section 2 instead of I_h . Now $\tilde{u}_b \in H^2(\Omega)$ and, according to the approximation properties of i_h (11) and (13), one has $\|u_h^0\|_{0,\Omega} \leq C \|u_0\|_{2,\Omega}$ and $\|\tilde{u}_{bh}\|_{1,\Omega} + |\tilde{u}_{bh}|_{1,3,\Omega} \leq C \|\tilde{u}_b\|_{2,\Omega}$.

Remark. It is worth remarking that, for the homogeneous case $u_b = 0$, instead of the direct proof presented in this manuscript, an analysis completely analogous to the one given in [8], Corollary 7, leads to the following stability result for $\theta \in [1/2, 1]$ and $N < N_T$

$$(70) \quad \frac{1}{2} \|u_h^N\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \left\{ \|u_h^{n+\theta}\|_{\text{LPS}}^2 + d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, u_h^{n+\theta}) \right\} \\ \leq e^{\frac{T}{T-\delta t}} \left\{ T \delta t \sum_{n=0}^{N-1} \|f^{n+\theta}\|_{0,\Omega}^2 + \frac{1}{2} \|u_h^0\|_{0,\Omega}^2 \right\}.$$

This result, very similar in form to the one in [8] (with the extra control on the nonlinear term, and a slightly smaller right-hand side), is independent of σ_0 , and hence represents an improvement over the way Lemma 12 is presented. The reason to present the direct proof here lies in the non-homogeneous case, where the presence of u_b is responsible for the dependency of the constant on the right-hand side on σ_0^{-1} . In the non-homogeneous case, both proofs lead to essentially equivalent results, the direct proof presented in this work being more straightforward.

Finally, if u_b would be supposed time dependent, then in the first line of the proof of stability there holds $u_h^{n+1} - u_h^n = \tilde{u}_h^{n+1} - \tilde{u}_h^n + \tilde{u}_{bh}^{n+1} - \tilde{u}_{bh}^n$, thus creating an extra right-hand side depending on the time derivative of u_b .

4.2. Error estimates. In this section, error estimates are derived for the solution of the discrete problem (64) with $\theta \in [1/2, 1]$. The error will be analyzed essentially in the quantity which is given by the stability estimate (65). Let us denote the error by $e^\alpha := u^\alpha - u_h^\alpha$ with $\alpha \in [0, N_T]$. Furthermore, to simplify the presentation of our results, we introduce the

quantities

$$\begin{aligned}
E^N &= \|e^N\|_{0,\Omega} + \left(\delta t \sum_{n=0}^{N-1} \|e^{n+\theta}\|_{\text{LPS}}^2 \right)^{1/2}, \\
Q^N &= h \left(|u_0|_{k+1,\Omega} + |u^N|_{k+1,\Omega} + \sigma_0^{-1/2} \|u_t\|_{L^2(0,t^N;H^{k+1}(\Omega))} \right) + \left(\delta t \sum_{n=0}^{N-1} \left(\varepsilon + h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \right. \right. \\
&\quad \left. \left. + h^2 \|\sigma^{n+\theta}\|_{0,\infty,\Omega} + h^2 \sigma_0^{-1} |\mathbf{b}^{n+\theta}|_{1,\infty,\Omega}^2 \right) \left(|u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\Omega}^2 \right) \right)^{1/2}, \\
R^N &= \left(\delta t \sum_{n=0}^{N-1} h^{k+1-d/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \left(|u^n|_{k+1,\Omega}^3 + |u^{n+1}|_{k+1,\Omega}^3 \right) \right)^{1/2}, \\
S^N &= \left(\delta t \sum_{n=0}^{N-1} h^{k+1} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \left(|u^n|_{k+1,\infty,\Omega} + |u^{n+1}|_{k+1,\infty,\Omega} \right) \left(|u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\Omega}^2 \right) \right)^{1/2}, \\
X^N &= \max_{n=0,\dots,N-1} \left(\varepsilon + h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} + \|\sigma^{n+\theta}\|_{0,\infty,\Omega} + \sigma_0^{-1} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega}^2 + \sigma_0^{-1} \|c^{n+\theta}\|_{0,\infty,\Omega}^2 \right)^{1/2}, \\
Y^N &= h^{1/2} \max_{n=0,\dots,N-1} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega}^{1/2},
\end{aligned}$$

where $N = 1, 2, \dots, N_T$.

Theorem 13. *Let $\theta \in [1/2, 1]$ be given. Let the weak solution of (3) satisfy $u, u_t \in L^2(0, T; H^{k+1}(\Omega))$ for some $k \in \{1, \dots, l\}$ and assume $u_{tt} \in L^2(0, T; L^2(\Omega))$. Let $\tilde{u}_b \in H^2(\Omega)$ be an extension of u_b and let $\tilde{u}_{bh} = i_h \tilde{u}_b$. Assume $u_0 \in H^{k+1}(\Omega)$ and let $u_h^0 = i_h u_0$. Let $\{u_h^n\}_{n=0}^{N_T}$ be the solution of the local projection discretization (64). If $\tilde{\tau}_M$ is defined by (22) and $u_t \in L^3(0, T; W^{1,3}(\Omega))$, then the error estimate*

$$\begin{aligned}
(71) \quad E^N &+ \left(\delta t \sum_{n=0}^{N-1} \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M(P_M^{n+\theta} \nabla e^{n+\theta})\|_{0,3,M}^3 \right)^{1/2} \\
&\leq C h^k Q^N + C \beta h^k R^N + C \delta t X^N \|u_t\|_{L^2(0,t^N;H^1(\Omega))} \\
&\quad + C \beta (\delta t)^{3/2} Y^N \|u_t\|_{L^3(0,t^N;W^{1,3}(\Omega))}^{3/2} + C \delta t \sigma_0^{-1/2} \|u_{tt}\|_{L^2(0,t^N;L^2(\Omega))}
\end{aligned}$$

is satisfied for $N = 1, 2, \dots, N_T$. Moreover, if $\theta = 1/2$, $u_{tt} \in L^3(0, T; W^{1,3}(\Omega))$, and $u_{ttt} \in L^2(0, T; L^2(\Omega))$, then

$$\begin{aligned} E^N + \left(\delta t \sum_{n=0}^{N-1} \sum_{M \in \mathcal{M}_h} \tilde{\tau}_M \|\kappa_M (P_M^{n+\theta} \nabla e^{n+\theta})\|_{0,3,M}^3 \right)^{1/2} \\ \leq C h^k Q^N + C \beta h^k R^N + C (\delta t)^2 X^N \|u_{tt}\|_{L^2(0,t^N;H^1(\Omega))} \\ + C \beta (\delta t)^3 Y^N \|u_{tt}\|_{L^3(0,t^N;W^{1,3}(\Omega))}^{3/2} + C (\delta t)^2 \sigma_0^{-1/2} \|u_{ttt}\|_{L^2(0,t^N;L^2(\Omega))}. \end{aligned}$$

If $u \in L^2(0, T; W^{k+1,\infty}(\Omega))$, then, in both estimates, R^N can be replaced by S^N .

If $\tilde{\tau}_M$ is defined by (23) and $u_t \in L^4(0, T; W^{1,4}(\Omega))$, then the following error estimate holds

$$(72) \quad E^N + \left(\delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(u_h^{n+\theta}; e^{n+\theta}, e^{n+\theta}) \right)^{1/2} \leq C h^k Q^N + C \delta t X^N \|u_t\|_{L^2(0,t^N;H^1(\Omega))} \\ + C \delta t T^{1/4} Y^N \|u_t\|_{L^4(0,t^N;W^{1,4}(\Omega))} + C \delta t \sigma_0^{-1/2} \|u_{tt}\|_{L^2(0,t^N;L^2(\Omega))}.$$

Moreover, if $\theta = 1/2$, $u_{tt} \in L^4(0, T; W^{1,4}(\Omega))$, and $u_{ttt} \in L^2(0, T; L^2(\Omega))$, then

$$\begin{aligned} E^N + \left(\delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(u_h^{n+\theta}; e^{n+\theta}, e^{n+\theta}) \right)^{1/2} \leq C h^k Q^N + C (\delta t)^2 X^N \|u_{tt}\|_{L^2(0,t^N;H^1(\Omega))} \\ + C (\delta t)^2 T^{1/4} Y^N \|u_{tt}\|_{L^4(0,t^N;W^{1,4}(\Omega))} + C (\delta t)^2 \sigma_0^{-1/2} \|u_{ttt}\|_{L^2(0,t^N;L^2(\Omega))}. \end{aligned}$$

Proof. Analogously to the steady-state case, the error will be split into an interpolation error and a remainder which belongs to the finite element space. The decomposition of the error e^α with any $\alpha \in [0, N_T]$ has the form

$$e^\alpha = \eta^\alpha - e_h^\alpha \quad \text{with} \quad \eta^\alpha := u^\alpha - \bar{r}_h^\alpha, \quad e_h^\alpha := u_h^\alpha - \bar{r}_h^\alpha \in V_h,$$

where we use the abbreviation $\bar{r}_h^\alpha = \tilde{r}_h^\alpha u$ with \tilde{r}_h^α given by (63). Using this decomposition, one obtains with the triangle inequality and with (60)

$$(73) \quad \|e^N\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \|e^{n+\theta}\|_{\text{LPS}}^2 + \delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(\gamma_0^{n+\theta}; e^{n+\theta}, e^{n+\theta}) \\ \leq 4 \left[\|\eta^N\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \|\eta^{n+\theta}\|_{\text{LPS}}^2 + \delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(\gamma_1^{n+\theta}; \eta^{n+\theta}, \eta^{n+\theta}) \right] \\ + 4 \left[\|e_h^N\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \|e_h^{n+\theta}\|_{\text{LPS}}^2 + \delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(\gamma_2^{n+\theta}; e_h^{n+\theta}, e_h^{n+\theta}) \right],$$

where $\gamma_0^{n+\theta} = e^{n+\theta}$, $\gamma_1^{n+\theta} = \eta^{n+\theta}$, $\gamma_2^{n+\theta} = e_h^{n+\theta}$ if $\tilde{\tau}_M$ is defined by (22) and $\gamma_0^{n+\theta} = \gamma_1^{n+\theta} = \gamma_2^{n+\theta} = u_h^{n+\theta}$ if $\tilde{\tau}_M$ is defined by (23).

First let us estimate the interpolation errors. The starting point is the identity

$$(74) \quad \eta^{n+\theta} = u^{n+\theta} - \theta u^{n+1} - (1-\theta) u^n + \theta (u^{n+1} - r_h u^{n+1}) + (1-\theta) (u^n - r_h u^n).$$

One has

$$(75) \quad u^{n+\theta} - \theta u^{n+1} - (1-\theta) u^n = (1-\theta) \int_{t^n}^{t^{n+\theta}} u_t(t) dt - \theta \int_{t^{n+\theta}}^{t^{n+1}} u_t(t) dt,$$

which, in view of (45), leads to

$$\begin{aligned} \|\eta^{n+\theta}\|_{0,\Omega} &\leq C h^{k+1} (|u^n|_{k+1,\Omega} + |u^{n+1}|_{k+1,\Omega}) + \sqrt{\delta t} \|u_t\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}, \\ |\eta^{n+\theta}|_{1,\Omega} &\leq C h^k (|u^n|_{k+1,\Omega} + |u^{n+1}|_{k+1,\Omega}) + \sqrt{\delta t} \|u_t\|_{L^2(t^n, t^{n+1}; H^1(\Omega))}. \end{aligned}$$

Using Taylor's formula with integral remainder or applying successively integration by parts gives

$$(76) \quad u^n = u^{n+\theta} - \theta \delta t u_t^{n+\theta} + \int_{t^{n+\theta}}^{t^n} u_{tt}(t) (t^n - t) dt,$$

$$(77) \quad u^{n+1} = u^{n+\theta} + (1-\theta) \delta t u_t^{n+\theta} + \int_{t^{n+\theta}}^{t^{n+1}} u_{tt}(t) (t^{n+1} - t) dt.$$

This may be used to derive improved interpolation estimates with respect to the time step provided that $u_{tt} \in L^2(0, T; H^1(\Omega))$. Indeed,

$$(78) \quad u^{n+\theta} - \theta u^{n+1} - (1-\theta) u^n = -(1-\theta) \int_{t^n}^{t^{n+\theta}} u_{tt}(t) (t - t^n) dt - \theta \int_{t^{n+\theta}}^{t^{n+1}} u_{tt}(t) (t^{n+1} - t) dt,$$

which leads to

$$\begin{aligned} \|\eta^{n+\theta}\|_{0,\Omega} &\leq C h^{k+1} (|u^n|_{k+1,\Omega} + |u^{n+1}|_{k+1,\Omega}) + (\delta t)^{3/2} \|u_{tt}\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}, \\ |\eta^{n+\theta}|_{1,\Omega} &\leq C h^k (|u^n|_{k+1,\Omega} + |u^{n+1}|_{k+1,\Omega}) + (\delta t)^{3/2} \|u_{tt}\|_{L^2(t^n, t^{n+1}; H^1(\Omega))}. \end{aligned}$$

Now let us estimate the norms of the interpolation error in (73). In view of (63), (45), (16), (18), and the geometrical hypotheses (5) and (4), one has

$$\begin{aligned} \|\eta^N\|_{0,\Omega} &= \|u^N - r_h u^N\|_{0,\Omega} \leq C h^{k+1} |u^N|_{k+1,\Omega}, \\ \|\eta^{n+\theta}\|_{\text{LPS}} &\leq (\varepsilon + C h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega})^{1/2} |\eta^{n+\theta}|_{1,\Omega} + \|\sigma^{n+\theta}\|_{0,\infty,\Omega}^{1/2} \|\eta^{n+\theta}\|_{0,\Omega}. \end{aligned}$$

Furthermore, analogously as in (54), for any $p \in [2, 6]$, one obtains

$$(79) \quad \begin{aligned} \|\kappa_M(P_M^{n+\theta} \nabla \eta^{n+\theta})\|_{0,p,M} &\leq C |u^{n+\theta} - \theta i_h u^{n+1} - (1-\theta) i_h u^n|_{1,p,M} \\ &\quad + C h_M^{\frac{d}{p} - \frac{d}{2}} (|\varrho_h(u^n - i_h u^n)|_{1,M} + |\varrho_h(u^{n+1} - i_h u^{n+1})|_{1,M}). \end{aligned}$$

If $\tilde{\tau}_M$ is defined by (22), this inequality implies that

$$d_h^{n+\theta}(\eta^{n+\theta}; \eta^{n+\theta}, \eta^{n+\theta}) \leq C \beta (I + II),$$

where

$$\begin{aligned} I &:= h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} |u^{n+\theta} - \theta u^{n+1} - (1-\theta)u^n|_{1,3,M}^3, \\ II &:= h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} (|u^{n+1} - i_h u^{n+1}|_{1,3,M}^3 + |u^n - i_h u^n|_{1,3,M}^3) \\ &\quad + h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} h_M^{-\frac{d}{2}} (|\varrho_h(u^n - i_h u^n)|_{1,M}^3 + |\varrho_h(u^{n+1} - i_h u^{n+1})|_{1,M}^3). \end{aligned}$$

Using (75) and (78), one obtains

$$I \leq C h (\delta t)^2 \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \|u_t\|_{L^3(t^n, t^{n+1}; W^{1,3}(\Omega))}^3,$$

resp.

$$I \leq C h (\delta t)^5 \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \|u_{tt}\|_{L^3(t^n, t^{n+1}; W^{1,3}(\Omega))}^3.$$

Furthermore, it follows from (13), (41), (11), (6), and (4) that

$$(80) \quad II \leq C h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} h_M^{3k-d/2} (|u^n|_{k+1,M}^3 + |u^{n+1}|_{k+1,M}^3),$$

which implies in view of (4) and (5) that

$$II \leq C h^{3k+1-d/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} (|u^n|_{k+1,\Omega}^3 + |u^{n+1}|_{k+1,\Omega}^3).$$

If $u \in L^2(0, T; W^{k+1,\infty}(\Omega))$, the inequality (80) together with (4) and (5) implies that

$$II \leq C h^{3k+1} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} (|u^n|_{k+1,\infty,\Omega} |u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\infty,\Omega} |u^{n+1}|_{k+1,\Omega}^2).$$

If $\tilde{\tau}_M$ is defined by (23), then, proceeding analogously as when deriving (61), but with (79) instead of (54), and applying (13) in addition, one gets

$$d_h^{n+\theta}(u_h^{n+\theta}; \eta^{n+\theta}, \eta^{n+\theta}) \leq C \tilde{I} + C \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} h^{2k+1} (|u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\Omega}^2),$$

where

$$\tilde{I} := h \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \sum_{M \in \mathcal{M}_h} h_M^{d/2} |u^{n+\theta} - \theta u^{n+1} - (1-\theta)u^n|_{1,4,M}^2.$$

Similarly as above, one obtains

$$\tilde{I} \leq C h (\delta t)^{3/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \|u_t\|_{L^4(t^n, t^{n+1}; W^{1,4}(\Omega))}^2,$$

resp.

$$\tilde{I} \leq C h (\delta t)^{7/2} \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} \|u_{tt}\|_{L^4(t^n, t^{n+1}; W^{1,4}(\Omega))}^2.$$

Now let us estimate the norms of the discrete part of the error on the right-hand side of (73). To derive an equation for this part of the error, the weak formulation (62) at $t = t^{n+\theta}$ is subtracted from (64) with $v = v_h = e_h^{n+\theta}$. Then, using the fact that $u_h^\alpha = e_h^\alpha + \bar{r}_h^\alpha$, one deduces that

$$(81) \quad (e_h^{n+1} - e_h^n, e_h^{n+\theta}) + \delta t \|e_h^{n+\theta}\|_{\text{LPS}}^2 + \delta t d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, e_h^{n+\theta}) \\ = \delta t \left[\left(u_t^{n+\theta} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t}, e_h^{n+\theta} \right) + a^{n+\theta}(\eta^{n+\theta}, e_h^{n+\theta}) - s_h^{n+\theta}(\bar{r}_h^{n+\theta}, e_h^{n+\theta}) \right].$$

Furthermore, one obtains

$$(82) \quad d_h^{n+\theta}(u_h^{n+\theta}; u_h^{n+\theta}, e_h^{n+\theta}) \geq \frac{1}{7} d_h^{n+\theta}(\gamma_2^{n+\theta}; e_h^{n+\theta}, e_h^{n+\theta}) + d_h^{n+\theta}(\gamma_3^{n+\theta}; \bar{r}_h^{n+\theta}, e_h^{n+\theta}),$$

where $\gamma_3^{n+\theta} = \bar{r}_h^{n+\theta}$ if $\tilde{\tau}_M$ is defined by (22) and $\gamma_3^{n+\theta} = u_h^{n+\theta}$ if $\tilde{\tau}_M$ is defined by (23) ($\gamma_2^{n+\theta}$ was defined below (73)). This estimate follows from (27) if $\tilde{\tau}_M$ is defined by (22) and simply by writing the second argument of $d_h^{n+\theta}$ as $e_h^{n+\theta} + \bar{r}_h^{n+\theta}$ and using the fact that $d_h^{n+\theta}(u_h^{n+\theta}; e_h^{n+\theta}, e_h^{n+\theta}) \geq 0$ if $\tilde{\tau}_M$ is defined by (23). Since $\theta \geq 1/2$, it follows from (69) with \tilde{u} replaced by e that

$$(83) \quad (e_h^{n+1} - e_h^n, e_h^{n+\theta}) \geq \frac{1}{2} (\|e_h^{n+1}\|_{0,\Omega}^2 - \|e_h^n\|_{0,\Omega}^2).$$

Substituting (82) and (83) into (81) and summing up over the discrete times yields an upper bound for the discrete part of the estimate (73)

$$(84) \quad \|e_h^N\|_{0,\Omega}^2 + \delta t \sum_{n=0}^{N-1} \|e_h^{n+\theta}\|_{\text{LPS}}^2 + \delta t \sum_{n=0}^{N-1} d_h^{n+\theta}(\gamma_2^{n+\theta}; e_h^{n+\theta}, e_h^{n+\theta}) \\ \leq \frac{7}{2} \|e_h^0\|_{0,\Omega}^2 + 7 \delta t \sum_{n=0}^{N-1} \left[\left(u_t^{n+\theta} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t}, e_h^{n+\theta} \right) + a^{n+\theta}(\eta^{n+\theta}, e_h^{n+\theta}) \right. \\ \left. - s_h^{n+\theta}(\bar{r}_h^{n+\theta}, e_h^{n+\theta}) - d_h^{n+\theta}(\gamma_3^{n+\theta}; \bar{r}_h^{n+\theta}, e_h^{n+\theta}) \right].$$

Using (42), the approximation property of i_h (11), (5), and (4), one obtains

$$\|e_h^0\|_{0,\Omega} = \|i_h u^0 - r_h u^0\|_{0,\Omega} = \|\varrho_h(u^0 - i_h u^0)\|_{0,\Omega} \leq C h^{k+1} |u^0|_{k+1,\Omega}.$$

Applying the Cauchy–Schwarz and Young inequalities gives

$$\left(u_t^{n+\theta} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t}, e_h^{n+\theta} \right) \leq \frac{1}{\sigma_0} \left\| u_t^{n+\theta} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t} \right\|_{0,\Omega}^2 + \frac{1}{4} \|e_h^{n+\theta}\|_{\text{LPS}}^2.$$

The last term can be hidden in the left-hand side of (84). The first term is a mixture of discretization errors in time and space. Elimination of $u^{n+\theta}$ from (76) and (77) yields

$$u_t^{n+\theta} = \frac{u^{n+1} - u^n}{\delta t} - \frac{1}{\delta t} \int_{t^n}^{t^{n+\theta}} u_{tt}(t) (t^n - t) dt - \frac{1}{\delta t} \int_{t^{n+\theta}}^{t^{n+1}} u_{tt}(t) (t^{n+1} - t) dt.$$

Since interpolation in space and differentiation in time commute, one has

$$u^{n+1} - \bar{r}_h^{n+1} - (u^n - \bar{r}_h^n) = \int_{t^n}^{t^{n+1}} (u_t - r_h u_t)(t) dt.$$

Thus, applying the Cauchy–Schwarz inequality, one derives

$$\left\| u_t^{n+\theta} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t} \right\|_{0,\Omega}^2 \leq \frac{2}{\delta t} \|u_t - r_h u_t\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}^2 + 2 \delta t \|u_{tt}\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}^2.$$

The first term on the right-hand side can be bounded using (45).

Assuming $u_{ttt} \in L^2(0, T; L^2(\Omega))$ and replacing (76) and (77) by

$$\begin{aligned} u^n &= u^{n+\theta} - \theta \delta t u_t^{n+\theta} + \frac{\theta^2}{2} (\delta t)^2 u_{tt}^{n+\theta} + \frac{1}{2} \int_{t^{n+\theta}}^{t^n} u_{ttt}(t) (t^n - t)^2 dt, \\ u^{n+1} &= u^{n+\theta} + (1 - \theta) \delta t u_t^{n+\theta} + \frac{(1 - \theta)^2}{2} (\delta t)^2 u_{tt}^{n+\theta} + \frac{1}{2} \int_{t^{n+\theta}}^{t^{n+1}} u_{ttt}(t) (t^{n+1} - t)^2 dt, \end{aligned}$$

one obtains

$$\begin{aligned} u_t^{n+\theta} &= \frac{u^{n+1} - u^n}{\delta t} + \frac{\delta t}{2} [\theta^2 - (1 - \theta)^2] u_{tt}^{n+\theta} \\ &\quad - \frac{1}{2 \delta t} \int_{t^n}^{t^{n+\theta}} u_{ttt}(t) (t^n - t)^2 dt - \frac{1}{2 \delta t} \int_{t^{n+\theta}}^{t^{n+1}} u_{ttt}(t) (t^{n+1} - t)^2 dt, \end{aligned}$$

which shows that an improved estimate with respect to δt follows for $\theta = 1/2$, i.e., for the Crank–Nicolson scheme. Indeed, one gets

$$\left\| u_t^{n+1/2} - \frac{\bar{r}_h^{n+1} - \bar{r}_h^n}{\delta t} \right\|_{0,\Omega}^2 \leq \frac{2}{\delta t} \|u_t - r_h u_t\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}^2 + (\delta t)^3 \|u_{ttt}\|_{L^2(t^n, t^{n+1}; L^2(\Omega))}^2.$$

Now let us consider the remaining three terms on the right-hand side of (84). According to (74) and (63), one has

$$\begin{aligned} a^{n+\theta}(\eta^{n+\theta}, e_h^{n+\theta}) - s_h^{n+\theta}(\bar{r}_h^{n+\theta}, e_h^{n+\theta}) &= a^{n+\theta}(u^{n+\theta} - \theta u^{n+1} - (1 - \theta) u^n, e_h^{n+\theta}) \\ &\quad + \theta \left[a^{n+\theta}(u^{n+1} - r_h u^{n+1}, e_h^{n+\theta}) - s_h^{n+\theta}(r_h u^{n+1}, e_h^{n+\theta}) \right] \\ &\quad + (1 - \theta) \left[a^{n+\theta}(u^n - r_h u^n, e_h^{n+\theta}) - s_h^{n+\theta}(r_h u^n, e_h^{n+\theta}) \right]. \end{aligned}$$

The last two terms can be estimated by (46) and the estimation of the first term on the right-hand side is performed using

$$\|u^{n+\theta} - \theta u^{n+1} - (1-\theta)u^n\|_{1,\Omega}^2 \leq \delta t \|u_t\|_{L^2(t^n, t^{n+1}; H^1(\Omega))}^2,$$

resp.

$$\|u^{n+\theta} - \theta u^{n+1} - (1-\theta)u^n\|_{1,\Omega}^2 \leq (\delta t)^3 \|u_{tt}\|_{L^2(t^n, t^{n+1}; H^1(\Omega))}^2,$$

which follows from (75), resp. (78). Finally, the last term on the right-hand side of (84) can be estimated analogously as (52), (56), and (59): if $\tilde{\tau}_M$ is defined by (22), one derives

$$d_h^{n+\theta}(\bar{r}_h^{n+\theta}; \bar{r}_h^{n+\theta}, \bar{r}_h^{n+\theta}) \leq C \beta \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} h^{3k+1-d/2} (|u^n|_{k+1,\Omega}^3 + |u^{n+1}|_{k+1,\Omega}^3),$$

if, in addition, $u \in L^2(0, T; W^{k+1,\infty}(\Omega))$, then

$$\begin{aligned} & d_h^{n+\theta}(\bar{r}_h^{n+\theta}; \bar{r}_h^{n+\theta}, \bar{r}_h^{n+\theta}) \\ & \leq C \beta \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} h^{3k+1} (|u^n|_{k+1,\infty,\Omega} + |u^{n+1}|_{k+1,\infty,\Omega}) (|u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\Omega}^2), \end{aligned}$$

and, if $\tilde{\tau}_M$ is defined by (23), then

$$d_h^{n+\theta}(u_h^{n+\theta}; \bar{r}_h^{n+\theta}, \bar{r}_h^{n+\theta}) \leq C \|\mathbf{b}^{n+\theta}\|_{0,\infty,\Omega} h^{2k+1} (|u^n|_{k+1,\Omega}^2 + |u^{n+1}|_{k+1,\Omega}^2).$$

These estimates together with analogs of (51) and (58) lead to an estimate of the term $d_h^{n+\theta}(\gamma_3^{n+\theta}; \bar{r}_h^{n+\theta}, e_h^{n+\theta})$.

Collecting all the above estimates proves the theorem. \square

At the end of this section, a semi-implicit (linearized) variant of the method (64) will be discussed: For $n = 0, 1, \dots, N_T - 1$, find $u_h^{n+1} \in W_h$ such that $u_h^{n+1} - \tilde{u}_{bh} \in V_h$ and

$$\begin{aligned} (85) \quad & \left(\frac{u_h^{n+1} - u_h^n}{\delta t}, v_h \right) + a^{n+\theta}(u_h^{n+\theta}, v_h) + s_h^{n+\theta}(u_h^{n+\theta}, v_h) + d_h^{n+\theta}(u_h^n; u_h^{n+\theta}, v_h) \\ & = (f^{n+\theta}, v_h) \quad \forall v_h \in V_h. \end{aligned}$$

The advantages of this linearized scheme over (64) in terms of computational complexity are clear. Indeed, for (85) only one linear system needs to be solved per time step. Moreover, the linearized problem is uniquely solvable for any non-negative integrable stabilization parameter τ_M^{sold} . If the parameter $\tilde{\tau}_M$ is defined by (23), the results of Lemma 12 and Theorem 13 remain essentially valid; the only difference is that in these results the first argument of $d_h^{n+\theta}$ is now u_h^n . The proofs of Lemma 12 and Theorem 13 can be repeated without any changes for $\tilde{\tau}_M$ defined by (23) since the estimates of the nonlinear term $d_h^{n+\theta}$ are based on (24) and hence are independent of the first argument of $d_h^{n+\theta}$. This is not the case if $\tilde{\tau}_M$ is defined by (22) and, therefore, we were able to prove only suboptimal convergence results and a

stability result depending on T in a similar way as in (70). Details of this analysis will be omitted here.

5. EXAMPLES OF SPACES AND PARTITIONS SATISFYING THE HYPOTHESES

This section is devoted to the presentation of some examples of spaces W_h and D_M and partitions \mathcal{M}_h satisfying the hypotheses from Section 2. For simplicity, the discussion is restricted to the two-dimensional case. In three dimensions, the spaces can be constructed analogously (for details, see [30]). Throughout this section, $\{\mathcal{T}_h\}_{h>0}$ stands for a regular family of triangulations of $\bar{\Omega}$. This family is formed either by closed triangles or by closed convex quadrilaterals K with diameters h_K and one has $h = \max_{K \in \mathcal{T}_h} h_K$. Note that the hypotheses from Section 2, e.g., (4), (6), and (7), do not allow the application of the analysis to anisotropic triangulations. In what follows, \hat{K} stands for a reference mesh cell, which is either a triangle or a square, depending on the type of elements in \mathcal{T}_h . For any $K \in \mathcal{T}_h$, there exists a bijective mapping $F_K : \hat{K} \rightarrow K$ that maps \hat{K} onto K and is affine if \hat{K} is a triangle and bilinear if \hat{K} is a square. For any integer $l \geq 0$, we denote by P_l the space of polynomials of total degree at most l and by Q_l the space of polynomials of degree at most l in each variable. Finally, we set $R_l(\hat{K}) = P_l(\hat{K})$ if \hat{K} is a triangle and $R_l(\hat{K}) = Q_l(\hat{K})$ if \hat{K} is a square.

i) The two-level approach. This is the approach considered in the original local projection stabilization method (cf. [2, 3]). The starting point is $\{\mathcal{M}_h\}_{h>0}$, a shape regular family of triangulations of $\bar{\Omega}$. Then, each triangle is divided into three triangles by connecting its vertices with the barycenter and each quadrilateral is divided into four quadrilaterals by connecting midpoints of opposite edges. The resulting triangulation is denoted by \mathcal{T}_h . Finally, given an integer $l \geq 1$, the spaces W_h and D_M are given by

$$(86) \quad W_h := \{v_h \in C(\bar{\Omega}); v_h|_K \circ F_K \in R_l(\hat{K}) \ \forall K \in \mathcal{T}_h\}, \quad D_M := P_{l-1}(M).$$

The inf-sup condition (9) is proved for this pair in [30].

Alternatively, for the quadrilateral case, the space D_M could be defined as the space of mapped polynomials. More precisely, we can present the following two alternative definitions for D_M :

$$\begin{aligned} D_M^1 &:= \{v \in L^2(M); v \circ F_M \in P_{l-1}(\hat{M})\}, \\ D_M^2 &:= \{v \in L^2(M); v \circ F_M \in Q_{l-1}(\hat{M})\}, \end{aligned}$$

where \hat{M} is a reference macro-cell and F_M is the analog of F_K . Both definitions lead to different methods (both different from the one presented so far) and have the advantage

that the computations can be done directly on the reference element, leading to simpler implementations. All the approximation and stability assumptions hold for D_M^2 , but for D_M^1 the approximation property (12) holds only on uniformly refined meshes (see [31, pp. 345–346] for a discussion on the topic).

ii) The one-level approach. This alternative was introduced in [30] and assumes $\mathcal{M}_h = \mathcal{T}_h$. Introducing a polynomial bubble function $b_{\widehat{K}} \in H_0^1(\widehat{K}) \setminus \{0\}$ (cubic if \widehat{K} is a triangle and biquadratic if \widehat{K} is a square), the spaces are given by

$$W_h := \{v_h \in C(\overline{\Omega}); v_h|_K \circ F_K \in R_l(\widehat{K}) + b_{\widehat{K}} \cdot R_{l-1}(\widehat{K}) \quad \forall K \in \mathcal{T}_h\}, \quad D_M := P_{l-1}(M).$$

The inf-sup condition (9) is proved for this pair in [30].

iii) The overlapping method. Let x_1, \dots, x_{N_h} be the inner vertices of the triangulation \mathcal{T}_h , introduce the neighborhoods $M_i := \text{int} \bigcup_{K \in \mathcal{T}_h, x_i \in K} K$ (where ‘int’ denotes the interior of the respective set), and define $\mathcal{M}_h := \{M_i\}_{i=1}^{N_h}$. The spaces W_h and D_M are given by (86). The inf-sup condition (9) is proved for this pair in [24].

In all of the examples above, i_h can be chosen to be the Lagrange interpolation operator and j_M to be the orthogonal L^2 projection of $L^2(M)$ onto D_M (see, e.g., [12]). The validity of the geometrical hypotheses (4)–(7) follows from the mesh regularity. The inverse inequality (8) arises from a local inverse inequality (cf. [12]) and the mesh regularity. Finally, if F_K is linear for any $K \in \mathcal{T}_h$, then the space G_M consists of functions that are polynomial on the mesh cells included in M and the inverse inequality (10) is standard (cf. [12]).

Note that if the set \mathcal{M}_h consists of nonoverlapping sets M , which is the case for both the one-level and two-level methods, then (significantly) more degrees of freedom are used for constructing the space W_h than in case of the method with overlapping sets M . This increase of the number of degrees of freedom is either due to an enrichment by bubble functions (in the one-level method) or due to a refinement of the given triangulation (in the two-level method). On the other hand, given a triangulation \mathcal{T}_h of $\overline{\Omega}$ and using \mathcal{M}_h consisting of overlapping sets M , the space W_h can be defined as a standard finite element space consisting of piecewise polynomials of degree l on \mathcal{T}_h , like in the Galerkin discretization.

6. NUMERICAL ILLUSTRATIONS

In this section, the theory of this paper is illustrated by results of numerical computations performed for both the steady-state problem (1) and the time-dependent problem (3). In addition, the reduction of spurious oscillations by applying the nonlinear crosswind diffusion is demonstrated. From the three possibilities for spaces and partitions proposed in the preceding section, we have chosen the overlapping version of the LPS method. This is mainly

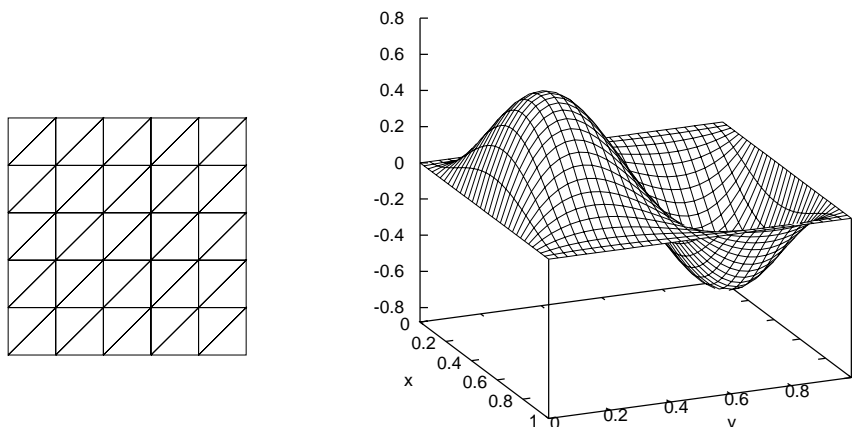


FIGURE 1. Type of the triangulations used in numerical computations (left) and solution for Example 1 (right).

due to the fact that, as shown in [24], the overlapping version is more robust with respect to the stabilization parameter than both the one- and two-level approaches. The overlapping version was applied with triangular meshes and conforming piecewise linear approximation spaces W_h (thus $l = 1$). Both possible definitions (22) and (23) of $\tilde{\tau}_M(u_h)$ were considered. The solution of the nonlinear system was performed using a fixed point iteration: Given an initial approximation $u_h^0 \in W_h$ of the solution of (19) satisfying $u_h^0 - \tilde{u}_{bh} \in V_h$, compute a sequence $\{u_h^k\} \subset W_h$ defined by

$$u_h^k = u_h^{k-1} + \omega (\tilde{u}_h^k - u_h^{k-1}), \quad k = 1, 2, \dots,$$

where $\omega \in (0, 1]$ is a damping factor and $\tilde{u}_h^k \in W_h$ satisfies $\tilde{u}_h^k - \tilde{u}_{bh} \in V_h$ and

$$a(\tilde{u}_h^k, v_h) + s_h(\tilde{u}_h^k, v_h) + d_h(u_h^{k-1}; \tilde{u}_h^k, v_h) = (f, v_h) \quad \forall v_h \in V_h.$$

The analysis of the convergence of this scheme remains an open problem. Its proof, based on the properties of the nonlinear operator from Section 3, does not seem an easy task. The actual behavior of the iteration in our numerical studies will be discussed in Example 2.

In all examples, $\Omega = (0, 1)^2$ and Friedrichs–Keller triangulations of the type depicted in Fig. 1 were used. It is worth mentioning that the mesh is not aligned with the considered convection fields.

Example 1. Smooth polynomial solution [20], support of error estimates. We considered problem (1) with $\varepsilon = 10^{-8}$, $\mathbf{b} = (3, 2)^T$, $c = 2$, and $u_b = 0$. The right-hand side f was chosen such that

$$u(x, y) = 100 x^2 (1 - x)^2 y (1 - y) (1 - 2y)$$

TABLE 1. Example 1, errors of the discrete solutions.

	parameter (22)				parameter (23)			
h	$\ \cdot\ _{\text{LPS}}$	$\ \cdot\ _{0,\Omega}$	$ \cdot _{1,\Omega}$	$\ \cdot\ _{0,\infty,h}$	$\ \cdot\ _{\text{LPS}}$	$\ \cdot\ _{0,\Omega}$	$ \cdot _{1,\Omega}$	$\ \cdot\ _{0,\infty,h}$
$8.84-2$	$4.74-2$	$1.83-2$	$4.20-1$	$6.46-2$	$4.30-2$	$1.47-2$	$4.00-1$	$5.04-2$
$4.42-2$	$1.48-2$	$3.54-3$	$1.88-1$	$1.52-2$	$1.41-2$	$2.93-3$	$1.84-1$	$1.13-2$
$2.21-2$	$5.02-3$	$7.24-4$	$9.02-2$	$3.40-3$	$4.93-3$	$6.57-4$	$8.96-2$	$2.44-3$
$1.10-2$	$1.76-3$	$1.58-4$	$4.45-2$	$7.63-4$	$1.75-3$	$1.57-4$	$4.44-2$	$5.57-4$
$5.52-3$	$6.19-4$	$3.63-5$	$2.21-2$	$1.77-4$	$6.18-4$	$3.83-5$	$2.21-2$	$1.44-4$
order	1.50	2.12	1.01	2.11	1.50	2.03	1.01	1.95

is the solution of (1), see Fig. 1.

In the stabilization parameters, the values $\tau_0 = 0.02$ and $\beta = 0.1$ were used. Table 1 shows errors of the discrete solutions measured in various norms for various mesh sizes. The notation $\|\cdot\|_{0,\infty,h}$ is used for the discrete L^∞ norm defined as the maximum of the errors at the vertices of the respective triangulation. The convergence orders were computed using values from the two finest triangulations. One can observe that the convergence order with respect to the LPS norm is $3/2$, as predicted by the theory, and that in other norms one obtains the usual optimal convergence orders.

Example 2. Solution with two interior layers [27], reduction of spurious oscillations. Equation (1) was considered with $\varepsilon = 10^{-8}$, $\mathbf{b}(x, y) = (-y, x)^T$, $c = f = 0$, and the boundary condition

$$u = u_b \quad \text{on } \Gamma^D, \quad \frac{\partial u}{\partial \mathbf{n}} = 0 \quad \text{on } \Gamma^N,$$

where $\Gamma^N = \{0\} \times (0, 1)$, $\Gamma^D = \partial\Omega \setminus \overline{\Gamma^N}$, \mathbf{n} is the outward pointing unit normal vector to the boundary of Ω , and

$$u_b(x, y) = \begin{cases} 1 & \text{for } (x, y) \in (1/3, 2/3) \times \{0\}, \\ 0 & \text{else on } \Gamma^D. \end{cases}$$

Results that were obtained on the triangulation having 33×33 vertices are presented. Figure 2 shows solutions computed by means of the LPS method with and without the nonlinear crosswind diffusion term d_h defined using the parameter (23). One can observe that the crosswind diffusion term manages to reduce the oscillations appearing in the solution of the linear LPS method. An increase of the parameter β does not only reduce the oscillations but also increases the smearing appearing at the layers. In this respect, the method behaves as expected. Two results obtained for d_h defined using the parameter (22) are shown in Fig. 3. A detailed comparison of the results in Figs. 2 and 3 reveals that the method with

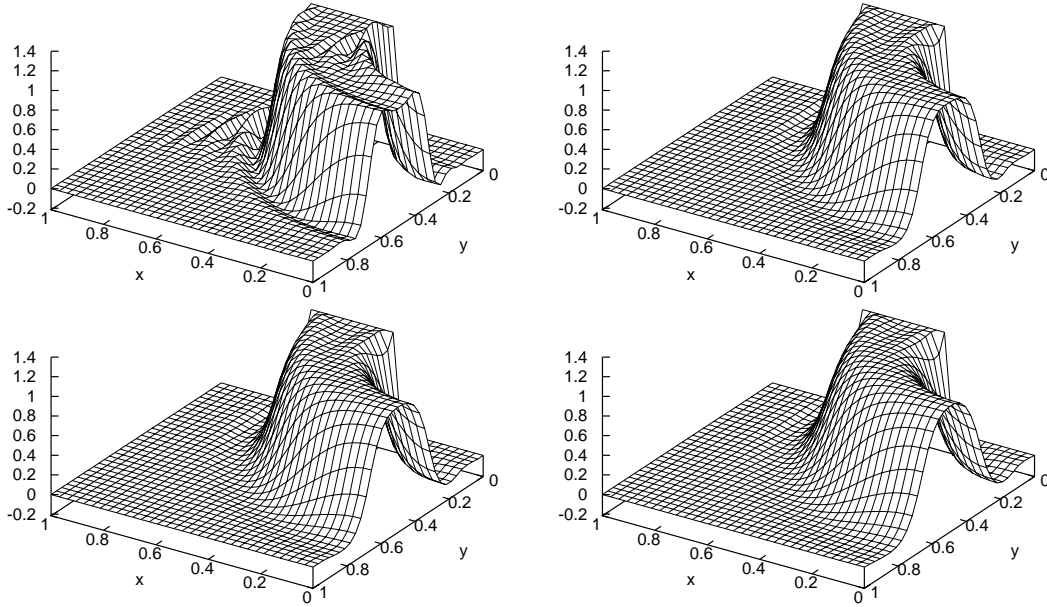


FIGURE 2. Example 2: solutions for the parameter (23) with $\tau_0 = 0.02$ and $\beta = 0, \beta = 0.03, \beta = 0.05, \beta = 0.1$, left to right, top to bottom.

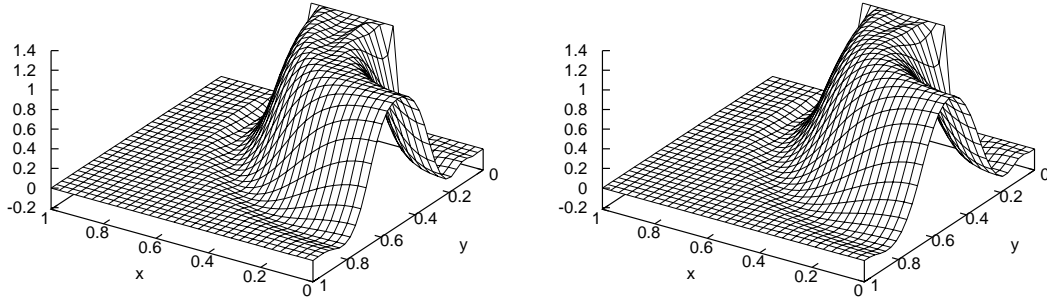


FIGURE 3. Example 2: solutions for the parameter (22) with $\tau_0 = 0.02, \beta = 0.03$ (left) and $\tau_0 = 0.02, \beta = 0.1$ (right).

the parameter (22) is less successful in suppressing spurious oscillations whereas it leads to a more pronounced smearing.

It is natural to ask whether similar results as presented above can be obtained using a linear crosswind diffusion term. To this end, the term d_h with

$$(87) \quad \tau_M^{\text{sold}} = \beta h_M |\mathbf{b}_M|$$

was considered. All other settings were the same as above. Since it is difficult to compare various solutions, we first concentrated on the outflow profile, i.e., the solution graph along the line $x = 0$. For $\beta \leq 0.02$, the outflow profile contains overshoots that decrease with increasing β . Fig. 4 shows that, for $\beta = 0.025$, the overshoots are not present in the outflow

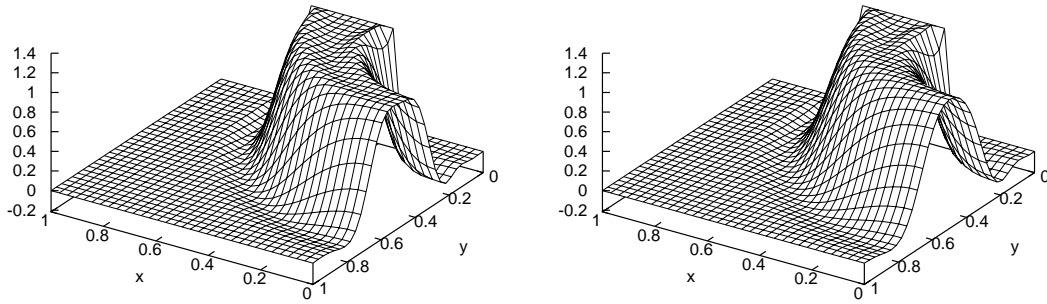


FIGURE 4. Example 2: solutions for the parameter (87) with $\tau_0 = 0.02$, $\beta = 0.025$ (left) and $\tau_0 = 0.02$, $\beta = 0.06$ (right).

profile but they can be still observed inside the computational domain. For this value of β , the outflow profile does not differ too much from the outflow profile in Fig. 2, top right. However, inside the computational domain, both overshoots and undershoots are larger for the linear method. A further increase of β leads to a reduction of the undershoots but also to a smearing of the solution whereas the magnitude of the undershoots does not change significantly. As an example, the solution for $\beta = 0.06$ is shown in Fig. 4. The smearing and the undershoots of this solution are more pronounced than in case of all the three solutions of the nonlinear method in Fig. 2. This study demonstrates that the method with linear crosswind diffusion was outperformed, with respect to the quality of the computed solution, by the nonlinear method with $\tilde{\tau}_M$ defined by (23).

From the discussion of the preceding paragraphs, the choice of the stabilization parameter β appears as an important issue. A good choice of user-chosen parameters in stabilized finite element methods is an open problem for all methods. In general, the parameters need to be chosen not constant but as functions (see [18] for the construction of an example). A non-constant choice, done automatically like in [19], will be the subject of future research.

Next, the computational cost connected with the solution of the nonlinear discrete problems will be briefly illustrated. Table 2 shows numbers of fixed-point iterations needed to solve Example 2 for $\tau_0 = 0.02$ and various values of β and the damping parameter ω . The iterative process was terminated if the Euclidean norm of the residual of the nonlinear algebraic system divided by the Euclidean norm of its right-hand side was smaller than 10^{-8} . The sequences of the residuals were monotonically decreasing, except for some of the computations with the parameter (22) for $\omega \in \{0.9, 1\}$ where oscillations of the residuals appeared at the beginning of the iterative process. One can observe that the number of iterations depends both on β and ω and that this dependence is more pronounced if the parameter $\tilde{\tau}_M$ is defined by (22). Since the optimal value of the damping parameter is usually not known,

TABLE 2. Example 2, number of fixed-point iterations.

	parameter (22)				parameter (23)			
	$\beta = 0.01$	$\beta = 0.03$	$\beta = 0.06$	$\beta = 0.10$	$\beta = 0.01$	$\beta = 0.03$	$\beta = 0.06$	$\beta = 0.10$
$\omega = 1.0$	82	163	305	494	16	27	39	51
$\omega = 0.9$	42	58	68	73	12	18	24	29
$\omega = 0.8$	25	30	32	33	12	13	16	19
$\omega = 0.7$	16	17	18	20	16	16	16	16
$\omega = 0.6$	20	20	20	20	21	21	21	21
$\omega = 0.5$	27	27	27	27	27	27	27	27

it can be expected that the numerical effort caused by the nonlinear crosswind diffusion term will be generally smaller if the parameter $\tilde{\tau}_M$ is defined by (23).

Example 3. Smooth time-dependent solution, support of error estimates. The setup of this example is very similar to Example 6.1 in [22]. Problem (3) was considered in the time interval $[0, 1]$ with $\varepsilon = 10^{-8}$, $\mathbf{b} = (3, 2)^T$, $c = 2$, and $u_b = 0$. The right-hand side f and the initial condition u_0 were chosen such that

$$u(x, y, t) = e^{\sin(2\pi t)} \sin(2\pi x) \sin(2\pi y)$$

is the solution of (3).

We considered the discrete problem (64) and its linearized variant (85) with $\theta = 1$ (i.e., the backward Euler scheme) for both choices of $\tilde{\tau}_M$. Like in Example 1, the values $\tau_0 = 0.02$ and $\beta = 0.1$ were used for the stabilization parameters. According to error estimates (71) and (72), one expects that the quantity E^N tends to zero with the convergence order $3/2$ if $\delta t \sim h^{3/2}$ and a nonlinear discretization is used (note the extra power of $h^{1/2}$ in Q^N and R^N). The same convergence behavior is expected for the linearized method if $\tilde{\tau}_M$ is defined by (23), see the discussion at the end of Section 4. These expectations are supported by the results presented in Fig. 5. In this figure, level 1 corresponds to the grid with mesh cells of diameter $h = \sqrt{2}\tilde{h}$ with $\tilde{h} = 1/8$. Uniform refinement in space was used and the length of the time step was set to be $\delta t = \tilde{h}^{3/2}$. If the final time was not obtained exactly with these time steps, the simulations were terminated at the last discrete time smaller than $T = 1$. It can be observed in Fig. 5 that the order of convergence $3/2$ was obtained for the error in the l^2 -LPS norm for all four methods. We could observe the same order of convergence also for $\|e^N\|_{0,\Omega}$. Using the time step $\delta t = \tilde{h}^2$, the error $\|e^N\|_{0,\Omega}$ showed even second order convergence, whereas the order of convergence of the error in the l^2 -LPS norm was still $3/2$. This result demonstrates the sharpness of the estimates (71) and (72).

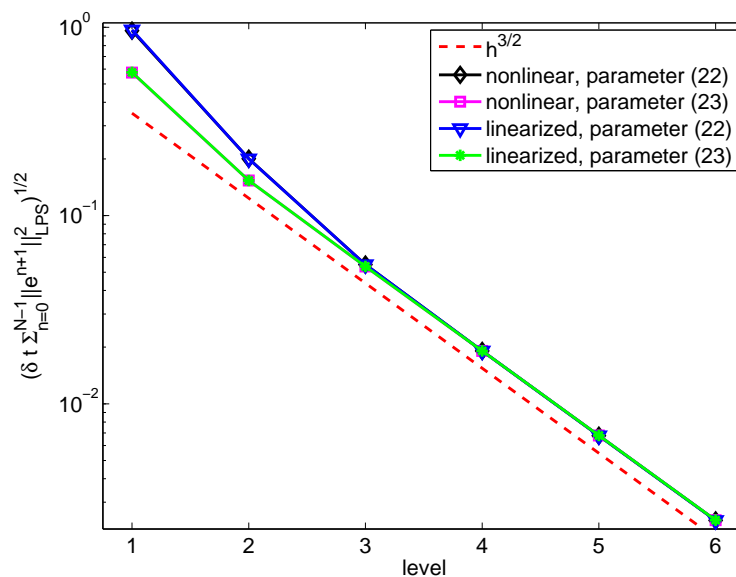


FIGURE 5. Example 3: order of convergence for piecewise linear finite elements, the backward Euler scheme, and $\delta t \sim h^{3/2}$. Note that the curves for the linearized methods are on top of the curves of the corresponding nonlinear method.

Concerning a comparison of the fully nonlinear and the linearized version of the methods, only very little differences can be seen in this example. On coarser grids, the solutions computed using the parameter (23) were more accurate compared with the solutions obtained using the parameter (22).

ACKNOWLEDGMENTS

The work of G.R. Barrenechea has been partially funded by The Leverhulme Trust, through the Research Project Grant RPG-2012-483. The work of P. Knobloch is a part of the research project MSM 0021620839 funded by the Ministry of Education, Youth and Sports of the Czech Republic and it was partly supported by the Grant Agency of the Czech Republic under the grant No. P201/11/1304.

REFERENCES

- [1] M. AUGUSTIN, A. CAIAZZO, A. FIEBACH, J. FUHRMANN, V. JOHN, A. LINKE, AND R. UMLA, An assessment of discretizations for convection-dominated convection–diffusion equations, *Comput. Methods Appl. Mech. Engrg.*, 200 (2011), pp. 3395–3409.
- [2] R. BECKER AND M. BRAACK, A finite element pressure gradient stabilization for the Stokes equations based on local projections, *Calcolo*, 38 (2001), pp. 173–199.

- [3] ———, A two-level stabilization scheme for the Navier–Stokes equations, in *Numerical Mathematics and Advanced Applications*, M. Feistauer, V. Dolejší, P. Knobloch, and K. Najzar, eds., Springer-Verlag, Berlin, 2004, pp. 123–130.
- [4] M. BRAACK AND E. BURMAN, Local projection stabilization for the Oseen problem and its interpretation as a variational multiscale method, *SIAM J. Numer. Anal.*, 43 (2006), pp. 2544–2566.
- [5] M. BRAACK, E. BURMAN, V. JOHN, AND G. LUBE, Stabilized finite element methods for the generalized Oseen problem, *Comput. Methods Appl. Mech. Engrg.*, 196 (2007), pp. 853–866.
- [6] 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, *Comput. Methods Appl. Mech. Engrg.*, 32 (1982), pp. 199–259.
- [7] E. BURMAN AND A. ERN, Stabilized Galerkin approximation of convection–diffusion–reaction equations: discrete maximum principle and convergence, *Math. Comp.*, 74 (2005), pp. 1637–1652.
- [8] E. BURMAN AND M. A. FERNÁNDEZ, Finite element methods with symmetric stabilization for the transient convection–diffusion–reaction equation, *Comput. Methods Appl. Mech. Engrg.*, 198 (2009), pp. 2508–2519.
- [9] E. BURMAN AND P. HANSBO, Edge stabilization for Galerkin approximations of convection–diffusion–reaction problems, *Comput. Methods Appl. Mech. Engrg.*, 193 (2004), pp. 1437–1453.
- [10] P. G. CIARLET, The finite element method for elliptic problems, North-Holland, Amsterdam, 1978.
- [11] R. CODINA, A discontinuity-capturing crosswind-dissipation for the finite element solution of the convection–diffusion equation, *Comput. Methods Appl. Mech. Engrg.*, 110 (1993), pp. 325–342.
- [12] A. ERN AND J.-L. GUERMOND, Theory and Practice of Finite Elements, Springer-Verlag, New York, 2004.
- [13] L. P. FRANCA, S. L. FREY, AND T. J. R. HUGHES, Stabilized finite element methods: I. Application to the advective–diffusive model, *Comput. Methods Appl. Mech. Engrg.*, 95 (1992), pp. 253–276.
- [14] L. P. FRANCA AND F. VALENTIN, On an improved unusual stabilized finite element method for the advective–reactive–diffusive equation, *Comput. Methods Appl. Mech. Engrg.*, 190 (2000), pp. 1785–1800.
- [15] S. GANESAN AND L. TOBISKA, Stabilization by local projection for convection–diffusion and incompressible flow problems, *J. Sci. Comput.*, 43 (2010), pp. 326–342.
- [16] T. J. R. HUGHES, L. P. FRANCA, AND G. M. HULBERT, A new finite element formulation for computational fluid dynamics. VIII. The Galerkin/least-squares method for advective–diffusive equations, *Comput. Methods Appl. Mech. Engrg.*, 73 (1989), pp. 173–189.
- [17] V. JOHN AND P. KNOBLOCH, On spurious oscillations at layers diminishing (SOLD) methods for convection–diffusion equations: Part I – A review, *Comput. Methods Appl. Mech. Engrg.*, 196 (2007), pp. 2197–2215.
- [18] ———, On spurious oscillations at layers diminishing (SOLD) methods for convection–diffusion equations: Part II – Analysis for P_1 and Q_1 finite elements, *Comput. Methods Appl. Mech. Engrg.*, 197 (2008), pp. 1997–2014.

- [19] V. JOHN, P. KNOBLOCH, AND S. B. SAVESCU, A posteriori optimization of parameters in stabilized methods for convection–diffusion problems – Part I, *Comput. Methods Appl. Mech. Engrg.*, 200 (2011), pp. 2916–2929.
- [20] V. JOHN, J. M. MAUBACH, AND L. TOBISKA, Nonconforming streamline–diffusion–finite–element–methods for convection–diffusion problems, *Numer. Math.*, 78 (1997), pp. 165–188.
- [21] V. JOHN, T. MITKOVA, M. ROLAND, K. SUNDMACHER, L. TOBISKA, AND A. VOIGT, Simulations of population balance systems with one internal coordinate using finite element methods, *Chem. Engrg. Sci.*, 64 (2009), pp. 733 – 741.
- [22] V. JOHN AND J. NOVO, Error analysis of the SUPG finite element discretization of evolutionary convection–diffusion–reaction equations, *SIAM J. Numer. Anal.*, 49 (2011), pp. 1149–1176.
- [23] V. JOHN AND E. SCHMEYER, Finite element methods for time-dependent convection–diffusion–reaction equations with small diffusion, *Comput. Methods Appl. Mech. Engrg.*, 198 (2008), pp. 475–494.
- [24] P. KNOBLOCH, A generalization of the local projection stabilization for convection–diffusion–reaction equations, *SIAM J. Numer. Anal.*, 48 (2010), pp. 659–680.
- [25] ———, Local projection method for convection–diffusion–reaction problems with projection spaces defined on overlapping sets, in *Numerical Mathematics and Advanced Applications 2009, Proceedings of ENUMATH 2009*, G. Kreiss, P. Lötstedt, A. Målqvist, and M. Neytcheva, eds., Springer-Verlag, Berlin, 2010, pp. 497–505.
- [26] P. KNOBLOCH AND G. LUBE, Local projection stabilization for advection–diffusion–reaction problems: One-level vs. two-level approach, *Appl. Numer. Math.*, 59 (2009), pp. 2891–2907.
- [27] T. KNOPP, G. LUBE, AND G. RAPIN, Stabilized finite element methods with shock capturing for advection–diffusion problems, *Comput. Methods Appl. Mech. Engrg.*, 191 (2002), pp. 2997–3013.
- [28] O. A. LADYZHENSKAYA, New equations for the description of motion of viscous incompressible fluids and solvability in the large of boundary value problems for them, *Tr. Mat. Inst. Steklova*, 102 (1967), pp. 85–104.
- [29] G. LUBE AND G. RAPIN, Residual-based stabilized higher-order FEM for advection-dominated problems, *Comput. Methods Appl. Mech. Engrg.*, 195 (2006), pp. 4124–4138.
- [30] G. MATTHIES, P. SKRZYPACZ, AND L. TOBISKA, A unified convergence analysis for local projection stabilizations applied to the Oseen problem, *M2AN Math. Model. Numer. Anal.*, 41 (2007), pp. 713–742.
- [31] H.-G. ROOS, M. STYNES, AND L. TOBISKA, Robust Numerical Methods for Singularly Perturbed Differential Equations. Convection–Diffusion–Reaction and Flow Problems. 2nd ed., Springer-Verlag, Berlin, 2008.
- [32] R. TEMAM, Navier-Stokes equations. Theory and numerical analysis, North-Holland, Amsterdam, 1977.

DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF STRATHCLYDE, 26 RICHMOND STREET, GLASGOW G1 1XH, SCOTLAND

E-mail address: gabriel.barrenechea@strath.ac.uk

WEIERSTRASS INSTITUTE FOR APPLIED ANALYSIS AND STOCHASTICS (WIAS), MOHRENSTR. 39, 10117 BERLIN, GERMANY AND FREE UNIVERSITY OF BERLIN, DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE, ARNIMALLEE 6, 14195 BERLIN, GERMANY.

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

DEPARTMENT OF NUMERICAL MATHEMATICS, FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY, SOKOLOVSKÁ 83, 18675 PRAHA 8, CZECH REPUBLIC.

E-mail address: knobloch@karlin.mff.cuni.cz