

## STABILIZED FINITE ELEMENT METHODS FOR NONSYMMETRIC, NONCOERCIVE, AND ILL-POSED PROBLEMS. PART II: HYPERBOLIC EQUATIONS\*

ERIK BURMAN<sup>†</sup>

**Abstract.** In this paper we consider stabilized finite element methods for hyperbolic transport equations without coercivity. Abstract conditions for the convergence of the methods are introduced and these conditions are shown to hold for three different stabilized methods: the Galerkin least squares method, the continuous interior penalty method, and the discontinuous Galerkin method. We consider both the standard stabilization methods and the optimization-based method introduced in [E. Burman, *SIAM J. Sci. Comput.*, 35 (2013), pp. A2752–A2780]. The main idea of the latter is to write the stabilized method in an optimization framework and select the discrete function for which a certain cost functional, in our case the stabilization term, is minimized. Some numerical examples illustrate the theoretical investigations.

**Key words.** finite elements, stabilization, noncoercive problems, hyperbolic equations, transport

**AMS subject classifications.** 65N12, 65N15, 65N20, 65N30, 35L02

**DOI.** 10.1137/130931667

**1. Introduction.** Several finite element methods have been proposed for the computation of hyperbolic problems, such as the SUPG method [5, 16], the discontinuous Galerkin (DG) method [19, 18, 17], and several different weakly consistent, symmetric stabilization methods for continuous approximation spaces [14, 11, 9, 3]. In most of these cases, however, the analysis relies on the satisfaction of a coercivity condition. Indeed if a scalar hyperbolic transport equation

$$(1.1) \quad \beta \cdot \nabla u + \sigma u = f$$

is considered, with data given on the inflow boundary, it is typically assumed that there exists  $\sigma_0 \in \mathbb{R}^+$  such that

$$(1.2) \quad \sigma_0 \leq \inf_{x \in \Omega} \left( \sigma - \frac{1}{2} \nabla \cdot \beta \right).$$

In, for instance, [16, 17, 1] the degenerate case  $\sigma_0 = 0$  is allowed using special exponentially weighted test functions, which we will also exploit in this paper.

In practice this condition is quite restrictive and rules out many important flow regimes such as exothermic reactions, compressible flow fields, or data assimilation problems with data given on the outflow boundary. Our objective in the present paper is to propose an analysis of stabilized finite element methods in the noncoercive case. Indeed similarly as in the elliptic case [20] the discrete solutions of standard stabilized finite element methods are shown to exist and have optimal convergence under a condition on the mesh size. Unlike the elliptic case there appears to be no

---

\*Submitted to the journal's Methods and Algorithms for Scientific Computing section August 2, 2013; accepted for publication (in revised form) April 24, 2014; published electronically August 19, 2014. This research was partially supported provided by EPSRC (award EP/J002313/1).

<http://www.siam.org/journals/sisc/36-4/93166.html>

<sup>†</sup>Department of Mathematics, University College London, London, UK–WC1E 6BT, United Kingdom (e.burman@ucl.ac.uk).

equivalent result, even suboptimal, for the standard Galerkin method. This part uses tools similar to those of [16, 17, 1]. Then we show how the method introduced in [6] can be applied to hyperbolic problems beyond the coercive regime of condition (1.2). The advantage of this latter method is that the mesh conditions under which the analysis holds are much less restrictive and boundary conditions may be imposed on the outflow boundary just as easily as on the inflow boundary, without modifying the parameters of the method. For a full motivation of the method and analysis in the elliptic case see [6].

We will consider problem (1.1) with smooth coefficients  $\beta \in [W^{2,\infty}(\Omega)]^d$  and  $\sigma \in W^{1,\infty}(\Omega)$ . Boundary data will be given on either the inflow or the outflow corresponding to solving either the standard transport problem or a model data assimilation problem. For such smooth physical parameters both cases can easily be solved using the method of characteristics, provided that for each  $x \in \Omega$  there exists a streamline leading, in finite time, to the boundary where data is imposed and  $|\beta(x)| \neq 0$  for all  $x \in \Omega$ . In the following we always assume that  $\beta$  satisfies these assumptions, unless otherwise stated, and that the stationary problem admits a unique, sufficiently smooth solution.

Problems on conservation form  $\nabla \cdot (\beta u)$  are cast on the form (1.1) by using the product rule and including the low order term with coefficient  $\nabla \cdot \beta$  in  $\sigma$ . The present paper has the following structure. In section 2 we propose an abstract analysis under certain assumptions on the discrete bilinear form. Then in section 3 we give a detailed description of how three different stabilization methods, the Galerkin least squares (GLS) method, the continuous interior penalty (CIP) method, and the DG method, satisfy the assumptions of the abstract theory for the case of the advection-reaction equation. In all cases we prove that the classical quasi-optimal estimate for stabilized methods holds,

$$\|u - u_h\|_{L^2(\Omega)} + \|h^{\frac{1}{2}}\beta \cdot \nabla(u - u_h)\|_{L^2(\Omega)} \leq Ch^{k+\frac{1}{2}}|u|_{H^{k+1}(\Omega)}.$$

We also show how to include a model problem for data assimilation in the analysis. Finally, in section 4 we illustrate the theory with some numerical examples.

**2. Abstract formulation.** Let  $\Omega$  be a polygonal/polyhedral subset of  $\mathbb{R}^d$ . The boundary of  $\Omega$  will be denoted by  $\partial\Omega$  and its outward pointing normal by  $n$ . We let  $V, W$  denote two Hilbert spaces with norms  $\|\cdot\|_V$  and  $\|\cdot\|_W$ . The abstract weak formulation of the continuous problem takes the following form: find  $u \in V$  such that

$$(2.1) \quad a(u, v) = (f, v) \quad \forall v \in W$$

with formal adjoint: find  $z \in W$  such that

$$(2.2) \quad a(w, z) = (g, w) \quad \forall w \in V.$$

The bilinear form  $a(\cdot, \cdot) : V \times W \rightarrow \mathbb{R}$  and the data  $f$  are assumed to satisfy the assumptions of Babuska's theorem [2] so that problems (2.1) and (2.2) are well-posed. (See [13] for an analysis of (1.1) in the coercive regime.) We denote the forward problem on strong form  $\mathcal{L}u = f$  and the adjoint problem on strong form  $\mathcal{L}^*z = g$ .

*Remark 1.* The analysis below never uses the full power of Babuska's theorem. We only need to assume that (2.1) admits a unique solution for the given data and that certain discrete stability conditions are satisfied by  $a(\cdot, \cdot)$  as specified below. For the problems considered here the solution of (2.2) will always be  $z = 0$ .

**2.1. Finite element discretization.** Let  $\{\mathcal{T}_h\}_h$  denote a family of quasi-uniform, shape regular triangulations  $\mathcal{T}_h := \{K\}$ , indexed by the maximum triangle radius  $h := \max_{K \in \mathcal{T}_h} h_K$ . The set of faces of the triangulation will be denoted by  $\mathcal{F}$  and  $\mathcal{F}_{int}$  denotes the subset of interior faces. Let  $X_h^k$  denote the finite element space of piecewise polynomial functions on  $\mathcal{T}_h$ ,

$$X_h^k := \{v_h \in L^2(\Omega) : v_h|_K \in \mathbb{P}_k(K) \quad \forall K \in \mathcal{T}_h\}.$$

Here  $\mathbb{P}_k(K)$  denotes the space of polynomials of degree less than or equal to  $k$  on a triangle  $K$ . The  $L^2$ -scalar product over some measurable  $X \subset \mathbb{R}^d$  is denoted  $(\cdot, \cdot)_X$  and the associated norm  $\|\cdot\|_X$ , and the subscript is dropped whenever  $X = \Omega$ . We will also use  $\langle \cdot, \cdot \rangle_Y$  to denote the  $L^2$ -scalar product over  $Y \subset \mathbb{R}^{d-1}$ . For the elementwise  $L^2$ -scalar product and norm over  $\Omega$  we will use the notation  $(\cdot, \cdot)_h := \sum_{K \in \mathcal{T}_h} (\cdot, \cdot)_K$ ,  $\|\cdot\|_h := (\cdot, \cdot)_h^{\frac{1}{2}}$ . In the estimates of the paper capital constants are generic, whereas lowercase constants are specific to the estimate. Sometimes capital constants will be given subscripts to point to the main dependencies on parameters. We will also use  $a \sim b$  to stress an important dependence in  $a$  on some parameter  $b$ , i.e.,  $a = Cb$ , with  $C$  assumed to be moderate.

We let  $\pi_L$  denote the standard  $L^2$ -projection onto  $X_h^k$  and  $i_h : C^0(\bar{\Omega}) \mapsto X_h^k$  the standard Lagrange interpolant. Recall that for any function  $u \in (V \cup W) \cap H^{k+1}(\Omega)$  there holds

$$(2.3) \quad \|u - i_h u\| + h \|\nabla(u - i_h u)\| + h^2 \|D^2(u - i_h u)\|_h \leq c_i h^{k+1} |u|_{H^{k+1}(\Omega)},$$

where  $D^2$  denotes the Hessian matrix, and the matrix norm used is the Frobenius norm. A similar result holds for  $\pi_L$ . If  $\pi_L$  projects onto  $X_h^k \cap C^0(\bar{\Omega})$  the same result holds under the assumption of local quasi regularity of the mesh. The following discrete commutator property follows by straightforward modifications of the result in [4] and holds for  $i_h$ , the elementwise  $L^2$ -projection onto  $X_h^k$ , and, under our assumptions on the mesh, for the  $L^2$ -projection onto continuous finite element functions. Here  $\varphi \in W^{2,\infty}(\Omega)$ ,  $0 \leq n \leq 2$ ,

$$(2.4) \quad \sum_{K \in \mathcal{T}_h} |\varphi u_h - i_h(\varphi u_h)|_{H^n(K)}^2 \leq c_{dc,n,\varphi}^2 h^{-2n+2} \|u_h\|_{L^2(\Omega)}^2.$$

We also note that the following inverse inequalities hold,  $\exists c_T, c_I \in \mathbb{R}^+$  such that

$$(2.5) \quad \|u\|_{\partial K} \leq c_T (h^{-\frac{1}{2}} \|u\|_K + h^{\frac{1}{2}} \|\nabla u\|_K) \quad \forall u \in H^1(K),$$

$$h_K^{-\frac{1}{2}} \|u_h\|_{\partial K} + h_K \|\nabla u_h\|_K \leq c_I \|u_h\|_K \quad \forall u_h \in \mathbb{P}_k(K).$$

Let  $V_h$  and  $W_h$  denote two finite element spaces such that  $\dim V_h = \dim W_h$  (in practice  $V_h = W_h$  herein). Now we introduce a discrete bilinear form  $a_h(\cdot, \cdot) : V_h \times W_h \mapsto \mathbb{R}$  associated to  $a(\cdot, \cdot)$  and a stabilization operator  $s_p(\cdot, \cdot) : V_h \times W_h \mapsto \mathbb{R}$ . The standard stabilized finite element formulation for problem (2.1) takes the following form: find  $u_h \in V_h$  such that

$$(2.6) \quad a_h(u_h, v_h) + s_p(u_h, v_h) = (f, v_h) + s_p(u, v_h) \quad \forall v_h \in W_h.$$

Observe that since  $s_p(u, v_h)$  appears in the right-hand side, we can only use stabilization operators such that this quantity is known. As we shall see below, the

noncoercivity of the form  $a_h(\cdot, \cdot)$  leads to problem-dependent mesh conditions for the well-posedness of (2.6). To alleviate the conditions on the mesh we propose the following finite element method for the approximation of (2.1): find  $(u_h, z_h) \in V_h \times W_h$  such that

$$(2.7) \quad \begin{aligned} a_h(u_h, w_h) + s_a(z_h, w_h) &= (f, w_h), \\ a_h(v_h, z_h) - s_p(u_h, v_h) &= -s_p(u, v_h) \end{aligned}$$

for all  $(v_h, w_h) \in V_h \times W_h$ . Here  $s_a(\cdot, \cdot)$  is a stabilization term related to the adjoint equation that will be discussed below. Observe that we here solve simultaneously (2.1) and (2.2) with  $g = 0$  in the latter equation. We will consider either continuous approximation spaces  $V_h := X_h^k \cap H^1(\Omega)$  or discontinuous approximation  $V_h := X_h^k$ . The bilinear form  $a_h(\cdot, \cdot)$  is a discrete realization of  $a(\cdot, \cdot)$ , typically modified to account for the effect of nonconformity, since in general  $V_h \not\subset V$  and  $W_h \not\subset W$ . Weakly imposed boundary conditions may be set in the form  $a_h(\cdot, \cdot)$ , but below we have chosen to impose them using  $s_p(\cdot, \cdot)$  and  $s_a(\cdot, \cdot)$  to obtain a more unified analysis. In (2.7) stabilization can also be added in  $a_h(\cdot, \cdot)$ . Our numerical experiments did not show any advantages of the addition and this approach will not be pursued here.

The bilinear forms  $s_a(\cdot, \cdot)$ ,  $s_p(\cdot, \cdot)$  in (2.7) are symmetric, positive semidefinite stabilization operators, defined on  $[V_h \cup W_h]^2$ . For simplicity we will always assume that  $u$  is sufficiently regular so that strong consistency holds, i.e.,  $s_p(u, v_h)$  is well defined. Note also that for the method to make sense  $s_p(u, v_h)$  must be known, either to be zero, or depending only on known data. This will be the case below. The modifications of the analysis to the case of weakly consistent stabilization are straightforward and not considered here. The seminorm on  $V_h \cup W_h$  associated to the stabilization is defined by

$$|x_h|_{S_y} := s_y(x_h, x_h)^{\frac{1}{2}}, \quad y = a, p.$$

We will assume that the following strong consistency property holds. If  $u$  is the solution of (2.1), then

$$(2.8) \quad a_h(u, \varphi) = (\mathcal{L}u, \varphi) = (f, \varphi) \quad \forall \varphi \in W_h.$$

Then  $u$  solution of (2.1) solves (2.6), and  $u$  solution of (2.1) and  $z \equiv 0$  solve the system (2.7).

We also assume that there are interpolation operators  $\pi_V : V \rightarrow V_h$  and  $\pi_W : W \rightarrow W_h$ , satisfying (2.3). We introduce the (semi)norm  $\|\cdot\|_+$  and assume that the following approximation estimates are satisfied:

$$(2.9) \quad \|v - \pi_V v\|_V + \|v - \pi_V v\|_+ + |v - \pi_V v|_{S_p} \leq c_{a\gamma} h^r |v|_{H^{k+1}(\Omega)} \quad \forall v \in V \cap H^{k+1}(\Omega),$$

where  $r > 0$ , depends on the approximation properties of the finite element space and the definition of the norms in the left-hand side. From the standard error estimates for stabilized methods we expect  $r = k + \frac{1}{2}$  for smooth exact solutions. The constant  $c_{a\gamma}$  depends on the form  $a(\cdot, \cdot)$  and stabilization parameter(s) of the method included in  $s_p(\cdot, \cdot)$  and  $s_a(\cdot, \cdot)$ , here denoted  $\gamma$ .

**2.2. Abstract assumptions on the formulation (2.6).** The assumptions made below constitutes sufficient conditions for the method (2.6) to converge. Here

we assume that  $\|\cdot\|_V \equiv \|\cdot\|_W$ . As usual the conditions are consistency, stability, and continuity of the forms. Galerkin orthogonality for (2.6) is a consequence of the consistency (2.8)

$$(2.10) \quad a_h(u - u_h, w_h) + s_p(u - u_h, w_h) = 0 \quad \forall w_h \in W_h.$$

We assume that there exists  $c_s, c_\eta \in \mathbb{R}^+$  such that for all  $h > 0$  and  $u_h \in V_h$  there exists  $v_a \in W_h$  satisfying

$$(2.11) \quad c_s(\|u_h\|_V^2 + |u_h|_{S_p}^2) \leq a_h(u_h, v_a(u_h)) + s_p(u_h, v_a(u_h)) + \epsilon(h)(\|u_h\|_V^2 + |u_h|_{S_p}^2),$$

where  $\epsilon(h)$  is a continuous function such that  $\epsilon(0) = 0$ , and

$$(2.12) \quad \|v_a(u_h)\|_V + |v_a(u_h)|_{S_p} \leq c_\eta(\|u_h\|_V + |u_h|_{S_p}).$$

These assumptions ensure that the stabilized formulation satisfies a discrete inf-sup condition for  $\epsilon(h)$  small enough. We also assume the following continuity:

$$(2.13) \quad a_h(v - \pi_V v, x_h) \leq \|v - \pi_V v\| + c_a(|x_h|_{S_p} + \|x_h\|_V) \quad \forall v \in V, x_h \in W_h.$$

**2.3. Abstract assumptions on formulation (2.7).** Observe that the following partial coercivity is obtained by taking  $v_h = u_h$  and  $w_h = z_h$  in (2.7):

$$(2.14) \quad |z_h|_{S_a}^2 + |u_h|_{S_p}^2 = (f, z_h) + s_p(u, u_h).$$

The following Galerkin orthogonality holds for (2.7) by (2.8):

$$(2.15) \quad \begin{aligned} a_h(u - u_h, w_h) &= s_a(z_h, w_h) \quad \forall w_h \in W_h, \\ a_h(v_h, z_h) &= s_p(u_h - u, v_h) \quad \forall v_h \in V_h. \end{aligned}$$

Let  $\tilde{\epsilon}(h)$  and  $\check{\epsilon}(h)$  denote continuous, monotonically increasing functions such that  $\tilde{\epsilon}(0) = 0$  and  $0 \leq \check{\epsilon}(h)$ . We assume that the following discrete stability holds for all  $u_h \in V_h, z_h \in W_h$ . For some  $\tilde{c}_s, \tilde{c}_\eta \in \mathbb{R}^+$ , for all  $u_h \in V_h$ , there exists  $v_a(u_h) \in W_h$  such that

$$(2.16) \quad \tilde{c}_s \|u_h\|_V^2 \leq a_h(u_h, v_a(u_h)) + \tilde{\epsilon}(h) \|u_h\|_V^2 + \tilde{c}_\eta |u_h|_{S_p}^2,$$

and similarly, for all  $z_h \in W_h$  there exists  $v_{a^*}(z_h) \in V_h$  such that

$$(2.17) \quad \tilde{c}_s \|z_h\|_W^2 \leq a_h(v_{a^*}(z_h), z_h) + \tilde{\epsilon}(h) \|z_h\|_W^2 + \tilde{c}_\eta |z_h|_{S_a}^2.$$

Moreover assume that the functions  $v_a$  and  $v_{a^*}$  satisfy the bounds

$$(2.18) \quad \|v_a(u_h)\|_W \leq \tilde{c}_\eta \|u_h\|_V, \quad |v_a(u_h)|_{S_a} \leq \check{\epsilon}(h) \|u_h\|_V + \tilde{c}_\eta |u_h|_{S_p},$$

$$(2.19) \quad \|v_{a^*}(z_h)\|_V \leq \tilde{c}_\eta \|z_h\|_W, \quad |v_{a^*}(z_h)|_{S_p} \leq \check{\epsilon}(h) \|z_h\|_W + \tilde{c}_\eta |z_h|_{S_a}.$$

Since we are interested in problems that are ill-conditioned, we here assume  $\tilde{c}_s < \tilde{c}_\eta$  without loss of generality. We finally assume that the following continuity relation holds:

$$(2.20) \quad a_h(v - \pi_V v, x_h) \leq \|v - \pi_V v\| + c_a(|x_h|_{S_a} + \|x_h\|_W) \quad \forall v \in V, x_h \in W_h.$$

**2.4. Convergence analysis for the abstract methods.** We will first prove a convergence result for the standard stabilized finite element method (2.6). Then we will consider (2.7).

PROPOSITION 2.1. *Assume that the solution of (2.1) is smooth and that the forms of (2.6) and the operators  $\pi_V$ ,  $\pi_W$  are such that (2.10)–(2.13) are satisfied. Also assume that  $\epsilon(h)$  satisfies the bound*

$$(2.21) \quad \epsilon(h) \leq \frac{c_s}{2}.$$

Then (2.6) admits a unique solution  $u_h$  for which there holds

$$\|u - u_h\|_V + |u - u_h|_{S_p} \leq c_{as\gamma} h^r |u|_{H^{k+1}(\Omega)},$$

where  $c_{as\gamma} \sim (c_a + 1) \frac{c_\eta}{c_s}$ .

*Proof.* Since the spaces  $W_h$  and  $V_h$  have the same dimension, the matrix is square and it is sufficient to prove uniqueness. Assume  $(f, v_h) + s_p(u, v_h) = 0$  for all  $v_h \in W_h$ . Under condition (2.21) there holds

$$\frac{1}{2} c_s (\|u_h\|_V^2 + |u_h|_{S_p}^2) \leq a_h(u_h, v_a(u_h)) + s_p(u_h, v_a(u_h)) = 0,$$

hence  $u_h = 0$  and existence and uniqueness follows. Let  $\xi_h := \pi_V u - u_h$ . By the stability assumption (2.11) we have

$$c_s (\|\xi_h\|_V^2 + |\xi_h|_{S_p}^2) \leq a_h(\xi_h, v_a(\xi_h)) + s_p(\xi_h, v_a(\xi_h)) + \epsilon(h) (\|\xi_h\|_V^2 + |\xi_h|_{S_p}^2).$$

It follows that under the condition (2.21) there holds

$$\frac{1}{2} c_s (\|\xi_h\|_V^2 + |\xi_h|_{S_p}^2) \leq a_h(\xi_h, v_a(\xi_h)) + s_p(\xi_h, v_a(\xi_h))$$

and by Galerkin orthogonality (2.10), the continuity (2.13), and the stability (2.12)

$$\begin{aligned} \frac{1}{2} c_s (\|\xi_h\|_V^2 + |\xi_h|_{S_p}^2) &\leq a_h(\pi_V u - u, v_a(\xi_h)) + s_p(\pi_V u - u, v_a(\xi_h)) \\ &\leq c_a \|\pi_V u - u\|_+ (|v_a(\xi_h)|_{S_p} + \|v_a(\xi_h)\|_V) + |\pi_V u - u|_{S_p} |v_a(\xi_h)|_{S_p} \\ &\leq (c_a + 1) (\|\pi_V u - u\|_+ + |\pi_V u - u|_{S_p}) c_\eta (\|\xi_h\|_V + |\xi_h|_{S_p}). \end{aligned}$$

We conclude by noting that  $\|u - u_h\|_V \leq \|u - \pi_V u\|_V + \|\xi_h\|_V$  and applying the approximation (2.9).  $\square$

We now turn to the analysis of (2.7). In this case the analysis is based on a combination of coercivity of the stabilization operators (2.14) and an inf-sup argument using (2.16) and (2.17). This allows us to exploit the strong stability property (2.14) enjoyed by the stabilization terms and thereby improve the robustness of the method.

THEOREM 2.2. *Assume that the solution of (2.1) is smooth, that the forms of (2.7) and the operators  $\pi_V$ ,  $\pi_W$  are such that (2.9), (2.15)–(2.20) are satisfied, and that*

$$(2.22) \quad \tilde{\epsilon}(h) \leq \frac{\tilde{c}_s}{2}.$$

Then (2.7) admits a unique solution  $u_h, z_h$  for which there holds

$$\|u - u_h\|_V + \|z_h\|_W + |u - u_h|_{S_p} + |z_h|_{S_a} \leq \tilde{c}_{as\gamma} h^r |u|_{H^{k+1}(\Omega)}.$$

The constant in the above estimate is given by

$$\tilde{c}_{as\gamma} \sim (c_a + 1) \frac{\tilde{c}_\eta}{\tilde{c}_s} \left( 1 + \frac{\check{\epsilon}(h)^2}{\tilde{c}_\eta \tilde{c}_s} \right).$$

Similarly, if  $s_p(u, w_h) = 0$ , there holds

$$|u_h|_{S_p} + |z_h|_{S_a} \leq \tilde{c}_{as\gamma} h^r |u|_{H^{k+1}(\Omega)}.$$

*Proof.* For the first inequality, let  $\xi_h = \pi_V u - u_h$ . As in the previous case it is enough to prove the claim for  $\xi_h$ . By the definition (2.7) there holds

$$|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2 = s_p(\xi_h, \xi_h) + s_a(z_h, z_h) = a_h(\xi_h, z_h) + s_a(z_h, z_h) - a_h(\xi_h, z_h) + s_p(\xi_h, \xi_h).$$

By the stabilities (2.16)–(2.18) there exists  $v_a(\xi_h)$  and  $v_{a^*}(z_h)$  such that

$$\begin{aligned} \tilde{c}_s (\|\xi_h\|_V^2 + \|z_h\|_W^2) &\leq a_h(\xi_h, v_a(\xi_h)) + s_a(v_a(\xi_h), z_h) \\ &\quad + a_h(v_{a^*}(z_h), z_h) - s_p(\xi_h, v_{a^*}(z_h)) + \tilde{\epsilon}(h) \|\xi_h\|_V^2 + \tilde{c}_\eta |\xi_h|_{S_p}^2 \\ &\quad + |z_h|_{S_a} (\check{\epsilon}(h) \|\xi_h\|_V + \tilde{c}_\eta |\xi_h|_{S_p}) + \tilde{\epsilon}(h) \|z_h\|_W^2 + \tilde{c}_\eta |z_h|_{S_a}^2 \\ &\quad + |\xi_h|_{S_p} (\check{\epsilon}(h) \|z_h\|_W + \tilde{c}_\eta |z_h|_{S_a}). \end{aligned}$$

It follows that for all  $\mu_V, \mu_S > 0$  we may write

$$\begin{aligned} \tilde{c}_s \mu_V (\|\xi_h\|_V^2 + \|z_h\|_W^2) + \mu_S (|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2) &\leq a_h(\xi_h, \mu_S z_h + \mu_V v_a(\xi_h)) \\ &\quad + s_a(\mu_S z_h + \mu_V v_a(\xi_h), z_h) - a_h(\mu_S \xi_h - \mu_V v_{a^*}(z_h), z_h) + s_p(\xi_h, \mu_S \xi_h - \mu_V v_{a^*}(z_h)) \\ &\quad + \mu_V \tilde{\epsilon}(h) (\|\xi_h\|_V^2 + \|z_h\|_W^2) + \mu_V \tilde{c}_\eta (|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2) \\ &\quad + \mu_V |z_h|_{S_a} (\check{\epsilon}(h) \|\xi_h\|_V + \tilde{c}_\eta |\xi_h|_{S_p}) + \mu_V |\xi_h|_{S_p} (\check{\epsilon}(h) \|z_h\|_W + \tilde{c}_\eta |z_h|_{S_a}). \end{aligned}$$

By arithmetic-geometric inequalities in the right-hand side

$$\begin{aligned} &\mu_V \tilde{\epsilon}(h) (\|\xi_h\|_V^2 + \|z_h\|_W^2) + \mu_V \tilde{c}_\eta (|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2) \\ &\quad + \mu_V |z_h|_{S_a} (\check{\epsilon}(h) \|\xi_h\|_V + \tilde{c}_\eta |\xi_h|_{S_p}) + \mu_V |\xi_h|_{S_p} (\check{\epsilon}(h) \|z_h\|_W + \tilde{c}_\eta |z_h|_{S_a}) \\ &\leq \mu_V \left( \tilde{\epsilon}(h) + \frac{1}{4} \tilde{c}_s \right) (\|\xi_h\|_V^2 + \|z_h\|_W^2) + \mu_V \left( 2\tilde{c}_\eta + \frac{\check{\epsilon}(h)^2}{\tilde{c}_s} \right) (|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2). \end{aligned}$$

Therefore under the condition (2.22) there holds

$$\begin{aligned} &\frac{1}{4} \tilde{c}_s \mu_V (\|\xi_h\|_V^2 + \|z_h\|_W^2) + \left( \mu_S - \mu_V \left( 2\tilde{c}_\eta + \frac{\check{\epsilon}(h)^2}{\tilde{c}_s} \right) \right) (|\xi_h|_{S_p}^2 + |z_h|_{S_a}^2) \\ &\leq a_h(\xi_h, \mu_S z_h + \mu_V v_a(\xi_h)) + s_a(\mu_S z_h + \mu_V v_a(\xi_h), z_h) \\ &\quad - a_h(\mu_S \xi_h - \mu_V v_{a^*}(z_h), z_h) + s_p(\xi_h, \mu_S \xi_h - \mu_V v_{a^*}(z_h)). \end{aligned}$$

Then, by choosing  $\mu_V = \frac{4}{\tilde{c}_s}$ ,  $\mu_S = \frac{9\tilde{c}_\eta}{\tilde{c}_s} + \frac{4\check{\epsilon}(h)^2}{\tilde{c}_s^2}$  and applying the Galerkin orthogonality of (2.15), we have, since by assumption  $\tilde{c}_s < \tilde{c}_\eta$ ,

$$\begin{aligned} &\|\xi_h\|_V^2 + \|z_h\|_W^2 + |\xi_h|_{S_p}^2 + |z_h|_{S_a}^2 \\ &\leq a_h(\pi_V u - u, \mu_S z_h + \mu_V v_a(\xi_h)) + s_p(\pi_V u - u, \mu_S \xi_h - \mu_V v_{a^*}(z_h)). \end{aligned}$$

We proceed by applying the continuity (2.20) in the first term of the right-hand side and the Cauchy–Schwarz inequality in the stabilization term,

$$\begin{aligned} & \|\xi_h\|_V^2 + \|z_h\|_W^2 + |\xi_h|_{S_p}^2 + |z_h|_{S_a}^2 \\ & \leq \|u - \pi_V u\|_+ c_a (|\mu_S z_h + \mu_V v_a(\xi_h)|_{S_a} + \|\mu_S z_h + \mu_V v_a(\xi_h)\|_W) \\ & \quad + |u - \pi_V u|_{S_p} |\mu_S \xi_h - \mu_V v_{a^*}(z_h)|_{S_p}. \end{aligned}$$

Using a triangle inequality followed by the stability of  $v_a$  (2.18) and  $v_{a^*}$  (2.19) and the bound  $\mu_V(\tilde{c}_\eta + \tilde{\epsilon}(h)) < \mu_S$ , which holds under the assumption  $\tilde{c}_s < \tilde{c}_\eta$ , we may conclude that

$$\begin{aligned} \|\xi_h\|_V^2 + \|z_h\|_W^2 + |\xi_h|_{S_p}^2 + |z_h|_{S_a}^2 & \leq (\|u - \pi_V u\|_+ + |u - \pi_V u|_{S_p}) \\ & \quad \times (c_a + 1) \mu_S (\|\xi_h\|_V + \|z_h\|_W + |\xi_h|_{S_p} + |z_h|_{S_a}). \end{aligned}$$

We conclude from this expression and (2.9) that the first claim holds. The second result is an immediate consequence of  $s_p(u, w_h) = 0$  and the symmetry of  $s_p(\cdot, \cdot)$ . Uniqueness of the discrete solution follows by taking  $f = 0$  in (2.1) and observing that since then  $u = \pi_V u = 0$  we have  $u_h = z_h = 0$  by which uniqueness follows using the same a priori estimates.  $\square$

**3. Stabilization methods.** We let  $\mathcal{L}$  denote the first order hyperbolic operator on nonconservation form,

$$(3.1) \quad \mathcal{L}u := \beta \cdot \nabla u + \sigma u.$$

Here  $\beta \in [W^{2,\infty}(\Omega)]^d$  is a nonsolenoidal velocity vectorfield and  $\sigma \in W^{1,\infty}(\Omega)$ . We assume that boundary conditions are set on the inflow boundary  $\partial\Omega^-$ ,

$$u|_{\partial\Omega^-} = g_{in}, \quad \partial\Omega^\pm := \{x \in \partial\Omega : \pm\beta(x) \cdot n > 0\}.$$

The adjoint operator takes the form

$$(3.2) \quad \mathcal{L}^*u := -\nabla \cdot (\beta u) + \sigma u.$$

We have assumed below that the reaction is moderately stiff so that the relevant time scale of the flow is given by  $h|\beta|^{-1}$ . In particular we will not track the influence of the size of  $\sigma$  in the error bounds below, assuming  $h^{\frac{1}{2}}(\|\sigma\|_{L^\infty(\Omega)} + \|\nabla \cdot \beta\|_{L^\infty(\Omega)})$  moderate. We will consider three different stabilized finite element methods below and show that they all satisfy the assumptions of the abstract theory. The bilinear form  $a_h(\cdot, \cdot)$  of (2.6) and (2.7) is defined as

$$(3.3) \quad a_h(u_h, v_h) := (\mathcal{L}u_h, v_h)_h - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [u_h] \{v_h\} ds,$$

where  $\{v_h\}$  denotes the average of  $v_h$  from the two element faces,

$$\{u_h\}(x)|_{\partial K} := \frac{1}{2} \lim_{\varepsilon \rightarrow 0^+} (u_h(x - \varepsilon n_{\partial K}) + u_h(x + \varepsilon n_{\partial K})),$$

the jump of  $u_h$  is defined as

$$[u_h](x)|_{\partial K} := \lim_{\varepsilon \rightarrow 0^+} (u_h(x - \varepsilon n_{\partial K}) - u_h(x + \varepsilon n_{\partial K})).$$



As usual the jump terms on  $u_h$  may be omitted when a continuous function is considered in the formulation. First we will prove a general stability result on  $a_h(u_h, v_h)$ .

LEMMA 3.1. *For the bilinear form (3.3) there holds for all  $\eta \in W^{1,\infty}(\Omega)$ , for all  $u_h, z_h \in X_h^k$ ,*

$$a_h(u_h, e^{\pm\eta}u_h) = \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) u_h^2 e^{\pm\eta} \, ds + \int_{\Omega} u_h^2 \left( \mp \frac{1}{2} \beta \cdot \nabla \eta - \frac{1}{2} \nabla \cdot \beta + \sigma \right) e^{\pm\eta} \, dx,$$

$$a_h(e^{\pm\eta}z_h, z_h) = \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) z_h^2 e^{\pm\eta} \, ds + \int_{\Omega} z_h^2 \left( \pm \frac{1}{2} \beta \cdot \nabla \eta - \frac{1}{2} \nabla \cdot \beta + \sigma \right) e^{\pm\eta} \, dx.$$

*Proof.* Consider the first inequality with the negative sign in the exponent. By definition we have

$$(3.4) \quad a_h(u_h, e^{-\eta}u_h) = (\beta \cdot \nabla u_h + \sigma u_h, e^{-\eta}u_h)_h - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [u_h] \{e^{-\eta}u_h\} \, ds$$

and note that an integration by parts in the advective term yields

$$\begin{aligned} & (\beta \cdot \nabla u_h, e^{-\eta}u_h)_h - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [u_h] \{e^{-\eta}u_h\} \, ds \\ &= (u_h, e^{-\eta}(\beta \cdot \nabla \eta - \nabla \cdot \beta)u_h) \\ & - (u_h, e^{-\eta}\beta \cdot \nabla u_h)_h + \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [u_h] \{e^{-\eta}u_h\} \, ds \\ & + \int_{\partial\Omega} (\beta \cdot n) u_h^2 e^{-\eta} \, ds. \end{aligned}$$

This equality implies the following well-known relation:

$$(3.5) \quad \begin{aligned} & (\beta \cdot \nabla u_h, e^{-\eta}u_h)_h - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [u_h] \{e^{-\eta}u_h\} \, ds \\ &= \frac{1}{2} \left( (u_h, e^{-\eta}(\beta \cdot \nabla \eta - \nabla \cdot \beta)u_h) + \int_{\partial\Omega} (\beta \cdot n) u_h^2 e^{-\eta} \, ds \right). \end{aligned}$$

The first stability result is obtained by applying this equality in (3.4). The inequality for the adjoint case is proven similarly by observing that after an integration by parts in the bilinear form

$$(3.6) \quad \begin{aligned} a_h(e^{-\eta}z_h, z_h) &= -(e^{-\eta}z_h, \beta \cdot \nabla z_h + (\nabla \cdot \beta - \sigma)z_h) \\ &+ \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \partial\Omega} \beta \cdot n_{\partial K} [z_h] \{e^{-\eta}z_h\} \, ds + \int_{\partial\Omega} (\beta \cdot n) z_h^2 e^{-\eta} \, ds \end{aligned}$$

and then applying (3.5). The case in which the power is positive follows similarly, observing that the change of sign has an effect only in the inner derivative  $\beta \cdot \nabla \eta$ .  $\square$

The importance of this lemma is a consequence of the existence of a particular function  $\eta$  that is given in the following result.

LEMMA 3.2. *Under the assumptions on  $\beta$  there exists  $\eta_0 \in W^{2,\infty}(\Omega)$  such that  $\beta \cdot \nabla \eta_0 \geq 1$  in  $\Omega$ . For the proof of this result see [1, Appendix A].*

It follows that the second term of the right-hand sides in the equations of Lemma 3.1 are nonnegative for

$$(3.7) \quad \eta := (1 + \|2\sigma - \nabla \cdot \beta\|_{L^\infty(\Omega)}) \eta_0.$$

Below we always assume that  $\eta$  is of this form. In general  $e^{-\eta}u_h \notin V_h$  and hence Lemma 3.1 is insufficient to prove (2.16) and (2.17). The trick is to chose  $v_a$  to be some suitable approximation of  $e^{-\eta}u_h$  in  $V_h$ ,  $\pi e^{-\eta}u_h$ , and control the approximation error using the stabilization. Since we are often required to estimate this error we introduce the notation  $\delta(e^{-\eta}u_h) := e^{-\eta}u_h - \pi e^{-\eta}u_h$ . Similarly  $v_{a^*}$  is chosen as an approximation of  $-e^{-\eta}z_h$ .

The stabilization terms may now be chosen as one of the following, where the first two assume  $H^1$ -conforming approximation and the last discontinuous approximation. In all three cases we have  $W_h \equiv V_h$ . Below  $\gamma_X \in \mathbb{R}^+$ ,  $X = GLS, CIP, DG$ , denotes a stabilization parameter associated to the method  $X$  and  $\gamma_{bc} \in \mathbb{R}^+$  a stabilization parameter associated to the weakly imposed boundary condition.

- The GLS method. In this case continuous finite element spaces are used,  $V_h = W_h := X_h^k \cap H^1(\Omega)$ , and the stabilization operators take the form

$$(3.8) \quad s_{p,GLS}(u_h, w_h) := (\gamma_{GLS}|\beta|^{-1}h\mathcal{L}u_h, \mathcal{L}w_h),$$

$$(3.9) \quad s_{a,GLS}(z_h, v_h) := (\gamma_{GLS}|\beta|^{-1}h\mathcal{L}^*z_h, \mathcal{L}^*v_h).$$

Note that  $s_{p,GLS}(u, w_h) = (f, \gamma_{GLS}|\beta|^{-1}h\mathcal{L}w_h)$ , showing that  $s_p(u, \cdot)$  can indeed be expressed using data.

- CIP stabilization. Here as well continuous finite element spaces are used,  $V_h = W_h := X_h^k \cap H^1(\Omega)$ , and the stabilization is given by

$$(3.10) \quad s_{CIP}(u_h, w_h) := \sum_{F \in \mathcal{F}_{int}} \int_F h_F^2 \gamma_{CIP} \|\beta_h \cdot n_F\|_{L^\infty(F)} \llbracket \nabla u_h \rrbracket \cdot \llbracket \nabla w_h \rrbracket \, dx$$

for both the primal and the adjoint equations, where  $\llbracket \nabla u_h \rrbracket|_F$  denotes the jump of the gradient over the face  $F$ .

- The DG method. In this case we do not impose any continuity constraints in the finite element space  $V_h := X_h^k$ . The method is stabilized by penalizing the jump of the solution over element faces for both the primal and the adjoint equations.

$$(3.11) \quad s_{DG}(u_h, w_h) := \sum_{F \in \mathcal{F}_{int}} \int_F \gamma_{DG} |\beta \cdot n_F| [u_h][w_h] \, dx,$$

where  $[u_h]|_F$  denotes the jump of the solution over the face  $F$ . The choice  $\gamma_{DG} = \frac{1}{2}$  is known to lead to the classical upwind formulation for the method (2.6).

To account for boundary conditions the above stabilizations are modified as follows:

$$(3.12) \quad \begin{aligned} s_p(u_h, w_h) &:= s_{p,X}(u_h, w_h) + s_{bc,-}(u_h, w_h), \\ s_a(z_h, v_h) &:= s_{a,X}(z_h, v_h) + s_{bc,+}(z_h, v_h) + s_{bc,-}(z_h, v_h), \end{aligned}$$

with  $X = GLS, CIP, DG$  and  $s_{bc,\pm} := \int_{\partial\Omega} \gamma_{bc} |(\beta \cdot n)_{\pm}| u_h v_h \, ds$ . Note that the value of  $z_h$  is penalized on the whole boundary. This is necessary to obtain robustness if no boundary conditions are set in  $a_h(\cdot, \cdot)$  and allows for the simple choice of test functions used in the analysis below. It should be noted that for problems where the adjoint solution satisfies  $z = 0$  the stabilization in the bulk or on the boundary can be changed to any form satisfying the assumptions (2.17)–(2.20). The consistency requirements are much weaker, since the exact solution is trivial. The variant where  $z_h$  is penalized only on the outflow boundary can also be shown to be stable using the arguments below, provided that weak boundary conditions are included also in  $a_h(\cdot, \cdot)$ . In this case different weight functions must be used for  $u_h$  and  $z_h$ . The present choice was motivated mainly by the use of a single exponential weight in all estimates and that it makes integration of data assimilation problems straightforward by changing the boundary contribution in  $s_p(\cdot, \cdot)$ .

Below we will consider the methods (2.6) and (2.7) one by one, in each case showing that the assumptions (2.10)–(2.13) are satisfied for method (2.6) as well as (2.15)–(2.20) for method (2.7). Clearly some arguments are very similar between the different methods and full details are given only for the GLS method. The conclusion is that all three schemes satisfy the assumptions necessary for the abstract analysis to hold. The dependence of the  $\epsilon(h)$ ,  $\tilde{\epsilon}(h)$ ,  $\check{\epsilon}(h)$  and  $c_\eta$  and  $\tilde{c}_\eta$  on the physical parameters and on  $h$  is specified in each case in the proofs. The natural norm for the analysis is

$$\|x\|_W = \|x\|_V := \|x\| + \|h^{\frac{1}{2}} \beta \cdot \nabla x\|_h + \| |\beta \cdot n|^{\frac{1}{2}} x \|_{\partial\Omega},$$

but to keep down the technical detail we will first prove the results in the reduced norm,

$$(3.13) \quad \|x\|_W = \|x\|_V := \|x\| + \| |\beta \cdot n|^{\frac{1}{2}} x \|_{\partial\Omega},$$

and then show how the control of the streamline derivative can be recovered separately. We also define the continuity norm for all three methods as

$$(3.14) \quad \|v\|_+ := \| (|\beta|^{\frac{1}{2}} h^{-\frac{1}{2}} + |\sigma_\beta|) v \| + \| |\beta \cdot n|^{\frac{1}{2}} v \|_{\mathcal{F}},$$

where  $\sigma_\beta = -\nabla \cdot \beta + \sigma$ . It is straightforward to show that in all cases the approximation estimate (2.9) holds with  $r = k + \frac{1}{2}$  for any interpolant in  $X_h^k$  with optimal approximation properties. The error estimate that results from the abstract analysis for the transport equation may be written in all cases, for both (2.6) and (2.7),

$$\|u - u_h\|_V + \|h^{\frac{1}{2}} \beta \cdot \nabla(u - u_h)\| + |u - u_h|_{S_p} \leq Ch^{k+\frac{1}{2}} |u|_{H^{k+1}(\Omega)}.$$

However, the condition (2.21) leads to a stronger constraint on the mesh for the formulation (2.6) than (2.22). We first prove a lemma, similar to the superapproximation result of [17], useful in all three cases.

LEMMA 3.3. *Let  $\pi$  be an interpolation operator that satisfies (2.3) and (2.4); then there holds*

$$\|e^{-\eta} u_h - \pi e^{-\eta} u_h\|_V + \|e^{-\eta} u_h - \pi e^{-\eta} u_h\|_+ + |e^{-\eta} u_h - \pi e^{-\eta} u_h|_{S_x} \leq \Pi(h) \|u_h\|_V,$$

where  $x = a, p$  and  $\Pi(h) = C_{\gamma\beta\sigma} c_{dc,e^{-\eta}} h^{\frac{1}{2}}$ . Here  $c_{dc,e^{-\eta}}$  refers to the maximum constant of (2.4) for  $n = 0, 1, 2$ . The result holds for all three methods presented above.

*Proof.* First observe that by inequality (2.4) we have

$$(3.15) \quad h^{-\frac{1}{2}}\|\delta(e^{-\eta}u_h)\| + h^{\frac{1}{2}}\|\nabla\delta(e^{-\eta}u_h)\| \leq C \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}} h^{\frac{1}{2}} \|u_h\|,$$

recalling that  $\delta(e^{-\eta}u_h) := e^{-\eta}u_h - \pi e^{-\eta}u_h$ . Similarly using (2.5) followed by (2.4) gives

$$\begin{aligned} \sum_{K \in \mathcal{T}_h} \|\delta(e^{-\eta}u_h)\|_{\partial K}^2 &\leq \sum_{K \in \mathcal{T}_h} c_T^2 (h^{-\frac{1}{2}}\|\delta(e^{-\eta}u_h)\|_K^2 + h^{\frac{1}{2}}\|\nabla\delta(e^{-\eta}u_h)\|_K^2) \\ &\leq C^2 \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}}^2 h \|u_h\|^2. \end{aligned}$$

Using these results in definitions (3.13) and (3.14) we obtain

$$\begin{aligned} \|\delta(e^{-\eta}u_h)\|_+ + \|\delta(e^{-\eta}u_h)\|_V &\leq C(\|\beta\|_{L^\infty}^{\frac{1}{2}} + h^{\frac{1}{2}}\|\sigma_\beta\|_{L^\infty}^{\frac{1}{2}} + h^{\frac{1}{2}}) \|h^{-\frac{1}{2}}\delta(e^{-\eta}u_h)\| \\ &\quad + \|\beta \cdot n\|^{\frac{1}{2}} \|\delta(e^{-\eta}u_h)\|_{\mathcal{F}} \leq C_{\beta\sigma} \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}} h^{\frac{1}{2}} \|u_h\|. \end{aligned}$$

For the stabilization norm we first consider the boundary term and the three methods separately. For the boundary terms we observe that

$$s_{bc,\pm}(\delta(e^{-\eta}u_h), \delta(e^{-\eta}u_h))^{\frac{1}{2}} \leq \gamma_{bc}^{\frac{1}{2}} \|\delta(e^{-\eta}u_h)\|_V \leq C_{\gamma\beta\sigma} \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}} h^{\frac{1}{2}} \|u_h\|.$$

Then note that for the GLS method

$$\begin{aligned} s_{p,GLS}(\delta(e^{-\eta}u_h), \delta(e^{-\eta}u_h))^{\frac{1}{2}} &\leq \gamma_{GLS}^{\frac{1}{2}} h^{\frac{1}{2}} (\|\beta\|_{L^\infty}^{\frac{1}{2}} \|\nabla\delta(e^{-\eta}u_h)\|_h + \|\sigma\|_{L^\infty} \|\delta(e^{-\eta}u_h)\|) \\ &\leq C_{\gamma\beta\sigma} \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}} h^{\frac{1}{2}} \|u_h\| \end{aligned}$$

and similarly  $s_{a,GLS}(\delta(e^{-\eta}u_h), \delta(e^{-\eta}u_h))^{\frac{1}{2}} \leq C_{\gamma\beta\sigma_\beta} \max_{n \in \{0,1\}} c_{dc,n,e^{-\eta}} h^{\frac{1}{2}} \|u_h\|$ .

For the CIP method we use elementwise trace inequalities followed by (2.4), with  $n = 1$  and  $n = 2$ ,

$$\begin{aligned} s_{CIP}(\delta(e^{-\eta}u_h), \delta(e^{-\eta}u_h))^{\frac{1}{2}} &\leq \gamma_{CIP}^{\frac{1}{2}} c_T h^{\frac{1}{2}} \|\beta\|_{L^\infty} \left( \sum_{K \in \mathcal{T}_h} (\|\nabla\delta(e^{-\eta}u_h)\|_K^2 + h^2 \|D^2\delta(e^{-\eta}u_h)\|_K^2) \right)^{\frac{1}{2}} \\ &\leq \gamma_{CIP}^{\frac{1}{2}} c_T h^{\frac{1}{2}} \|\beta\|_{L^\infty} (c_{dc,1,e^{-\eta}} + c_{dc,2,e^{-\eta}}) \|u_h\|. \end{aligned}$$

Finally for the DG method, we simply observe that

$$s_{DG}(\delta(e^{-\eta}u_h), \delta(e^{-\eta}u_h))^{\frac{1}{2}} \leq C_{\gamma DG} \|\delta(e^{-\eta}u_h)\|_+. \quad \square$$

**3.1. GLS stabilization.** We assume that

$$V_h = X_h^k \cap H^1(\Omega), \quad W_h = V_h.$$

Let  $\pi_V, \pi_W$  be defined by the Lagrange interpolator  $i_h$ . It follows by the construction of the stabilization operator and (2.8) that (2.10) and (2.15) hold (recalling that  $z \equiv 0$ .) It is also straightforward to show that (2.9) holds with  $r = k + \frac{1}{2}$ . We collect

the proof of the remaining assumptions of Proposition 2.1 and Theorem 2.2 in two propositions.

PROPOSITION 3.4 (satisfaction of assumptions for (2.6) with GLS). *Let the bilinear forms of (2.6) be defined by (3.3) and (3.8) with  $\gamma_{bc} \geq 1$ . Then (2.11)–(2.13) are satisfied, with  $\epsilon(h) = C_{\gamma\beta\sigma\eta}h^{\frac{1}{2}}$ .*

*Proof.* To show (2.11) we take  $v_a := \pi_V(e^{-\eta}u_h)$  with  $\eta$  defined by (3.7) and use the first inequality of Lemma 3.1 to obtain

$$(3.16) \quad \begin{aligned} a_h(u_h, \pi_V(e^{-\eta}u_h)) &= a_h(u_h, e^{-\eta}u_h) - a_h(u_h, \delta(e^{-\eta}u_h)) \\ &\geq -\gamma_{GLS}^{-\frac{1}{2}}|u_h|_{S_p} \|\delta(e^{-\eta}u_h)\|_+ \\ &\quad + \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n)_+ u_h^2 e^{-\eta} \, ds + \frac{1}{2} \|u_h e^{-\frac{\eta}{2}}\|^2. \end{aligned}$$

Using Lemma 3.3 we have

$$(3.17) \quad \begin{aligned} \frac{1}{2} \|u_h e^{-\frac{\eta}{2}}\|^2 + \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n)_+ u_h^2 e^{-\eta} \, ds &\leq a_h(u_h, \pi_V(e^{-\eta}u_h)) \\ &\quad - \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n)_- u_h^2 e^{-\eta} + \frac{1}{2} \gamma_{GLS}^{-\frac{1}{2}} \Pi(h) (|u_h|_{S_p}^2 + \|u_h\|_V^2). \end{aligned}$$

We need a similar bound for the stabilization operator using the function  $v_a(u_h)$ . This is straightforward observing that

$$\begin{aligned} s_p(u_h, v_a(u_h)) &= (\mathcal{L}u_h, \gamma_{GLS}|\beta|^{-1}h\mathcal{L}(u_h e^{-\eta})) + s_{bc,-}(u_h, e^{-\eta}u_h) \\ &\quad - (\mathcal{L}u_h, \gamma_{GLS}|\beta|^{-1}h\mathcal{L}\delta(e^{-\eta}u_h)) - s_{bc,-}(u_h, \delta(e^{-\eta}u_h)) \\ &\geq \|(\gamma_{GLS}h|\beta|^{-1})^{\frac{1}{2}}\mathcal{L}u_h e^{-\frac{\eta}{2}}\|^2 + \gamma_{bc} \|(\beta \cdot n)_-|^{\frac{1}{2}}u_h e^{-\frac{\eta}{2}}\|_{\partial\Omega}^2 \\ &\quad - |u_h|_{S_p} \left( \|\delta(e^{-\eta}u_h)\|_{S_p} + \|(\gamma_{GLS}h|\beta|^{-1})^{\frac{1}{2}}(\mathcal{L}e^{-\frac{\eta}{2}})u_h\| \right). \end{aligned}$$

Combining this result with (3.17), using (3.3) it follows that for  $\gamma_{bc}$  large enough

$$(3.18) \quad \begin{aligned} \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} (\|u_h\|_V^2 + |u_h|_{S_p}^2) &\leq a_h(u_h, \pi_V(e^{-\eta}u_h)) + s_{p, GLS}(u_h, \pi_V(e^{-\eta}u_h)) \\ &\quad + (C_\gamma \Pi(h) + (\gamma_{GLS}h|\beta|^{-1})^{\frac{1}{2}} \sup_{x \in \Omega} |\mathcal{L}e^{-\frac{\eta}{2}}|) (|u_h|_{S_p}^2 + \|u_h\|_V^2). \end{aligned}$$

We conclude that (2.11) holds with  $c_s = \frac{1}{2} \inf_{x \in \Omega} e^{-\eta}$  and

$$\epsilon(h) = (C_\gamma \Pi(h) + (\gamma_{GLS}h|\beta|^{-1})^{\frac{1}{2}} \sup_{x \in \Omega} |\mathcal{L}e^{-\frac{\eta}{2}}|) \sim C_{\gamma\beta\sigma\eta}h^{\frac{1}{2}}.$$

Considering now (2.12) we have

$$(3.19) \quad \|v_a(u_h)\|_V \leq \|e^{-\eta}u_h\|_V + \|\delta(e^{-\eta}u_h)\|_V \leq (\sup_{x \in \Omega} e^{-\eta} + \Pi(h)) \|u_h\|_V$$

and for the stabilization part,

$$(3.20) \quad \begin{aligned} |v_a(u_h)|_{S_p} &\leq \sup_{x \in \Omega} |\mathcal{L}e^{-\frac{\eta}{2}}| h^{\frac{1}{2}} C_\gamma \|u_h\|_V + \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_p} + |\delta(e^{-\eta}u_h)|_{S_p} \\ &\leq (\sup_{x \in \Omega} e^{-\eta} h^{\frac{1}{2}} C_{\gamma\beta\sigma\eta} + \Pi(h)) \|u_h\|_V + \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_p}. \end{aligned}$$

It follows that (2.12) holds for any

$$c_\eta \geq \max(\sup_{x \in \Omega} |\mathcal{L}e^{-\frac{\eta}{2}}| h^{\frac{1}{2}} C_{\gamma\beta\sigma}, \sup_{x \in \Omega} e^{-\eta} + \Pi(h)) \sim C_{\gamma\beta\sigma\eta} h^{\frac{1}{2}} + \sup_{x \in \Omega} e^{-\eta}.$$

For the continuity (2.13) we first use an integration by parts and the Cauchy–Schwarz inequality to obtain

$$(3.21) \quad \begin{aligned} a_h(v - \pi_V v, x_h) &= (u - \pi_V u, \mathcal{L}^* x_h) + \int_{\partial\Omega} (\beta \cdot n)(v - \pi_V v)x_h \, ds \\ &\leq \|u - \pi_V u\|_+ (|\beta|^{-1}h)^{\frac{1}{2}} \|\mathcal{L}^* x_h\| + \|x_h\|_V. \end{aligned}$$

To conclude we need to express the norm over the adjoint operator in the right-hand side by the stabilization of the primal operator. Observe that for all  $x_h \in V_h$  there holds

$$(3.22) \quad \|(|\beta|^{-1}h)^{\frac{1}{2}} \mathcal{L}^* x_h\| \leq |x_h|_{S_p} + C_{\gamma\beta} h^{\frac{1}{2}} (2\|\sigma\|_{L^\infty(\Omega)} + \|\nabla \cdot \beta\|_{L^\infty(\Omega)}) \|x_h\|_V.$$

Collecting the results of (3.21) and (3.22) we see that  $c_a \geq 1 + C_{\gamma\beta\sigma} h^{\frac{1}{2}}$ .  $\square$

PROPOSITION 3.5 (satisfaction of the assumptions for (2.7) with GLS). *Let the bilinear forms of (2.7) be defined by (3.3), (3.8), and (3.9). Then the inequalities (2.15)–(2.20) hold with  $\tilde{\epsilon}(h) = 0$ .*

*Proof.* Starting from the inequality (3.16) with  $v_a(u_h) := \pi_W(e^{-\eta}u_h)$  we immediately get

$$\begin{aligned} \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} \|u_h\|_V^2 &\leq a_h(u_h, \pi_W(e^{-\eta}u_h)) + \gamma_{GLS}^{-\frac{1}{2}} \Pi(h) |u_h|_{S_p} \|u_h\|_V \\ &\quad + \sup_{x \in \Omega} e^{-\eta} \gamma_{bc}^{-1} |u_h|_{S_p}^2, \end{aligned}$$

from which we deduce, using  $(\inf_{x \in \Omega} e^{-\eta})^{-1} = \sup_{x \in \Omega} e^\eta$ ,

$$\frac{1}{4} \inf_{x \in \Omega} e^{-\eta} \|u_h\|_V^2 \leq a_h(u_h, \pi_W(e^{-\eta}u_h)) + (\sup_{x \in \Omega} e^\eta \gamma_{GLS}^{-1} \Pi(h)^2 + \sup_{x \in \Omega} e^{-\eta} \gamma_{bc}^{-1}) |u_h|_{S_p}^2,$$

which is the required inequality with  $\tilde{\epsilon}(h) = 0$ ,  $\tilde{c}_s = \frac{1}{4} \inf_{x \in \Omega} e^{-\eta}$ , and

$$\tilde{c}_\eta \geq \sup_{x \in \Omega} e^\eta \gamma_{GLS}^{-1} \Pi(h)^2 + \sup_{x \in \Omega} e^{-\eta} \gamma_{bc}^{-1}.$$

In a similar fashion we may show that (2.17) holds, also with the weight  $e^{-\eta}$ , and corresponding test function  $v_{a^*}(z_h) = -\pi_V(e^{-\eta}z_h)$ . First observe that in this case using Lemma 3.1 (second equation),

$$\begin{aligned} \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} \|z_h\|^2 &- \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) z_h^2 e^{-\eta} \, ds \leq -a_h(e^{-\eta}z_h, z_h) \\ &= a_h(-\pi_V(e^{-\eta}z_h), z_h) - a_h(\delta(e^{-\eta}z_h), z_h). \end{aligned}$$

For the second term in the right-hand side we have after integration by parts and application of Lemma 3.3

$$\begin{aligned} a_h(\delta(e^{-\eta}z_h), z_h) &= \int_{\partial\Omega} (\beta \cdot n) \delta(e^{-\eta}z_h) z_h \, ds + (\delta(e^{-\eta}z_h), \mathcal{L}^* z_h) \\ &\leq C \|\delta(e^{-\eta}z_h)\|_+ |z_h|_{S_a} \leq C \Pi(h) \|z_h\|_V |z_h|_{S_a}. \end{aligned}$$

Here we used that the boundary penalty on  $z_h$  is active on the whole boundary. We may then conclude as before that

$$\frac{1}{4} \inf_{x \in \Omega} e^{-\eta} \|z_h\|_W^2 \leq a_h(-\pi_V(e^{-\eta} z_h), z_h) + (\sup_{x \in \Omega} e^\eta C_\gamma \Pi(h))^2 + \sup_{x \in \Omega} e^{-\eta} \gamma_{bc}^{-1} |z_h|_{S_a}^2$$

with similar constants as before.

The inequalities of (2.18) and (2.19) follow by similar arguments as (3.19) and (3.20). The only differences occur in the right inequalities.

$$\begin{aligned} |v_a(u_h)|_{S_a} &\leq h^{\frac{1}{2}} \gamma_{GLS}^{\frac{1}{2}} \sup_{x \in \Omega} \mathcal{L}^* e^{-\eta} \|u_h\|_V + \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_a} + |\delta(e^{-\eta} u_h)|_{S_a} \\ &\leq (C_{\gamma\beta\sigma\eta} h^{\frac{1}{2}} \sup_{x \in \Omega} e^{-\eta} + \Pi(h)) \|u_h\|_V + \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_a}. \end{aligned}$$

We then use an inequality similar to (3.22), this time adding the boundary penalty term that is included in the stabilization in formulation (2.7) (see (3.12)):

(3.23)

$$|u_h|_{S_a} \leq |u_h|_{S_p} + C_{\beta\gamma} h^{\frac{1}{2}} (2\|\sigma\|_{L^\infty(\Omega)} + \|\nabla \cdot \beta\|_{L^\infty(\Omega)}) \|u_h\|_V + \gamma_{bc}^{\frac{1}{2}} \|\beta \cdot n\|_{\frac{1}{2}} \|u_h\|_{\partial\Omega}.$$

Note that the boundary contribution cannot be controlled by  $|u_h|_{S_p}$  as one would like but must be controlled using the  $V$ -norm. This adds an  $O(\gamma_{bc}^{\frac{1}{2}})$  contribution to the constant in front of  $\|u_h\|_V$ :

$$(3.24) \quad |u_h|_{S_a} \leq |u_h|_{S_p} + (\gamma_{bc}^{\frac{1}{2}} + C_{\beta\gamma} h^{\frac{1}{2}} (2\|\sigma\|_{L^\infty(\Omega)} + \|\nabla \cdot \beta\|_{L^\infty(\Omega)})) \|u_h\|_V.$$

The proof of (2.19) is similar, but here the stronger adjoint boundary penalty can control the boundary term, leading to

$$|z_h|_{S_p} \leq |z_h|_{S_a} + C_{\beta\gamma} h^{\frac{1}{2}} (2\|\sigma\|_{L^\infty(\Omega)} + \|\nabla \cdot \beta\|_{L^\infty(\Omega)}) \|z_h\|_W.$$

We conclude that the inequalities (2.18) and (2.19) hold with

$$\tilde{c}_\eta \geq \sup_{x \in \Omega} e^{-\eta} + \Pi(h) \text{ and } \check{\epsilon}(h) \geq C_{\beta\sigma\gamma\eta} h^{\frac{1}{2}} + \sup_{x \in \Omega} e^{-\eta} \gamma_{bc}^{\frac{1}{2}}.$$

The continuity (2.20) is immediate by integration by parts and the Cauchy–Schwarz inequality,

$$\begin{aligned} a_h(v - \pi_V v, x_h) &= (u - \pi_V u, \mathcal{L}^* x_h) + \int_{\partial\Omega} (\beta \cdot n)(v - \pi_V v) x_h \, ds \\ &\leq C_\gamma \|u - \pi_V u\|_+ (|x_h|_{S_a} + \|x_h\|_W). \quad \square \end{aligned}$$

*Remark 2.* Note that for the GLS method  $\tilde{\epsilon}(h) = 0$  in (2.16) and (2.17), indicating that the scheme is unconditionally stable. This follows from the fact that the whole residual is considered in the stabilization term. This nice feature, however, only holds under exact quadrature. When the integrals are approximated, the quadrature error once again gives rise to oscillation terms from data that introduces a nonzero contribution to  $\tilde{\epsilon}(h)$ .

**3.2. CIP.** In this case also  $W_h = V_h := X_h^k \cap H^1(\Omega)$ , but the stabilization added to the standard Galerkin formulation is a penalty on the jump of the gradient over element faces [12, 9]. The key observation is that the following discrete approximation result holds for  $\gamma_{CIP}$  large enough (see [7, 8]):

$$(3.25) \quad \|h^{\frac{1}{2}}|\beta_h|^{-\frac{1}{2}}(\beta_h \cdot \nabla u_h - I_{os}\beta_h \cdot \nabla u_h)\|^2 \leq s_{CIP}(u_h, u_h).$$

Here  $\beta_h$  is some piecewise affine interpolant of the velocity vector field  $\beta$  and  $I_{os}$  is the quasi-interpolation operator defined in each node of the mesh as a straight average of the function values from triangles sharing that node,

$$(I_{os}\beta_h \cdot \nabla u_h)(x_i) = N_i^{-1} \sum_{\{K: x_i \in K\}} (\beta_h \cdot \nabla u_h)(x_i)|_K,$$

with  $N_i := \text{card}\{K : x_i \in K\}$ . Stability is then a consequence of the following lemma.

LEMMA 3.6. *The following inequalities hold:*

$$(3.26) \quad \inf_{v_h \in V_h} \|h^{\frac{1}{2}}(\mathcal{L}u_h - v_h)\| \leq C_{\gamma\beta} s_{CIP}(u_h, u_h)^{\frac{1}{2}} + \epsilon_{CIP}(h)\|u_h\|$$

and

$$(3.27) \quad \inf_{w_h \in W_h} \|h^{\frac{1}{2}}(\mathcal{L}^*z_h - w_h)\| \leq C_{\gamma\beta} s_{CIP}(z_h, z_h)^{\frac{1}{2}} + \epsilon_{CIP}(h)\|z_h\|$$

with  $\epsilon_{CIP}(h) \sim h^{\frac{3}{2}}(\|\beta\|_{W^{2,\infty}(\Omega)} + c_{dc,0,\sigma})$ .

*Proof.* Since the proofs of the two results are similar we only detail the arguments for (3.26). First note that

$$\begin{aligned} \inf_{v_h \in V_h} \|h^{\frac{1}{2}}(\mathcal{L}u_h - v_h)\| &\leq \|h^{\frac{1}{2}}(i_h\beta \cdot \nabla u_h - I_{os}(i_h\beta \cdot \nabla u_h))\| \\ &\quad + h^{\frac{1}{2}}\|\beta - i_h\beta\|_{L^\infty(\Omega)}\|\nabla u_h\| + h^{\frac{1}{2}}\|\sigma u_h - i_h(\sigma u_h)\|. \end{aligned}$$

Using (3.25), interpolation in  $L^\infty$ , an inverse inequality, and the discrete commutator property (2.4) we conclude

$$\inf_{v_h \in V_h} \|h^{\frac{1}{2}}(\mathcal{L}u_h - v_h)\| \leq C_{\beta\gamma} s_{CIP}(u_h, u_h)^{\frac{1}{2}} + h^{\frac{3}{2}}(\|\beta\|_{W^{2,\infty}(\Omega)} + c_{dc,0,\sigma})\|u_h\|. \quad \square$$

For the CIP method we choose the  $\pi_V$  and  $\pi_W$  as the  $L^2$ -projection in order to exploit orthogonality to “filter” the element residual. Observe that if  $u \in H^{\frac{3}{2}+\epsilon}(\Omega)$ ,  $\epsilon > 0$ , then  $s_{CIP}(u, \cdot) = 0$ . The consistencies (2.10) and (2.15) hold from the consistency of (3.3). The approximation result (2.9), with  $r = k + \frac{1}{2}$  is a consequence of standard results for the CIP method (see, for instance, [8].) We now prove that the remaining assumptions for Proposition 2.1 and Theorem 2.2 hold.

PROPOSITION 3.7 (satisfaction of assumptions for (2.6) with CIP). *Let the bilinear forms of (2.6) be defined by (3.3) and (3.10). Let  $\gamma_{bc} \geq 1$ . Then (2.10)–(2.13) are satisfied, with  $\epsilon(h) \sim h^{\frac{1}{2}}$ .*

*Proof.* To prove the stability (2.11) take  $v_a = \pi_V(e^{-\eta}u_h)$  and use lemma 3.1, the orthogonality of the  $L^2$ -projection, and lemma 3.6 to obtain

$$\begin{aligned} (3.28) \quad a_h(u_h, \pi_V(e^{-\eta}u_h)) &= a_h(u_h, e^{-\eta}u_h) - (\mathcal{L}u_h - w_h, \delta(e^{-\eta}u_h)) \\ &\geq -C_\gamma|u_h|_{S_p}\|\delta(e^{-\eta}u_h)\|_+ - \epsilon_{CIP}(h)\|u_h\|\|h^{-\frac{1}{2}}\delta(e^{-\eta}u_h)\| \\ &\quad + \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n)u_h^2 e^{-\eta} \, ds + \frac{1}{2}\|u_h e^{-\frac{\eta}{2}}\|^2. \end{aligned}$$



We also observe that for the stabilization

$$(3.29) \quad s_p(u_h, \pi_V(e^{-\eta}u_h)) \geq s_p(u_h, u_h e^{-\eta}) - |u_h|_{S_p} |\delta(e^{-\eta}u_h)|_{S_p}.$$

Now observe that since the jump of  $\nabla e^{-\eta}$  is zero we have, using (3.28) and (3.29),

$$\begin{aligned} \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} (\|u_h\|_V^2 + |u_h|_{S_p}^2) &\leq \frac{1}{2} \|u_h e^{-\frac{\eta}{2}}\|^2 + \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) u_h^2 e^{-\eta} \, ds + s_p(u_h, u_h e^{-\eta}) \\ &\leq a_h(u_h, \pi_V(e^{-\eta}u_h)) + s_p(u_h, \pi_V(e^{-\eta}u_h)) \\ &\quad + |u_h|_{S_p} (C_\gamma \|\delta(e^{-\eta}u_h)\|_+ + |\delta(e^{-\eta}u_h)|_{S_p}) + \epsilon_{CIP}(h) \|u_h\| \|h^{-\frac{1}{2}} \delta(e^{-\eta}u_h)\|. \end{aligned}$$

Using Lemma 3.3 we deduce that (2.11) holds with

$$c_s = \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} \text{ and } \epsilon(h) \geq \Pi(h)(C_\gamma + \epsilon_{CIP}(h)).$$

For (2.12) only the stabilization part differs from the GLS case. Since the jump of  $\nabla e^{-\eta}$  is zero we immediately get

$$|v_a(u_h)|_{S_p} \leq \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_p} + |\delta(e^{-\eta}u_h)|_{S_p} \leq \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_p} + \Pi(h) \|u\|_V$$

and hence  $c_\eta \geq \sup_{x \in \Omega} e^{-\eta} + \Pi(h)$ . The continuity (2.13) follows by observing that by (3.6) there holds

$$(3.30) \quad \begin{aligned} a_h(v - \pi_V v, x_h) &= \inf_{w_h \in V_h} (v - \pi_V v, \mathcal{L}^* x_h - w_h) + \int_{\partial\Omega} (\beta \cdot n)(v - \pi_V v) x_h \, ds \\ &\leq \|v - \pi_V v\|_+ (C_\gamma |x_h|_{S_p} + (C_\beta h^{\frac{1}{2}} \epsilon_{CIP}(h) + 1) \|x_h\|_V), \end{aligned}$$

where we observe that the boundary part must be controlled using the norm  $\|\cdot\|_V$ .  $\square$

**PROPOSITION 3.8** (satisfaction of assumptions for (2.7) with CIP). *Let the bilinear forms of (2.7) be defined by (3.3) and (3.10) for both  $s_p(\cdot, \cdot)$  and  $s_a(\cdot, \cdot)$ , together with the respective boundary penalty terms of (3.12). Then the inequalities (2.15)–(2.20) hold with  $\tilde{\epsilon}(h) \sim h^2$ .*

*Proof.* Starting from (3.28) with  $v_a(u_h) := \pi_W(e^{-\eta}u_h)$  we have using Lemma 3.3,

$$(3.31) \quad \begin{aligned} \frac{1}{2} \|u_h e^{-\frac{\eta}{2}}\|^2 + \frac{1}{2} \int_{\partial\Omega} |\beta \cdot n| u_h^2 e^{-\eta} \, ds &\leq a_h(u_h, \pi_W(e^{-\eta}u_h)) - \int_{\partial\Omega} (\beta \cdot n)_- u_h^2 e^{-\eta} \, ds \\ &\quad + (C_\gamma |u_h|_{S_p} + \epsilon_{CIP}(h) \|u_h\|) \Pi(h) \|u_h\| \\ &\leq a_h(u_h, \pi_W(e^{-\eta}u_h)) + (\gamma_{bc}^{-\frac{1}{2}} \sup_{x \in \Omega} e^{-\eta} + C_\gamma^2 \sup_{x \in \Omega} e^\eta \Pi(h)^2) |u_h|_{S_p} \\ &\quad + \left( \frac{1}{4} \inf_{x \in \Omega} e^{-\eta} + \epsilon_{CIP}(h) \Pi(h) \right) \|u_h\|^2. \end{aligned}$$

The last inequality is due to an arithmetic-geometric inequality. Hence we see that (2.16) holds with  $\tilde{\epsilon}(h) = \epsilon_{CIP}(h) \Pi(h) \sim h^2$  and

$$\tilde{c}_s = \frac{1}{4} \inf_{x \in \Omega} e^{-\eta}, \quad \tilde{c}_\eta \geq C_\gamma^2 \sup_{x \in \Omega} e^\eta \Pi(h)^2 + \gamma_{bc}^{-1} \sup_{x \in \Omega} e^{-\eta}.$$

The inequality (2.17) is proved similarly as in the GLS case, taking this time  $v_{a^*}(z_h) := -\pi_V(e^{-\eta}z_h)$ , with  $\pi_V$  the  $L^2$ -projection and using the second inequality of Lemma 3.1 and Lemma 3.3 after integration by parts:

$$\begin{aligned}
 a_h(\delta(e^{-\eta} z_h), z_h) &= \int_{\partial\Omega} (\beta \cdot n) \delta(e^{-\eta} z_h) z_h \, ds + \inf_{w_h \in W_h} (\delta(e^{-\eta} z_h), \mathcal{L}^* z_h - w_h) \\
 &\leq C \|\delta(e^{-\eta} z_h)\|_+ (|z_h|_{S_a} + \epsilon_{CIP}(h) \|z_h\|) \leq C \Pi(h) \|z_h\| (|z_h|_{S_a} + \epsilon_{CIP}(h) \|z_h\|).
 \end{aligned}$$

Then we conclude as before. For the stabilities (2.18) and (2.19) we proceed as in Proposition 3.4 and we only detail the second inequality of (2.18). When using the CIP method the primal and adjoint stabilization terms differ only in the boundary contributions; therefore, by symmetry, the second inequality of (2.19) follows identically. Since the jump of  $\nabla e^{-\eta}$  is zero we get

$$(3.32) \quad |v_a(u_h)|_{S_a} \leq \sup_{x \in \Omega} e^{-\eta} |u_h|_{S_p} + |\delta(e^{-\eta} u_h)|_{S_p} + \gamma_{bc}^{\frac{1}{2}} \|\beta \cdot n\|_{\frac{1}{2}} |v_a(u_h)|_{\partial\Omega}.$$

The boundary term is controlled by adding and subtracting  $e^{-\eta} u_h$  and then applying a triangle inequality followed by Lemma 3.3, leading to

$$\begin{aligned}
 (3.33) \quad \gamma_{bc}^{\frac{1}{2}} \|\beta \cdot n\|_{\frac{1}{2}} |v_a(u_h)|_{\partial\Omega} &\leq \Pi(h) \|u_h\|_V + \gamma_{bc}^{\frac{1}{2}} \|\beta\|_{\frac{1}{2}} |u_h e^{-\eta}|_{\partial\Omega} \\
 &\leq (\Pi(h) + \gamma_{bc}^{\frac{1}{2}} \sup_{x \in \Omega} e^{-\eta}) \|u_h\|_V.
 \end{aligned}$$

Therefore (2.18) and (2.19) hold with

$$(3.34) \quad \tilde{c}_\eta \geq \sup_{x \in \Omega} e^{-\eta} + \Pi(h) \text{ and } \check{\epsilon}(h) \geq \gamma_{bc}^{\frac{1}{2}} \sup_{x \in \Omega} e^{-\eta} + 2\Pi(h).$$

The proof of continuity (2.20) follows as in (3.30).  $\square$

**3.3. The discontinuous Galerkin method.** In the case where discontinuous elements are used, i.e.,  $V_h = W_h := X_h^k$ , the analysis is simplified by the fact that  $\beta_h \cdot \nabla u_h \in V_h$ . Here we let  $\pi_V$  and  $\pi_W$  denote the elementwise  $L^2$ -projection onto  $X_h^k$ . The analysis is essentially the same as for the CIP method and when appropriate we will refer to the previous analysis. Thanks to the local character of the DG method the results hold without assuming any quasi regularity of the meshes. The consistency results (2.10) and (2.15) are standard, as well as the approximation result (2.9), with  $r = k + \frac{1}{2}$  (see [13]). As before we collect the proofs of the remaining assumption in a proposition.

**PROPOSITION 3.9** (satisfaction of assumptions for (2.6) with DG). *Let the bilinear forms of (2.6) be defined by (3.3) and (3.11). Then (2.10)–(2.13) are satisfied with  $\epsilon(h) \sim h^{\frac{1}{2}}$ .*

*Proof.* Let  $i_h \beta \in X_h^1$  be the Lagrange interpolant of  $\beta$  with and  $\pi_0 \sigma \in X_h^0$  the projection of  $\sigma$  on piecewise constant functions. For (2.11) take  $v_a := \pi_V(e^{-\eta} u_h)$ , use  $L^2$ -orthogonality, and apply Lemma 3.1 to obtain for  $\gamma_{bc}$  large enough,

$$\begin{aligned}
 (3.35) \quad &a_h(u_h, \pi_V(e^{-\eta} u_h)) + s_p(u_h, \pi_V(e^{-\eta} u_h)) = a_h(u_h, e^{-\eta} u_h) + s_p(u_h, e^{-\eta} u_h) \\
 &\quad - a_h(u_h, \delta(e^{-\eta} u_h)) - s_p(u_h, \delta(e^{-\eta} u_h)) \\
 &\geq ((i_h \beta - \beta) \nabla u_h + (\pi_0 \sigma - \sigma) u_h, \delta(e^{-\eta} u_h)) \\
 &\quad - 2 \sum_{K \in \mathcal{T}_h} \langle |\beta \cdot n| [u_h], (1 + \gamma_{DG}) |\delta(e^{-\eta} u_h)| \rangle_{\partial K \setminus \partial\Omega} + \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} (\|u_h\|_V^2 + |u_h|_{S_p}^2) \\
 &\geq -|u_h|_{S_p} C_\gamma \|\delta(e^{-\eta} u_h)\|_+ - \epsilon_{DG}(h) \|u_h\| \|h^{-\frac{1}{2}} \delta(e^{-\eta} u_h)\| \\
 &\quad + \frac{1}{2} \inf_{x \in \Omega} e^{-\eta} (\|u_h\|_V^2 + |u_h|_{S_p}^2),
 \end{aligned}$$

where

$$\epsilon_{DG}(h) = h^{-\frac{1}{2}} \|i_h \beta - \beta\|_{L^\infty(\Omega)} + h^{\frac{1}{2}} \|\pi_0 \sigma - \sigma\|_{L^\infty(\Omega)} \sim (\|\beta\|_{W^{2,\infty}(\Omega)} + \|\sigma\|_{W^{1,\infty}(\Omega)}) h^{\frac{3}{2}}.$$

It follows that (2.11) holds with  $c_s = \frac{1}{2} \inf_{x \in \Omega} e^{-\eta}$  and  $\epsilon(h) \geq (C_\gamma + \epsilon_{DG}(h)) \Pi(h)$ . The proof of (2.12) is analogous with the CIP case with similar constants. Considering finally the continuity (2.13) we have after an integration by parts

(3.36)

$$\begin{aligned} a_h(v - \pi_V v, x_h) &= (v - \pi_V v, \mathcal{L}^* x_h) + \frac{1}{2} \sum_K \langle (\beta \cdot n) \{v - \pi_V v\}, [x_h] \rangle_{\partial K \setminus \partial \Omega} \\ &\quad + \langle (\beta \cdot n)(v - \pi_V v), x_h \rangle_{\partial \Omega} \\ &= (v - \pi_V v, (i_h \beta - \beta) \nabla x_h + (\sigma - \pi_0 \sigma) x_h) + \frac{1}{2} \sum_K \langle (\beta \cdot n) \{v - \pi_V v\}, [x_h] \rangle_{\partial K \setminus \partial \Omega} \\ &\quad + \langle (\beta \cdot n)(v - \pi_V v), x_h \rangle_{\partial \Omega} \\ &\leq \|v - \pi_V v\|_+ (C_\gamma |x_h|_{S_a} + (C_\beta h^{\frac{1}{2}} \epsilon_{DG}(h) + 1) \|x_h\|_V). \quad \square \end{aligned}$$

**PROPOSITION 3.10** (satisfaction of assumptions for (2.7) with DG). *Let the bilinear forms of (2.7) be defined by (3.3) and (3.11) for both  $s_p(\cdot, \cdot)$  and  $s_a(\cdot, \cdot)$  together with the respective boundary penalty terms of (3.12). Then the inequalities (2.15)–(2.20) hold with  $\tilde{\epsilon}(h) \sim h^2$ .*

*Proof.* The stability (2.16) and (2.17) follows by taking  $v_a := \pi_W(e^{-\eta} u_h)$  and  $v_{a*} := -\pi_V(e^{-\eta} z_h)$ , using (3.35) and the manipulations of Proposition 3.8. The proof of the inequalities (2.18) and (2.19) uses the same techniques as the corresponding results for the CIP-method and results in similar constants. Finally (2.20) follows from (3.36).  $\square$

**3.4. Convergence of the error in the streamline derivative.** As mentioned the natural norm for the above analysis would include the  $L^2$ -norm of the  $h^{\frac{1}{2}}$ -weighted streamline derivative. Given the results of the previous section it is straightforward to prove optimal convergence of the streamline derivative for both (2.6) and (2.7). We only give the result for the method (2.7) below. The proof of the result for (2.6) is identical.

**PROPOSITION 3.11.** *Let  $u_h, z_h$  be the solution of (2.7) with bilinear form (3.3) stabilized with one of the methods presented in sections 3.1–3.3. Assume that the conditions of Theorem 2.2 are satisfied. Then there holds*

$$\|\beta \cdot \nabla(u - u_h)\|_h \leq C_{\eta\beta\sigma\gamma} h^k |u|_{H^{k+1}(\Omega)}.$$

*Proof.* First consider the GLS method. Add and subtract  $\sigma(u - u_h)$  inside the streamline derivative norm and use a triangle inequality to obtain, using the previously obtained error estimates,

$$\|\beta \cdot \nabla(u - u_h)\| \leq C_\gamma \|\beta\|_{L^\infty}^{-\frac{1}{2}} h^{-\frac{1}{2}} (|u - u_h|_{S_p} + \|\sigma\|_{L^\infty(\Omega)} h^{\frac{1}{2}} \|u - u_h\|) \leq C_{\gamma\beta\sigma} h^k |u|_{H^{k+1}(\Omega)}.$$

For the CIP method we may write  $\xi_h := \pi_V u - u_h$ , where  $\pi_V$  is any interpolation operator with optimal approximation properties, and note that by Galerkin orthogonality, interpolation in  $L^\infty$ , and inverse inequalities, we have

$$\begin{aligned}
\|\beta \cdot \nabla(u - u_h)\|^2 &= (\beta \cdot \nabla(u - u_h), \beta_h \cdot \nabla \xi_h - I_{os} \beta_h \cdot \nabla \xi_h) - (\sigma(u - u_h), I_{os} \beta_h \cdot \nabla \xi_h) \\
&\quad - s_a(z_h, I_{os} \beta_h \cdot \nabla \xi_h) \\
&\quad + (\beta \cdot \nabla(u - u_h), (\beta - \beta_h) \cdot \nabla \xi_h) - (\beta \cdot \nabla(u - u_h), \beta \cdot \nabla(\pi_V u - u)) \\
&\leq C_\gamma \|\beta \cdot \nabla(u - u_h)\| (h^{-\frac{1}{2}} |\xi_h|_{S_p} + \|\beta\|_{W^{1,\infty}(\Omega)} \|\xi_h\| + \|\beta \cdot \nabla(u - \pi_V u)\|) \\
&\quad + (C_{\gamma\beta\sigma} h^{-\frac{1}{2}} |z_h|_{S_a} + \|\sigma\|_{L^\infty(\Omega)} \|u - u_h\|) \|\beta_h \cdot \nabla \xi_h\|.
\end{aligned}$$

Here we have used the  $L^2$ -stability of the interpolation operator  $I_{os}$  and the inequality

$$|s_a(z_h, I_{os} \beta_h \cdot \nabla \xi_h)| \leq |z_h|_{S_a} C_{\gamma\beta\sigma} h^{-\frac{1}{2}} \|\beta_h \cdot \nabla \xi_h\|.$$

Observing that

$$\|\beta_h \cdot \nabla \xi_h\| \leq C \|\beta\|_{W^{1,\infty}(\Omega)} \|\xi_h\| + \|\beta \cdot \nabla(u - u_h)\| + \|\beta \cdot \nabla(u - \pi_V u)\|$$

and using suitable arithmetic-geometric inequalities to absorb factors  $\|\beta \cdot \nabla(u - u_h)\|$  in the left-hand side we conclude that

$$\begin{aligned}
\|\beta \cdot \nabla(u - u_h)\|^2 &\leq C_{\gamma\beta\sigma} \left( h^{-1} |\xi_h|_{S_p}^2 + \|\xi_h\|^2 + \|u - u_h\|^2 \right. \\
&\quad \left. + h^{-1} |z_h|_{S_a}^2 + \|\beta \cdot \nabla(u - \pi_V u)\|^2 \right) \leq C_{\gamma\beta\sigma} h^{2k} |u|_{H^{k+1}(\Omega)}^2.
\end{aligned}$$

The last inequality is a consequence of the estimate

$$\|u - u_h\|_V + |u - u_h|_{S_p} + |z_h|_{S_a} \leq C_{\gamma\beta\sigma} h^{k+\frac{1}{2}} |u|_{H^{k+1}(\Omega)}$$

of Theorem 2.2 and standard approximation results on  $\|u - \pi_V u\|$  and  $\|\beta \cdot \nabla(u - \pi_V u)\|$ . The proof for the discontinuous Galerkin method is similar and is left to the reader.  $\square$

**3.5. The data assimilation case.** The aim of the methods presented in [6] is to introduce a framework where also ill-posed problems such as those arising in inverse problems or data assimilation problems can be included, without modifying the method. We will therefore in this section discuss the case where *data is given on the outflow boundary* in (1.1) as a model case of data assimilation. By the reversibility of the transport equation under our assumptions on  $\beta$  this problem is not ill-posed on the continuous level. However, on the discrete level methods based on unwinding are likely to experience difficulties. Since our framework relies on neither unwinding nor coercivity, this case can be included with only minor modifications in the formulations without any loss of stability. Consider the problem (1.1) with the boundary condition  $u = g$  on  $\partial\Omega_+$ . Let the formulation (2.7) be defined by the bilinear form (3.3) and the stabilization term  $s_p(\cdot, \cdot)$  for  $X = GLS, CIP, DG$ ,

$$(3.37) \quad s_p(u_h, v_h) := s_{p,X}(u_h, v_h) + s_{bc,+}(u_h, v_h).$$

The term  $s_a(\cdot, \cdot)$  is unchanged. The data assimilation problem then typically consists in finding  $u|_{\partial\Omega_-}$ , which amounts to solving the backward transport equation. Observe that the boundary penalty for the primal equation now acts on the outflow boundary. The stabilization may then be chosen as any of the three methods considered in sections 3.1–3.3 and Theorem 2.2 holds under the same conditions as before, but the stability will be given by a different weight function. Once the functions  $v_a$  and  $v_{a^*}$  have been identified the rest of the analysis is identical to that of sections 3.1–3.3. We recall the following inequalities from Lemma 3.1.

LEMMA 3.12. *For the bilinear form (3.3) there holds, for all  $\eta \in W^{1,\infty}(\Omega)$ ,*

$$a_h(u_h, -e^\eta u_h) = -\frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) u_h^2 e^\eta \, ds + \int_{\Omega} u_h^2 \left( \frac{1}{2} \beta \cdot \nabla \eta + \frac{1}{2} \nabla \cdot \beta - \sigma \right) e^\eta \, dx,$$

$$a_h(e^\eta z_h, z_h) = \frac{1}{2} \int_{\partial\Omega} (\beta \cdot n) z_h^2 e^\eta \, ds + \int_{\Omega} z_h^2 \left( \frac{1}{2} \beta \cdot \nabla \eta + \frac{1}{2} \nabla \cdot \beta - \sigma \right) e^\eta \, dx.$$

It follows that apart from the form of the exponential dependencies in the constants nothing changes for the method (2.7). The situation is different for method (2.6), since here the same test function must be used in the forms  $a_h(\cdot, \cdot)$  and  $s_p(\cdot, \cdot)$ . We see that the choice  $v_a(u_h) := -\pi_V(e^\eta u_h)$  is necessary in  $a_h(\cdot, \cdot)$ ; however, due to the least squares character of  $s_p(\cdot, \cdot)$ , the term can never have a stabilizing effect for positive stabilization parameter when this weight function is used. If instead the stabilization parameters in (2.6) are chosen *negative* it is straightforward to show that the assumptions for Proposition 2.1 hold. This corresponds to using downwind fluxes instead of upwind fluxes. For more general problems, however, data are provided at some points along the characteristics and it is therefore not possible for any given point in the domain to decide whether the data will arrive from the upwind or the downwind side unless the characteristic equations are solved for each given data. Therefore the strategy of changing the sign of the stabilization parameter inside the domain to match the location of given data is not so attractive. In contrast the method (2.7) does not use the flow direction for stability and can therefore be applied in a much wider context, without tuning the stabilization parameters.

**4. Numerical examples.** Here we will give some simple numerical examples illustrating the above theory. All computations were made using Freefem++ [15]. We will only consider the CIP method and compare the results obtained by (2.6) with those of (2.7) and in some cases with the standard Galerkin method. We use an exact solution from [10] adapted for the case of vanishing viscosity with some different velocity fields. We consider pure transport on conservation form and with a nonsolenoidal velocity field,

$$(4.1) \quad \nabla \cdot (\beta u) = f \quad \text{on } \Omega.$$

Three different velocity fields will be used:

$$(4.2) \quad \beta_1 := \begin{pmatrix} -(x+1)^4 + y \\ -8(y-x) \end{pmatrix},$$

$$(4.3) \quad \beta_2 := -100 \begin{pmatrix} x+y \\ y-x \end{pmatrix},$$

or

$$(4.4) \quad \beta_3 = \begin{pmatrix} 10 \arctan\left(\frac{y-\frac{1}{2}}{\varepsilon}\right) - \frac{x^2}{\varepsilon} \\ \sin(x/\varepsilon) + \sin(y/\varepsilon) \end{pmatrix}.$$

We will consider two different exact solutions, one smooth given by

$$(4.5) \quad u(x, y) = 30x(1-x)y(1-y),$$

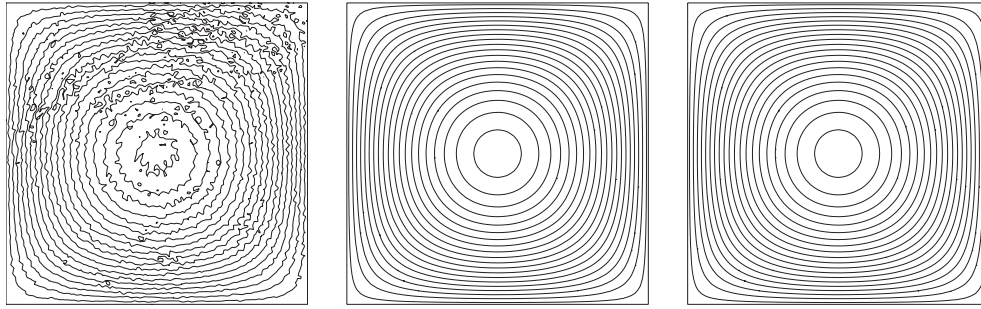


FIG. 1. Contour plots of approximations of the smooth solution (4.5),  $64 \times 64$  mesh, affine approximation. From left to right: standard Galerkin, method (2.6), method (2.7).

obtained by choosing a suitable right-hand-side  $f$ , and one nonsmooth obtained by setting  $f = 0$ , but introducing a discontinuous function for the boundary data. The smooth solution (4.5) satisfies homogeneous Dirichlet boundary conditions both on the inflow and the outflow boundary and has  $\|u\| = 1$ . Unless otherwise stated, we use the stabilization parameters  $\gamma_{CIP} = 0.01$  for piecewise affine approximation and  $\gamma_{CIP} = 0.001$  for piecewise quadratic approximation. The boundary penalty term is taken as  $\gamma_{bc} = 0.5$  for (2.7) and  $\gamma_{bc} = 1.0$  for (2.6).

We have first considered the velocity field (4.2) and the solution (4.5). Note that  $\inf_{x \in \Omega} \nabla \cdot \beta_1 = -40$ , making the problem strongly noncoercive, since then  $\sigma_0 = \frac{1}{2} \inf_{x \in \Omega} \nabla \cdot \beta_1 = -20$ . In our experience the standard Galerkin method performs relatively well for the coercive case when approximating smooth solutions in two space dimensions. As can be seen in Figure 1, this is not the case here. Three contour plots are presented representing computations using the standard Galerkin method, the method (2.6), and (2.7) on a  $64 \times 64$  unstructured mesh. Note the oscillations that persist in the standard Galerkin solution, despite the smoothness of the solution. These oscillations remained on all the meshes considered, up to a finest mesh with  $256 \times 256$  elements, although their amplitude decreased. This highlights the increased need of stabilization for noncoercive problems. In Table 1 we present the errors in both the  $L^2$ -norm and the streamline derivative norm,

$$(4.6) \quad \|h^{\frac{1}{2}} |\beta|^{-\frac{1}{2}} \beta \cdot \nabla(u - u_h)\|,$$

on six consecutive unstructured meshes with  $2^N$ ,  $N = 3, \dots, 8$ , elements on each side and piecewise affine approximation. We note that the stabilized methods both have (and sometimes exceed) the expected convergence orders. Indeed the  $L^2$ -error converges as  $O(h^{k+1})$  and the error in the streamline derivative (4.6) as  $O(h^{k+\frac{1}{2}})$ . As expected the convergence of the standard Galerkin method is very uneven. It is unclear if the error in the streamline derivative converges at all. In Table 2 the same sequence of computations is reported using piecewise quadratic elements. The stability of the standard Galerkin method is noticeably improved. Nevertheless the errors of the stabilized methods are two orders of magnitude smaller. The errors of formulation (2.7) are slightly smaller than those of (2.6), but on the other hand the former method uses twice as many degrees of freedom as the latter.

Both methods (2.6) and (2.7) control spurious oscillations in nonsmooth exact solutions, as can be seen in Figure 2, where the contour plots of a computation with nonsmooth exact solution created by using the velocity field (4.3) in (4.1) setting  $f = 0$

TABLE 1

Errors of estimated quantities for the smooth solution approximated using piecewise affine elements. SG means standard Galerkin and equations refer to methods used.  $L^2$  denotes the error in the  $L^2$ -norm, and SD denotes the error in the streamline derivative norm defined in (4.6).

$N$	SG, $L^2$	SG, SD	(2.6), $L^2$	(2.6), SD	(2.7), $L^2$	(2.7), SD
3	0.041	1.0	0.029	0.58	0.028	0.58
4	0.025	0.88	7.2E-3	0.20	6.5E-3	0.20
5	0.010	0.48	1.7E-3	0.071	1.5E-3	0.069
6	0.015	1.1	4.5E-4	0.026	4.0E-4	0.025
7	7.8E-3	0.76	1.1E-4	9.1E-3	1.0E-4	8.7E-3
8	1.9E-3	1.1	2.5E-5	3.0E-3	2.4E-5	3.0E-3

TABLE 2

Errors of estimated quantities for the smooth solution approximated using piecewise quadratic elements.

$N$	SG, $L^2$	SG, SD	(2.6), $L^2$	(2.6), SD	(2.7), $L^2$	(2.7), SD
3	0.028	0.58	9.3E-4	0.060	7.5E-4	0.045
4	4.6E-3	0.25	1.7E-4	0.014	1.1E-4	8.7E-3
5	1.9E-3	0.17	2.7E-5	3.1E-3	1.4E-5	1.7E-3
6	3.0E-4	0.042	3.3E-6	5.1E-4	1.7E-6	2.7E-4
7	3.3E-5	6.1E-3	4.4E-7	9.2E-5	2.1E-7	4.7E-5

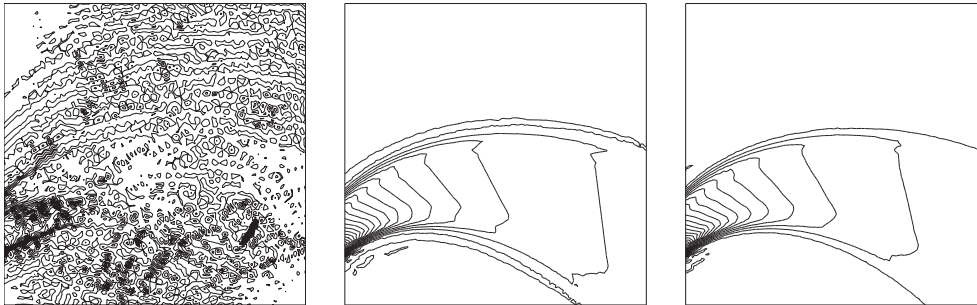


FIG. 2. Discontinuous solution,  $64 \times 64$  mesh, affine approximation. From left to right: standard Galerkin, method (2.6), method (2.7).

and the boundary data equal one wherever  $x > 0.8$  and  $y < 0.5$  and zero elsewhere. To show the increased robustness of the formulation (2.7), we propose to study the problem (4.1) with the velocity field (4.4). This velocity field is strictly speaking not covered by the analysis, since for some values on  $\varepsilon$  there may be points in the domain where  $\beta_3$  vanishes. Nevertheless the right-hand side is chosen such that the exact solution is given by (4.5). We consider a fixed  $64 \times 64$  unstructured mesh and vary  $\varepsilon$ , creating a series of increasingly ill-posed problems where the divergence and the maximum derivatives of  $\beta$  behaves as  $-\frac{1}{\varepsilon}$ . The error in the streamline derivative (4.6) for varying  $\varepsilon$  is plotted in the left graphic of Figure 3. It is fair to say that the method (2.7) (circle markers) outperforms (2.6) (square markers). As  $\varepsilon$  becomes small the error for the approximations computed using (2.7) exhibits moderate growth of order  $O(\varepsilon^{-\frac{1}{3}})$  but remains below 0.06, whereas over half the approximations computed using (2.6) has an error larger than 0.5 and none below 0.1. For  $\varepsilon = 0.05$ , the error is 120 and the computed solution bears no resemblance to the exact one. In the right plot of Figure 3 we study how the error depends on the choice of the stabilization parameter



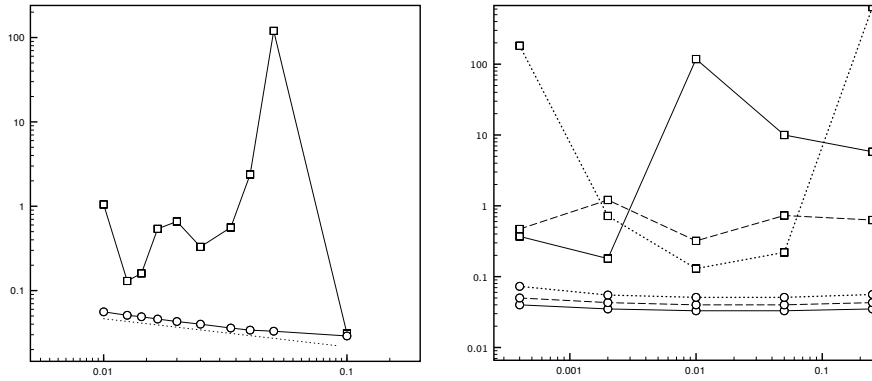


FIG. 3. Study of the error in the SD-norm error (4.6). Circles: method (2.7), squares: method (2.6). Left: under variation of  $\epsilon$  in (4.4), with  $\gamma_{CIP} = 0.01$ , dotted line  $O(\epsilon^{-\frac{1}{3}})$ . Right: under variation of  $\gamma_{CIP}$  for different  $\epsilon$  (full line,  $\epsilon = 0.05$ ; dashed line,  $\epsilon = 0.025$ ; dotted line,  $\epsilon = 0.0125$ ).

TABLE 3

Data assimilation using (2.7). Errors of estimated quantities for the smooth solution (4.5) computed with data given on the outflow boundary. Approximation using piecewise affine ( $\mathbb{P}_1$ ) and quadratic ( $\mathbb{P}_2$ ) elements.

$N$	$\mathbb{P}_1, L^2$	$\mathbb{P}_1, SD$	$\mathbb{P}_2, L^2$	$\mathbb{P}_2, SD$
3	0.033	0.75	1.1E-3	0.052
4	7.1E-3	0.23	1.5E-4	9.6E-3
5	1.6E-3	0.075	1.8E-5	1.8E-3
6	4.1E-4	0.026	2.0E-6	2.8E-4
7	1.0E-4	8.9E-3	2.4E-7	4.8E-5
8	2.4E-5	3.0E-3	—	—

$\gamma_{CIP}$ . We plot the error defined by (4.6), this time varying the parameter  $\gamma_{CIP}$  for three different  $\epsilon$ . Even when accounting for the increased number of degrees of freedom in method (2.7) the error of (2.6) is more than 50% large in all the computations and where (2.6) fails it is more than a factor 1000 larger.

**4.1. A data assimilation example.** Finally we consider a model problem for data assimilation where the boundary conditions of the problem (4.1) are imposed on the outflow boundary instead of the inflow boundary. Method (2.7) with the bilinear form (3.3) and the stabilizing term (3.37) with  $X = CIP$  was applied. We consider the test case with smooth solution (4.5) and velocity field (4.2). In Table 3 we give the computational errors in the  $L^2$ -norm and the streamline norm (4.6), using either piecewise affine or piecewise quadratic elements. Recalling the results in Tables 1 and 2 we see that the errors are comparable. This is not surprising since the use of the adjoint equation makes the two cases similar. Attempts to use (2.6) with weakly imposed boundary conditions on the outflow and  $\gamma_{CIP} > 0$  were not fruitful. This is expected since the stabilized methods on the form (2.6) all are based on upwinding, which is unphysical in this setting. Indeed the standard unstabilized Galerkin method performs better than the standard stabilized method for this smooth solution. When the stabilization parameter is chosen negative we recover the expected behavior of the



TABLE 4

Data assimilation using the method (2.6) with the forms (3.3) and (3.37), piecewise affine elements,  $\gamma_{bc} = -1$ , and three different choices of  $\gamma_{CIP}$  denoted by  $\gamma_1$ ,  $\gamma_2$ , and  $\gamma_3$ . The CIP stabilization parameters are assigned the values  $\gamma_1 = 10^{-3}$ ,  $\gamma_2 = 0$ , and  $\gamma_3 = -10^{-2}$ . Errors of estimated quantities for the smooth solution (4.5) computed with data given on the outflow boundary.

$N$	$\gamma_1, L^2$	$\gamma_1, SD$	$\gamma_2, L^2$	$\gamma_2, SD$	$\gamma_3, L^2$	$\gamma_3, SD$
3	0.044	3.48	0.034	2.8	0.029	2.25
4	0.027	2.96	0.01	1.2	6.7E-3	0.74
5	0.27	31.0	2.7E-3	0.44	1.6E-3	0.26
6	2.74	455	1.1E-3	0.26	4.2E-4	0.094
7	6170	1.8E6	3.7E-4	0.11	1.1E-4	0.033
8	67471	3.4E7	9.9E-5	0.041	2.5E-3	0.011

stabilized method. We give the results of (2.6) using  $\gamma_{bc} = -1.0$  and  $\gamma_{CIP} = 0.001$ ,  $\gamma_{CIP} = 0$ ,  $\gamma_{CIP} = -0.01$  in Table 4.

**5. Concluding remarks.** We have extended the methods proposed in [6] to include hyperbolic equations and have shown how three stabilization methods known from the literature can be used to obtain stable and (quasi-) optimally convergent approximations. Compared to the standard stabilized method we show that the new method yields existence of discrete solutions and (quasi-) optimal error estimates under much weaker assumptions on the mesh parameter (“ $h^2$  small enough” compared to “ $h^{\frac{1}{2}}$  small enough”). We would like to stress that the method proposed here will not necessarily yield a more accurate solution than the standard stabilized methods in cases where both methods work. The new method, however, has increased robustness for noncoercive problems. It also makes it easier to incorporate data other than classical inflow boundary data. The idea of recasting the problem in an optimization framework opens interesting perspectives for optimal control, inverse problems, and data assimilation using observers.

## REFERENCES

- [1] B. AYUSO AND L. D. MARINI, *Discontinuous Galerkin methods for advection-diffusion-reaction problems*, SIAM J. Numer. Anal., 47 (2009), pp. 1391–1420.
- [2] I. BABUSKA, *Error-bounds for finite element method*, Numer. Math., 16 (1971), pp. 322–333.
- [3] R. BECKER AND M. BRAACK, *A two-level stabilization scheme for the Navier-Stokes equations*, in Numerical Mathematics and Advanced Applications, Springer, Berlin, 2004, pp. 123–130.
- [4] S. BERTOLUZZA, *The discrete commutator property of approximation spaces*, C. R. Acad. Sci. Paris Ser. I Math., 329 (1999), pp. 1097–1102.
- [5] 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.
- [6] E. BURMAN, *Stabilized finite element methods for nonsymmetric, noncoercive, and ill-posed problems. Part I: Elliptic equations*, SIAM J. Sci. Comput., 35 (2013), pp. A2752–A2780.
- [7] E. BURMAN, *A unified analysis for conforming and nonconforming stabilized finite element methods using interior penalty*, SIAM J. Numer. Anal., 43 (2005), pp. 2012–2033.
- [8] E. BURMAN, M. A. FERNÁNDEZ, AND P. HANSBO, *Continuous interior penalty finite element method for Oseen’s equations*, SIAM J. Numer. Anal., 44 (2006), pp. 1248–1274.
- [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] C. CHAINAIS-HILLAIRET AND J. DRONIOU, *Finite-volume schemes for noncoercive elliptic problems with Neumann boundary conditions*, IMA J. Numer. Anal., 31 (2011), pp. 61–85.

- [11] R. CODINA, *Stabilization of incompressibility and convection through orthogonal sub-scales in finite element methods*, Comput. Methods Appl. Mech. Engrg., 190 (2000), pp. 1579–1599.
- [12] J. DOUGLAS AND T. DUPONT, *Interior penalty procedures for elliptic and parabolic Galerkin methods*, in Computing Methods in Applied Sciences, Lecture Notes in Phys. 58, Springer, Berlin, 1976, pp. 207–216.
- [13] A. ERN AND J.-L. GUERMOND, *Discontinuous Galerkin methods for Friedrichs' systems. I. General theory*, SIAM J. Numer. Anal., 44 (2006), pp. 753–778.
- [14] J.-L. GUERMOND, *Stabilization of Galerkin approximations of transport equations by subgrid modeling*, M2AN Math. Model. Numer. Anal., 33 (1999), pp. 1293–1316.
- [15] F. HECHT, *New development in FreeFem++*, J. Numer. Math., 20 (2012), pp. 251–265.
- [16] C. JOHNSON, U. NÄVERT, AND J. PITKÄRANTA, *Finite element methods for linear hyperbolic problems*, Comput. Methods Appl. Mech. Engrg., 45 (1984), pp. 285–312.
- [17] C. JOHNSON AND J. PITKÄRANTA, *An analysis of the discontinuous Galerkin method for a scalar hyperbolic equation*, Math. Comp., 46 (1986), pp. 1–26.
- [18] P. LESAINTE AND P.-A. RAVIART, *On a finite element method for solving the neutron transport equation*, in Mathematical Aspects of Finite Elements in Partial Differential Equations, Academic Press, New York, 1974, pp. 89–123.
- [19] W. H. REED AND T. R. HILL, *Triangular Mesh Methods for the Neutron Transport Equation*, Tech. report LA-UR-73-479, Los Alamos National Laboratory, Los Alamos, NM, 1973.
- [20] A. SCHATZ, *An observation concerning Ritz-Galerkin methods with indefinite bilinear forms*, Math. Comp., 28 (1974), pp. 959–962.