

# FORMULATION AND ANALYSIS OF A PARAMETER-FREE STABILIZED FINITE ELEMENT METHOD\*

P. BOCHEV<sup>†</sup>, M. PEREGO<sup>†</sup>, AND K. PETERSON<sup>†</sup>

**Abstract.** We present a stabilized finite element method for the scalar advection-diffusion equation, which does not require tunable mesh-dependent parameters. Stabilization is achieved by using diffusive fluxes extracted from an edge element lifting of Scharfetter-Gummel edge fluxes into the elements. Although the method is formally first-order accurate, qualitative numerical studies suggest that it occupies a middle ground between an artificial diffusion and a streamline-upwind Petrov-Galerkin formulations. The method is substantially less dissipative than the former, while having much smaller overshoots and undershoots than the latter.

**Key words.** Advection-diffusion, stabilization, Scharfetter-Gummel upwinding, finite elements, edge elements.

**AMS subject classifications.** 65N30, 65N12, 65N15, 65L11, 65L60

**1. Introduction.** We consider the scalar advection-diffusion equation

$$(1.1) \quad \begin{cases} -\nabla \cdot F(\phi) = f & \text{in } \Omega \\ F(\phi) = (\varepsilon \nabla \phi - \mathbf{u} \phi) & \text{in } \Omega \end{cases} \quad \text{and} \quad \phi = g \text{ on } \Gamma$$

where  $\varepsilon$  is a diffusion coefficient,  $\mathbf{u}$  is the advective velocity, and  $f$  and  $g$  are given functions. When  $\varepsilon$  is small relative to  $\mathbf{u}$ , solutions of (1.1) can develop internal and/or boundary layers. If the grid is not fine enough to resolve these layers, Galerkin solutions of (1.1) exhibit spurious oscillations.

To stabilize the solution one can use a variety of tools ranging from artificial diffusion [14] to consistently stabilized methods such as SUPG [12, 13], and multiscale and enriched methods [11, 18]. However, most if not all of the existing stabilized methods for (1.1) require a mesh-dependent stabilization parameter. The choice of this parameter is critical for the accuracy and stability of the corresponding finite element solution. Yet, finding the best possible stabilization parameter for a given problem is difficult and remains an open question [15, 7]. The principal reasons for this are (i) the dependence of the stabilization parameter on constants that are known exactly only in special cases [10], and (ii) the fact that different solution features, such as internal and boundary layers may require different stabilization techniques; see [15]. As a result, stabilization parameters may have to be individually tailored to different applications using heuristic arguments [5], numerical calibration or a combination thereof.

This paper presents and studies a new stabilized finite element method for (1.1), which does not require tunable mesh-dependent stabilization parameters. Stabilization is achieved by augmenting the standard Galerkin formulation of (1.1) with an artificial diffusion tensor defined by a product of two diffusive fluxes. We extract these fluxes from an approximation  $F_E \in H(\text{curl}, \Omega)$  of the total flux  $F$ , defined by an edge element lifting of Scharfetter-Gummel edge fluxes into the elements [3, 4]. Previously, we used  $F_E(\phi_h)$  to replace the standard nodal flux  $F(\phi_h)$  in a Galerkin [3] and a control volume finite element (CVFEM) [4] formulations of the governing equations (1.1). These methods provide multidimensional extensions of the Scharfetter-Gummel scheme [17] to unstructured grids and have been implemented in Sandia's device modeling code Charon with encouraging results.

---

\*Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energys National Nuclear Security Administration under contract DE-AC04-94AL85000

<sup>†</sup>Computational Mathematics, Sandia National Laboratories, Mail Stop 1320 Albuquerque, New Mexico, 87185-1320 (`{pbboche,mperego,kjpeter}@sandia.gov`)

In this paper the edge element flux  $F_E$  is a starting point for the development of a new stabilized Galerkin formulation that improves upon the robustness and accuracy of the method in [3]. A careful analysis reveals that  $F_E$  comprises a projection of the nodal flux  $F(\phi_h)$  onto an edge element space, and a stabilizing diffusive flux. The pairing of the latter with the gradient of a nodal finite element function  $\phi_h$  produces a non-symmetric diffusion tensor whose stabilizing effect may deteriorate on unstructured grids. To correct this drawback we (i) discard the part of the  $F_E$  that corresponds to the projection of  $F(\phi_h)$  onto the edge element space, and (ii) use the remaining diffusive flux to define a symmetric diffusion tensor. In a nutshell, the new method comprises a standard nodal Galerkin formulation of (1.1), augmented by this symmetric artificial diffusion tensor. Because edge elements are used only locally, assembly of the diffusion tensor does not require global edge data structured and can be easily incorporated into an existing Galerkin code for (1.1).

The rest of this section introduces the notation. Sections 2–3 present the formulation and the analysis of the method, respectively. Section 4 illustrates the method numerically, including a demonstration of its improved stability on unstructured grids.

**1.1. Notation.** In this paper  $\Omega \subset \mathbb{R}^n$ ,  $n = 2, 3$  is a bounded domain with Lipschitz-continuous boundary  $\Gamma = \partial\Omega$ ,  $W^{k,p}(\Omega)$  is a Sobolev space of order  $k$ ,  $W_0^{k,p}(\Omega)$  is the subspace of its functions with vanishing trace,  $L^p(\Omega) = W^{0,p}(\Omega)$ , and  $H^k(\Omega) = W^{k,2}(\Omega)$ . The space  $H(\text{curl}, \Omega)$  contains vector fields in  $L^2(\Omega)^n$  whose curl belongs in  $L^2(\Omega)^{2n-3}$ . The meaning of the symbol  $|\cdot|$  depends on the context and can be Euclidean length, semi-norm, domain measure, or cardinality of a finite set. We use lower case bold face for vectors, upper case Roman for matrices,  $\underline{\lambda}(A)$  and  $\overline{\lambda}(A)$  are the smallest and largest eigenvalues of  $A$ , respectively, and  $\kappa(A)$  is the condition number of  $A$ . The versors of a Cartesian coordinate system are  $\mathbf{i} = (1, 0, 0)$ ,  $\mathbf{j} = (0, 1, 0)$ , and  $\mathbf{k} = (0, 0, 1)$  in  $\mathbb{R}^3$ .

We use shape-regular finite element partitions  $K(\Omega)$  of  $\Omega$  into elements  $\mathbf{k}$  with size  $h_{\mathbf{k}}$ . The average element size is denoted by  $h$ . For brevity we restrict attention to triangles and quadrilaterals in 2D and tetrahedrons and hexahedrons in 3D. We assume that the simplicial elements are affine, the quadrilaterals are bilinear images of the unit square, and the hexahedrals are trilinear images of the unit cube. The symbols  $V(\star)$  and  $E(\star)$  stand for the sets of all vertices and edges in  $K(\Omega)$  belonging to an entity  $\star$ . For instance,  $V(\Omega)$  is the set of all vertices,  $E(\Omega)$  is the set of all edges,  $V(\mathbf{k})$  are the vertices of element  $\mathbf{k}$ ,  $E(\mathbf{v}_i)$  are all edges having  $\mathbf{v}_i$  as a vertex, and so on. We label edges by a multi-index  $\alpha = (\alpha_1, \alpha_2)$  comprising the indices of their endpoints, i.e.,  $\mathbf{e}_\alpha$  is an edge with endpoints  $\mathbf{v}_{\alpha_1}$  and  $\mathbf{v}_{\alpha_2}$ . We also need the node-to-edge incidence matrix  $G$  with element  $g_{\alpha,j}$  and dimension  $|E(\Omega)| \times |V(\Omega)|$ . Assuming that every edge  $\mathbf{e}_\alpha$  in the mesh is oriented by choosing the order of its vertices,  $g_{\alpha,j} = 0$  if  $\mathbf{v}_j$  is not a vertex of  $\mathbf{e}_\alpha$  and

$$(1.2) \quad g_{\alpha,j} = \begin{cases} -1 & \text{if } \mathbf{v}_j \text{ is the first vertex of } \mathbf{e}_\alpha \\ 1 & \text{if } \mathbf{v}_j \text{ is the second vertex of } \mathbf{e}_\alpha \end{cases}$$

For most “reasonable” finite element partitions, e.g., meshes satisfying the assumptions of [6, p.51] the matrix  $G$  has a one dimensional kernel comprising the constant vector. Thus,  $G$  can be interpreted as a “topological” gradient operator that depends only on the mesh connectivity.

The midpoint, the length, and the unit tangent of an oriented edge  $\mathbf{e}_\alpha$  are

$$\mathbf{m}_\alpha = \frac{\mathbf{v}_{\alpha_1} + \mathbf{v}_{\alpha_2}}{2}, \quad h_\alpha = |\mathbf{v}_{\alpha_1} - \mathbf{v}_{\alpha_2}|, \quad \text{and} \quad \mathbf{t}_\alpha = g_{\alpha,\alpha_1} \frac{\mathbf{v}_{\alpha_1} - \mathbf{v}_{\alpha_2}}{|\mathbf{v}_{\alpha_1} - \mathbf{v}_{\alpha_2}|},$$

respectively. Definition of  $\mathbf{t}_\alpha$  implies that it always points towards the second vertex of  $\mathbf{e}_\alpha$ .

In what follows,  $N^h(\Omega)$  is the isoparametric  $C^0$  piecewise linear, bilinear or trilinear *nodal* finite element space and  $E^h(\Omega)$  is the lowest-order Nedelec *edge element* space [16].  $N_0^h(\Omega)$  is the

subspace of all functions in  $N^h(\Omega)$  that vanish on  $\Gamma$ . The standard nodal basis  $\{N_i\}$ ,  $\mathbf{v}_i \in V(\Omega)$  of  $N^h(\Omega)$  has the property  $N_i(\mathbf{v}_j) = \delta_i^j$ . The coefficient vector of  $\phi_h \in N^h(\Omega)$  is  $\boldsymbol{\phi} = (\phi_i, \dots, \phi_m)$ ,  $m = |V(\Omega)|$ . The standard basis of  $E^h(\Omega)$  is  $\{\vec{W}_\xi\}$ ,  $\mathbf{e}_\xi \in E(\Omega)$  such that

$$(1.3) \quad \int_{\mathbf{e}_\eta} \vec{W}_\xi \cdot \mathbf{t}_\eta dl = \delta_\xi^\eta \quad \text{and} \quad \vec{W}_\xi \cdot \mathbf{t}_\eta \Big|_{\mathbf{e}_\eta} = \text{const} = \frac{\delta_\xi^\eta}{h_\eta}.$$

Since  $\vec{W}_\xi \cdot \mathbf{t}_\xi > 0$  the orientation of  $\vec{W}_\xi$  matches the orientation of its associated edge  $\mathbf{e}_\xi$ .  $N^h(\Omega)$  and  $E^h(\Omega)$  belong to an exact sequence [1] and  $\nabla N^h(\Omega) \subset E^h(\Omega)$ . In particular, there holds

$$(1.4) \quad \nabla N_i = \sum_{\mathbf{e}_\xi \in E(\mathbf{v}_i)} g_{\xi,i} \vec{W}_\xi \quad \text{and} \quad \nabla \phi_h = \sum_{\mathbf{e}_\xi \in E(\Omega)} (g_{\xi,\xi_1} \phi_{\xi_1} + g_{\xi,\xi_2} \phi_{\xi_2}) \vec{W}_\xi = \sum_{\mathbf{e}_\xi \in E(\Omega)} (\mathbf{g}_\xi \boldsymbol{\phi}) \vec{W}_\xi,$$

where  $\mathbf{g}_\xi$  is the row of  $G$  corresponding to edge  $\mathbf{e}_\xi$ . The last identity implies the following factorization of the weak Laplace operator; see e.g. [2, p.108]:

$$(1.5) \quad \int_{\Omega} \nabla \phi_h \cdot \nabla \varphi_h dx = \boldsymbol{\phi}^T (G^T M G) \boldsymbol{\varphi},$$

where  $M$  is the Gramm matrix of the edge element basis  $\{\vec{W}_\xi\}$ . Since  $G$  has one-dimensional nullspace spanned by the constant vector, it follows that for  $\phi_h \in N_0^h(\Omega)$  with coefficient vector  $\boldsymbol{\phi}$ ,

$$(1.6) \quad G \boldsymbol{\phi} = 0 \quad \text{if and only if} \quad \phi_h = 0.$$

We recall the nodal interpolation operator  $\mathcal{I}_N : H^1(\Omega) \cap C^0(\Omega) \mapsto N^h(\Omega)$ ,

$$(1.7) \quad \mathcal{I}_N(\phi) = \sum_{\mathbf{v}_i \in V(\Omega)} \phi(\mathbf{v}_i) N_i(\mathbf{x}),$$

and, for  $p > 2$ , the edge interpolation operator  $\mathcal{I}_E : H(\text{curl}, \Omega) \cap (W^{1,p}(\Omega))^n \mapsto E^h(\Omega)$ ,

$$(1.8) \quad \mathcal{I}_E(\mathbf{u}) = \sum_{\mathbf{e}_\xi \in E(\Omega)} \vec{W}_\xi \int_{\mathbf{e}_\xi} \mathbf{u} \cdot \mathbf{t}_\xi dl.$$

We will also need the local interpolation result [9, Theorem 1.103]

$$(1.9) \quad |\mathcal{I}_N(\phi) - \phi|_{m,p,\mathbf{k}} \leq C h_{\mathbf{k}}^{l+1-m} |\phi|_{l+1,p,\mathbf{k}},$$

which holds for all  $\phi \in W^{l+1,p}(\mathbf{k})$  with  $\mathbf{k} \in K(\Omega)$  and  $1 \leq p \leq \infty$ ,  $0 \leq l \leq 1$ , and the fact that

$$(1.10) \quad \|\mathcal{I}_N(\phi)\|_1 \leq C \|\phi\|_2 \quad \forall \phi \in H^2(\Omega).$$

**2. Formulation.** Consider the Galerkin method for (1.1): find  $\phi_h \in N^h(\Omega)$  such that

$$(2.1) \quad a(\phi_h, \varphi_h) = b(\varphi_h) \quad \forall \varphi_h \in N_0^h(\Omega),$$

where the bilinear form  $a(\cdot, \cdot)$  and the linear functional  $b(\cdot)$  are given by

$$(2.2) \quad a(\phi_h, \varphi_h) = \int_{\Omega} F(\phi_h) \cdot \nabla \varphi_h dV \quad \text{and} \quad b(\varphi_h) = \int_{\Omega} f \varphi_h dV,$$

respectively. We will stabilize (2.1) using a diffusive flux  $\tilde{\Theta}(\phi_h)$  extracted from an edge element approximation  $F_E(\phi_h)$  of the total flux  $F(\phi_h)$ . We define the former by an edge element lifting of one-dimensional, Scharfetter-Gummel edge fluxes  $F_\alpha$  into the elements.

**2.1. The  $H(\text{curl}, \Omega)$  edge element flux.** To define the edge fluxes  $F_\alpha$  and their  $H(\text{curl}, \Omega)$  lifting  $F_E$ , we temporarily adopt the convention that the first vertex of  $\mathbf{e}_\alpha$  is always  $\mathbf{v}_{\alpha_1}$  and so,  $g_{\alpha, \alpha_1} = -1$  and  $g_{\alpha, \alpha_2} = 1$ . Then, along  $\mathbf{e}_\alpha$  we consider the following problem:

$$(2.3) \quad \begin{cases} -\frac{d}{ds} \left( \bar{\varepsilon}_\alpha \frac{d\phi(s)}{ds} - \bar{u}_\alpha \phi(s) \right) = 0 & \text{for } 0 < s < h_\alpha \\ \phi(0) = \phi_{\alpha_1} & \text{and } \phi(h_\alpha) = \phi_{\alpha_2} \end{cases}$$

where the unknown nodal coefficients  $\phi_{\alpha_1}$  and  $\phi_{\alpha_2}$  of  $\phi_h$  specify the boundary data in (2.3), and

$$\bar{u}_\alpha = \frac{1}{h_\alpha} \int_{\mathbf{e}_\alpha} \mathbf{u} \cdot \mathbf{t}_\alpha dl \quad \text{and} \quad \bar{\varepsilon}_\alpha = \frac{1}{h_\alpha} \int_{\mathbf{e}_\alpha} \varepsilon dl$$

are the mean edge velocity and diffusion, respectively. The exact solution of (2.3) is

$$(2.4) \quad \phi(s) = \frac{\exp(2p_\alpha)\phi_{\alpha_1} - \phi_{\alpha_2}}{\exp(2p_\alpha) - 1} + \frac{\phi_{\alpha_2} - \phi_{\alpha_1}}{\exp(2p_\alpha) - 1} \exp(s\bar{u}_\alpha/\bar{\varepsilon}_\alpha); \quad p_\alpha = \frac{\bar{u}_\alpha h_\alpha}{2\bar{\varepsilon}_\alpha}.$$

The value of  $p_\alpha$  relates the rates of advection and diffusion along  $\mathbf{e}_\alpha$  and so we call it *the edge Péclet number*. The edge flux corresponding to the exact solution (2.4) is given by

$$(2.5) \quad F_\alpha = h_\alpha \left( \bar{\varepsilon}_\alpha \frac{d\phi(s)}{ds} - \bar{u}_\alpha \phi(s) \right) = h_\alpha \bar{u}_\alpha \frac{\phi_{\alpha_2} - \exp(2p_\alpha)\phi_{\alpha_1}}{\exp(2p_\alpha) - 1}.$$

After discarding the temporary orientation convention, multiplying and dividing (2.5) by  $\exp(-p_\alpha)$ , and performing some simple algebraic manipulations, the edge fluxes assume the form:

$$(2.6) \quad F_\alpha = \sigma_\alpha \mathbf{g}_\alpha \phi - \frac{h_\alpha \bar{u}_\alpha}{2} (\phi_{\alpha_1} + \phi_{\alpha_2}); \quad \sigma_\alpha = \frac{h_\alpha \bar{u}_\alpha}{2} \coth(p_\alpha) = \bar{\varepsilon}_\alpha p_\alpha \coth(p_\alpha).$$

The lifting of the edge fluxes (2.6) into the edge element space  $E^h(\Omega)$  defines the  $H(\text{curl}, \Omega)$  flux

$$(2.7) \quad F_E(\phi_h) = \sum_{\mathbf{e}_\xi \in E(\Omega)} F_\xi \vec{W}_\xi = \sum_{\mathbf{e}_\xi \in E(\Omega)} \left[ \sigma_\xi \mathbf{g}_\xi \phi - \frac{h_\xi \bar{u}_\xi}{2} (\phi_{\xi_1} + \phi_{\xi_2}) \right] \vec{W}_\xi.$$

**2.2. The stabilizing diffusive flux.** We start by examining the structure of  $F_E$ .

LEMMA 2.1. *Let  $\phi_h \in N^h(\Omega)$ . Then,*

$$(2.8) \quad F_E(\phi_h) = \mathcal{I}_E(\varepsilon \nabla \phi_h - \mathbf{u}_E \phi_h) + \Theta(\phi_h),$$

where  $\mathbf{u}_E = \mathcal{I}_E \mathbf{u}$  and

$$(2.9) \quad \Theta(\phi_h) = \sum_{\mathbf{e}_\xi \in E(\Omega)} \theta_\xi (\mathbf{g}_\xi \phi) \vec{W}_\xi \quad \text{with} \quad \theta_\xi = \sigma_\xi - \bar{\varepsilon}_\xi = \bar{\varepsilon}_\xi (p_\xi \coth(p_\xi) - 1)$$

is a stabilizing diffusive flux.

*Proof.* Adding and subtracting  $\bar{\varepsilon}_\xi \mathbf{g}_\xi \phi$  to (2.6) allows us to write the  $H(\text{curl}, \Omega)$  flux as

$$F_E = \sum_{\mathbf{e}_\xi \in E(\Omega)} \left[ \bar{\varepsilon}_\xi \mathbf{g}_\xi \phi - \frac{h_\xi \bar{u}_\xi}{2} (\phi_{\xi_1} + \phi_{\xi_2}) \right] \vec{W}_\xi + \sum_{\mathbf{e}_\xi \in E(\Omega)} (\sigma_\xi - \bar{\varepsilon}_\xi) (\mathbf{g}_\xi \phi) \vec{W}_\xi.$$

The second term is the stabilizing diffusive flux  $\Theta$ . To complete the proof it remains to show that the first term equals  $\mathcal{I}_E(\varepsilon \nabla \phi_h - \mathbf{u}_E \phi_h)$ . Using (1.8) and the exactness (1.4) yields

$$\begin{aligned} \mathcal{I}_E(\varepsilon \nabla \phi_h) &= \sum_{e_\xi \in E(\Omega)} \left[ \int_{e_\xi} \varepsilon \nabla \phi_h \cdot \mathbf{t}_\xi dl \right] \vec{W}_\xi = \sum_{e_\xi \in E(\Omega)} \left[ \int_{e_\xi} \varepsilon \left( \sum_{e_\eta \in E(\Omega)} (\mathbf{g}_\eta \phi) \vec{W}_\eta \right) \cdot \mathbf{t}_\xi dl \right] \vec{W}_\xi \\ &= \sum_{e_\xi \in E(\Omega)} \left[ \int_{e_\xi} \frac{\varepsilon}{h_\xi} (\mathbf{g}_\xi \phi) dl \right] \vec{W}_\xi = \sum_{e_\xi \in E(\Omega)} \bar{\varepsilon}_\xi (\mathbf{g}_\xi \phi) \vec{W}_\xi. \end{aligned}$$

From the definition of  $\bar{u}_\alpha$  and the linearity of  $\phi_h$  along any edge, it follows that

$$\mathbf{u}_E = \sum_{e_\eta \in E(\Omega)} \left[ \int_{e_\eta} \mathbf{u} \cdot \mathbf{t}_\eta dl \right] \vec{W}_\eta = \sum_{e_\eta \in E(\Omega)} h_\eta \bar{u}_\eta \vec{W}_\eta \quad \text{and} \quad \int_{e_\xi} \phi_h dl = \frac{h_\xi}{2} (\phi_{\xi_1} + \phi_{\xi_2}),$$

respectively. As a result, using (1.3)

$$\begin{aligned} \mathcal{I}_E(\mathbf{u}_E \phi_h) &= \sum_{e_\xi \in E(\Omega)} \left[ \int_{e_\xi} (\mathbf{u}_E \phi_h) \cdot \mathbf{t}_\xi dl \right] \vec{W}_\xi = \sum_{e_\xi \in E(\Omega)} \left[ \int_{e_\xi} \left[ \sum_{e_\eta \in E(\Omega)} h_\eta \bar{u}_\eta \vec{W}_\eta \right] \cdot \mathbf{t}_\xi \phi_h dl \right] \vec{W}_\xi \\ &= \sum_{e_\xi \in E(\Omega)} \left[ \bar{u}_\xi \int_{e_\xi} \phi_h dl \right] \vec{W}_\xi = \sum_{e_\xi \in E(\Omega)} \frac{h_\xi \bar{u}_\xi}{2} (\phi_{\xi_1} + \phi_{\xi_2}) \vec{W}_\xi, \end{aligned}$$

which completes the proof.  $\square$

The following result confirms that  $\Theta$  is responsible for the stabilizing effect of  $F_E$ .

LEMMA 2.2. *For  $\varepsilon > 0$  and any  $\mathbf{u}$  there holds*

$$(2.10) \quad \max \left\{ \bar{\varepsilon}_\alpha, \frac{h_\alpha |u_\alpha|}{2} \right\} \leq \sigma_\alpha \leq \bar{\varepsilon}_\alpha + \frac{h_\alpha |u_\alpha|}{2} \quad \text{and} \quad \max \left\{ 0, \frac{h_\alpha |u_\alpha|}{2} - \bar{\varepsilon}_\alpha \right\} \leq \theta_\alpha \leq \frac{h_\alpha |u_\alpha|}{2}.$$

*In particular,  $\theta_\alpha \geq 0$ ,  $\theta_\alpha \rightarrow 0$  in the diffusive limit, and  $\theta_\alpha \rightarrow h_\alpha \bar{u}_\alpha / 2$  in the advective limit.*

*Proof.* The lemma follows from the inequality

$$\max \{1, |x|\} \leq x \coth(x) \leq 1 + |x| \quad \forall x \in \mathfrak{R}.$$

$\square$

**2.3. The stabilized Galerkin formulation.** In [3] we considered the formulation

$$(2.11) \quad \int_{\Omega} F_E(\phi_h) \cdot \nabla \varphi_h dV = b(\varphi_h) \quad \forall \varphi_h \in N_0^h(\Omega),$$

which stabilizes (2.1) by replacing the nodal flux  $F(\phi_h)$  with the edge element flux  $F_E(\phi_h)$ . The decomposition of  $F_E$  in Lemma 2.1 prompts a better stabilization strategy. Specifically, Lemma 2.2 asserts that in the advective limit the stabilizing effect of  $F_E$  is solely due to the diffusive flux  $\Theta$ , whereas the first part of  $F_E$  is just an approximation of the (unstable) nodal flux  $F(\phi_h)$ .

Since this part is not important for the stabilization we may as well switch back to  $F(\phi_h)$ , i.e., the standard Galerkin formulation (2.1), and use  $\Theta$  to stabilize the latter. This would result in the following modification of (2.11): seek  $\phi_h \in N^h(\Omega)$  such that

$$(2.12) \quad a(\phi_h, \varphi_h) + \int_{\Omega} \Theta(\phi_h) \cdot \nabla \varphi_h dV = b(\varphi_h) \quad \forall \varphi_h \in N_0^h(\Omega).$$

However, the pairing of  $\Theta$  and  $\nabla\varphi_h$  produces a non-symmetric artificial diffusion. Numerical examples in Section 4 show that on unstructured grids this may reduce the stability of (2.12) and (2.11). Thus, we consider instead the following symmetric diffusion form:

$$(2.13) \quad Q(\phi_h, \varphi_h) = \int_{\Omega} \tilde{\Theta}_h(\phi_h) \cdot \tilde{\Theta}_h(\varphi_h) dV,$$

where

$$\tilde{\Theta}_h(\phi_h) = \sum_{e_{\xi} \in E(\Omega)} \tilde{\theta}_{\xi}(\mathbf{g}_{\xi}\phi) \vec{W}_{\xi} \quad \text{and} \quad \tilde{\theta}_{\xi} = \sqrt{\theta_{\xi}}.$$

Using (2.13) to stabilize (2.1) yields the new method, which is a ‘‘symmetrized’’ version of (2.12) and reads: seek  $\phi_h \in N^h(\Omega)$  such that

$$(2.14) \quad a(\phi_h, \varphi_h) + Q(\phi_h, \varphi_h) = b(\varphi_h) \quad \forall \varphi_h \in N_0^h(\Omega).$$

**3. Analysis.** Let  $a^h(\cdot, \cdot) := a(\cdot, \cdot) + Q(\cdot, \cdot)$  denote the bilinear form of the stabilized method (2.14). The analysis in this section utilizes a non-standard approach based on the algebraic representation of this form. In order to avoid non-essential technical details and focus instead on the key junctures of the proofs we assume that (i)  $g = 0$ , i.e., the solution of (2.14)  $\phi_h \in N_0^h(\Omega)$ , and (ii)  $\varepsilon$  and  $\mathbf{u}$  are constant on  $\Omega$ . Appendix A extends the analysis to non-constant diffusion and velocity. We first establish some useful matrix representations of the forms involved in (2.14).

LEMMA 3.1. *Let  $\phi$  and  $\varphi$  denote the coefficient vectors of  $\phi_h, \varphi_h \in N_0^h(\Omega)$ . There holds*

$$(3.1) \quad Q(\phi_h, \varphi_h) = \varphi^T (G^T D M D G) \phi$$

where  $M$  is the Gramm matrix of  $\{\vec{W}_{\xi}\}$  and  $D$  is diagonal matrix with element  $d_{\xi, \xi} = \tilde{\theta}_{\xi}$ .

*Proof.* The formula follows by factoring out the finite element coefficient vectors in (2.13)

$$Q(\phi_h, \varphi_h) = \varphi^T \left( \sum_{e_{\xi}, e_{\eta}} \mathbf{g}_{\xi}^T \tilde{\theta}_{\xi} \left( \int_{\Omega} \vec{W}_{\xi} \cdot \vec{W}_{\eta} dV \right) \tilde{\theta}_{\eta} \mathbf{g}_{\eta} \right) \phi$$

and noting that the integral above is element  $M_{\xi, \eta}$  of the Gramm matrix and  $G^T D M D G = \sum_{\xi, \eta} \mathbf{g}_{\xi}^T \tilde{\theta}_{\xi} M_{\xi, \eta} \tilde{\theta}_{\eta} \mathbf{g}_{\eta}$ , where  $\mathbf{g}_{\xi}$  and  $\mathbf{g}_{\eta}$  are the rows of  $G$  corresponding to edges  $e_{\xi}$  and  $e_{\eta}$ .  $\square$

**3.1. Stability.** We establish a baseline property of the stabilized formulation.

LEMMA 3.2. *For  $\varepsilon > 0$  the form  $a^h(\cdot, \cdot)$  defines an inner product on  $N_0^h(\Omega) \times N_0^h(\Omega)$ .*

*Proof.* Obviously,  $a^h(\cdot, \cdot)$  is a symmetric bilinear form. Therefore, we only need to show that this form is positive definite. Since  $\phi_h \in N_0^h(\Omega)$  and  $\nabla \cdot \mathbf{u} = 0$ ,

$$a^h(\phi_h, \phi_h) = \|\sqrt{\varepsilon} \nabla \phi_h\|_0^2 + Q(\phi_h, \phi_h).$$

The lemma follows from the fact that  $\nabla \phi_h = 0$  for  $\phi_h \in N_0^h(\Omega)$  iff  $\phi_h = 0$ , and

$$Q(\phi_h, \phi_h) = (M^{1/2} D G \phi)^T (M^{1/2} D G \phi) = |M^{1/2} D G \phi|^2 \geq 0,$$

i.e., the stabilizing term is non-negative.  $\square$

Lemma 3.2 implies that

$$(3.2) \quad \|\phi_h\| := \sqrt{a^h(\phi_h, \phi_h)}$$

is a mesh-dependent “energy” norm on  $N_0^h(\Omega)$ . We examine its properties in the following theorem.

**THEOREM 3.3.** *Assume that  $K(\Omega)$  is shape-regular finite element partition of  $\Omega$ , and that  $\varepsilon$  and  $\mathbf{u}$  are constant in  $\Omega$ . There exists a positive constant,  $\rho$ , independent of  $\varepsilon$ ,  $\mathbf{u}$ , and  $h$ , such that (i) for linear finite elements (on triangular or tetrahedral meshes)*

$$(3.3) \quad a^h(\phi_h, \phi_h) \geq \rho h \sum_{\mathbf{k} \in K(\Omega)} |\mathbf{u}| \|\nabla \phi_h\|_{0,\mathbf{k}}^2,$$

(ii) for bilinear and trilinear elements (on quadrilaterals and hexahedrals)

$$(3.4) \quad a^h(\phi_h, \phi_h) \geq \rho h \sum_{\mathbf{k} \in K(\Omega)} \|\mathbf{u}_{\mathbf{k}} \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2,$$

where  $\mathbf{u}_{\mathbf{k}}$  is the effective advective velocity field on element  $\mathbf{k}$ , defined in (B.4).

*Proof.* The stabilized bilinear form can be written as a sum of element forms

$$a^h(\phi_h, \varphi_h) = \sum_{\mathbf{k} \in K(\Omega)} a_{\mathbf{k}}^h(\phi_h, \varphi_h) \quad \text{where} \quad a_{\mathbf{k}}^h(\phi_h, \varphi_h) := a^h(\phi_h, \varphi_h)|_{\mathbf{k}}.$$

We will prove that for every element  $\mathbf{k}$  there is  $\rho_{\mathbf{k}} > 0$  such that

$$(3.5) \quad a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \rho_{\mathbf{k}} h_{\mathbf{k}} \begin{cases} |\mathbf{u}| \|\nabla \phi_h\|_{0,\mathbf{k}}^2 & \text{on simplices} \\ \|\mathbf{u}_{\mathbf{k}} \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 & \text{on quads and hexes} \end{cases}$$

Then, the theorem will follow with  $\rho = \min_{\mathbf{k} \in K(\Omega)} \rho_{\mathbf{k}}$  and  $h = \min_{\mathbf{k} \in K(\Omega)} h_{\mathbf{k}}$ .

Using the same arguments as in Lemmae 3.1 and 3.2 one can show that

$$(3.6) \quad a_{\mathbf{k}}^h(\phi_h, \phi_h) = \phi_{\mathbf{k}}^T (\varepsilon G_{\mathbf{k}}^T M_{\mathbf{k}} G_{\mathbf{k}} + G_{\mathbf{k}}^T D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}} G_{\mathbf{k}}) \phi_{\mathbf{k}},$$

where  $\phi_{\mathbf{k}}$  is the restriction of  $\phi$  to the nodes of  $\mathbf{k}$ ,  $M_{\mathbf{k}}$  is the Gramm matrix of the edge element basis on  $\mathbf{k}$ , and  $G_{\mathbf{k}}$  and  $D_{\mathbf{k}}$  are matrices containing the rows of  $G$  and  $D$ , respectively, corresponding to the edges of  $\mathbf{k}$ . Similar to  $G$ , the “topological” element gradient matrix  $G_{\mathbf{k}}$  has a one-dimensional nullspace spanned by the constant vector. As a result, if  $\phi_h|_{\mathbf{k}} = \text{const}$ , then  $a_{\mathbf{k}}^h(\phi_h, \phi_h) = 0$ ,  $\nabla \phi_h|_{\mathbf{k}} = 0$ , and (3.5) is trivially satisfied on  $\mathbf{k}$ . Thus, without a loss of generality we may assume that  $\phi_h|_{\mathbf{k}} \neq \text{const}$ , or what is the same –  $G_{\mathbf{k}} \phi_{\mathbf{k}} \neq 0$ .

*Algebraic bounds for element forms.* Let  $\mathbf{z}_{\mathbf{k}} = G_{\mathbf{k}} \phi_{\mathbf{k}}$  and  $\mathbf{y}_{\mathbf{k}} = D_{\mathbf{k}} G_{\mathbf{k}} \phi_{\mathbf{k}}$ . Using (3.6)

$$\begin{aligned} a_{\mathbf{k}}^h(\phi_h, \phi_h) &= \varepsilon \mathbf{z}_{\mathbf{k}}^T M_{\mathbf{k}} \mathbf{z}_{\mathbf{k}} + \mathbf{y}_{\mathbf{k}}^T M_{\mathbf{k}} \mathbf{y}_{\mathbf{k}} \geq \underline{\lambda}(M_{\mathbf{k}}) \left( \varepsilon \mathbf{z}_{\mathbf{k}}^T \mathbf{z}_{\mathbf{k}} + \mathbf{y}_{\mathbf{k}}^T \mathbf{y}_{\mathbf{k}} \right) \\ &= \underline{\lambda}(M_{\mathbf{k}}) \left( \varepsilon \mathbf{z}_{\mathbf{k}}^T \mathbf{z}_{\mathbf{k}} + \mathbf{z}_{\mathbf{k}}^T D_{\mathbf{k}}^2 \mathbf{z}_{\mathbf{k}} \right) = \underline{\lambda}(M_{\mathbf{k}}) \sum_{\mathbf{e}_{\xi} \in E(\mathbf{k})} (\bar{\varepsilon}_{\xi} + \theta_{\xi}) z_{\xi}^2 = \underline{\lambda}(M_{\mathbf{k}}) \sum_{\mathbf{e}_{\xi} \in E(\mathbf{k})} \sigma_{\xi} z_{\xi}^2, \end{aligned}$$

where  $\sigma_\xi$  is the number defined in (2.6). From the lower bound (2.10) in Lemma 2.2

$$a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \underline{\lambda}(M_{\mathbf{k}}) \sum_{e_\xi \in E(\mathbf{k})} \frac{h_\xi |\bar{u}_\xi|}{2} z_\xi^2.$$

Because  $\mathbf{u}$  is constant  $\bar{u}_\xi = |\mathbf{u}|v_\xi$ , where  $v_\xi = (\mathbf{u}/|\mathbf{u}|) \cdot \mathbf{t}_\xi$ . Furthermore, there exists a constant  $C_{\mathbf{k}}$  such that  $0 < C_{\mathbf{k}} h_{\mathbf{k}} \leq h_\xi$  for every edge  $e_\xi \in E(\mathbf{k})$ . It follows that

$$(3.7) \quad a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \underline{\lambda}(M_{\mathbf{k}}) C_{\mathbf{k}} \frac{h_{\mathbf{k}} |\mathbf{u}|}{2} \sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi|.$$

We estimate the sum on the right hand side of (3.7) separately for simplices (triangles and tets) and tensor product elements (quads and hexes).

*Simplicial elements.* On triangles  $\mathbf{u} \cdot \mathbf{t}_\xi = 0$  for at most one edge  $e_\alpha \in E(\mathbf{k})$ , in which case  $v_\alpha = 0$ . On tetrahedrons  $\mathbf{u}$  can be orthogonal to either 3 planar edges  $e_\alpha, e_\beta$  and  $e_\gamma$ , or at most two non-coplanar edges  $e_\alpha$  and  $e_\beta$ . In the first case  $v_\alpha = v_\beta = v_\gamma = 0$  and in the second case  $v_\alpha = v_\beta = 0$ . In general, let  $e_\alpha, e_\beta \in E(\mathbf{k})$  be the edges corresponding to the smallest two values  $0 \leq |v_\alpha| \leq |v_\beta|$  on element  $\mathbf{k}$ . If  $\mathbf{k}$  is a triangle, define the set  $E(\mathbf{k}/\mathbf{u})$  to be the set containing the edge  $e_\alpha$ . On tetrahedrons, define this set as follows. If  $e_\alpha$  and  $e_\beta$  are coplanar, set  $E(\mathbf{k}/\mathbf{u}) = \{e_\alpha, e_\beta, e_\gamma\}$ , where  $e_\gamma$  is the edge coplanar to  $e_\alpha$  and  $e_\beta$ . If  $e_\alpha$  and  $e_\beta$  are not coplanar, then define  $E(\mathbf{k}/\mathbf{u}) = \{e_\alpha, e_\beta\}$ . Finally, let  $v_{\mathbf{k}} = \min_{e_\xi \notin E(\mathbf{k}/\mathbf{u})} |v_\xi|$ . Then<sup>1</sup>,

$$\sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq v_{\mathbf{k}} \sum_{e_\xi \notin E(\mathbf{k}/\mathbf{u})} z_\xi^2 = v_{\mathbf{k}} \phi_{\mathbf{k}}^T \left( \sum_{e_\xi \notin E(\mathbf{k}/\mathbf{u})} \mathbf{g}_\xi^T \mathbf{g}_\xi \right) \phi_{\mathbf{k}}.$$

It is easy to check that for all possible configurations of  $E(\mathbf{k}/\mathbf{u})$  the matrix  $\sum_{e_\xi \notin E(\mathbf{k}/\mathbf{u})} \mathbf{g}_\xi^T \mathbf{g}_\xi$  is positive semidefinite with a single zero eigenvalue corresponding to the constant eigenvector, a case that is ruled out by the assumption  $\phi_{\mathbf{k}} \neq \text{const}$ . Its smallest eigenvalue is either 2, when  $\mathbf{k}$  is a tetrahedron and  $E(\mathbf{k}/\mathbf{u})$  contains two non coplanar edges, or 1 in all other cases. As a result,

$$(3.8) \quad \sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq v_{\mathbf{k}} |\phi_{\mathbf{k}}|^2 \quad \text{and} \quad a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \underline{\lambda}(M_{\mathbf{k}}) C_{\mathbf{k}} \frac{h_{\mathbf{k}} |\mathbf{u}|}{2} v_{\mathbf{k}} |\phi_{\mathbf{k}}|^2.$$

The next step is to use (1.5) to bound the norm of  $\nabla \phi_h$  by  $|\phi_{\mathbf{k}}|^2$ :

$$(3.9) \quad \|\nabla \phi_h\|_{0,\mathbf{k}}^2 = \mathbf{z}_{\mathbf{k}}^T M_{\mathbf{k}} \mathbf{z}_{\mathbf{k}} \leq \bar{\lambda}(M_{\mathbf{k}}) |\mathbf{z}_{\mathbf{k}}|^2 = \bar{\lambda}(M_{\mathbf{k}}) \phi_{\mathbf{k}}^T (G_{\mathbf{k}}^T G_{\mathbf{k}}) \phi_{\mathbf{k}} \leq \bar{\lambda}(M_{\mathbf{k}}) \bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) |\phi_{\mathbf{k}}|^2.$$

Taking into account that  $\bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}})$  equals 3 for triangles and 4 for tets, (3.9) yields

$$|\phi_{\mathbf{k}}|^2 \geq \frac{1}{4\bar{\lambda}(M_{\mathbf{k}})} \|\nabla \phi_h\|_{0,\mathbf{k}}^2.$$

Combining this estimate with (3.8) shows that on simplicial elements

$$a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \rho_{\mathbf{k}} h_{\mathbf{k}} |\mathbf{u}| \|\nabla \phi_h\|_{0,\mathbf{k}}^2, \quad \text{with} \quad \rho_{\mathbf{k}} = \frac{1}{8} \kappa^{-1} (M_{\mathbf{k}}) C_{\mathbf{k}} v_{\mathbf{k}}.$$

REMARK 1. *The value of  $v_{\mathbf{k}}$  depends on the direction of  $\mathbf{u}$  with respect to the element edges. For example, on a triangular grid  $v_{\mathbf{k}} \geq \sin\left(\frac{\zeta_{\mathbf{k}}}{2}\right)$ , where  $\zeta_{\mathbf{k}}$  is the smallest angle in  $K(\Omega)$ .*

<sup>1</sup>We remind that  $\mathbf{g}_\xi$  is the row of  $G_{\mathbf{k}}$  corresponding to edge  $e_\xi \in E(\mathbf{k})$ .



*Quadrilateral elements.* We refer to Appendix B for the relevant notations. According to (B.3)  $|v_i| \leq |v_{\alpha^i}|$  and  $|v_j| \leq |v_{\beta^i}|$ . It follows that

$$(3.10) \quad \sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq |v_i| (z_{\alpha^1}^2 + z_{\alpha^2}^2) + |v_j| (z_{\beta^1}^2 + z_{\beta^2}^2) = |v_i| \phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} + |v_j| \phi_{\mathbf{k}}^T G_\beta \phi_{\mathbf{k}}$$

where  $G_\alpha = \mathbf{g}_{\alpha^1}^T \mathbf{g}_{\alpha^1} + \mathbf{g}_{\alpha^2}^T \mathbf{g}_{\alpha^2}$  and  $G_\beta = \mathbf{g}_{\beta^1}^T \mathbf{g}_{\beta^1} + \mathbf{g}_{\beta^2}^T \mathbf{g}_{\beta^2}$ . We bound (3.10) from below by the norm of the streamline derivative  $\mathbf{u}_{\mathbf{k}} \cdot \nabla \phi_h$  along the *effective advective velocity*. Since the terms involving  $G_\alpha$  and  $G_\beta$  have the same structure we give the details for the first term only.

*Step 1. Algebraic lower bound for  $\phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}}$ .* Let  $\{\mathbf{q}_i\}$  be an orthonormal basis for the (two-dimensional) kernel of  $G_\alpha$  and consider the orthogonal decomposition of the coefficient vector  $\phi_{\mathbf{k}}$ :

$$\phi_{\mathbf{k}} = \phi_{\mathbf{q}}^\perp + \phi_{\mathbf{q}} \quad \text{with} \quad \phi_{\mathbf{q}}^\perp = \left( I - \sum_{i=1}^2 \mathbf{q}_i \mathbf{q}_i^T \right) \phi_{\mathbf{k}}.$$

It is easy to check that all non-zero eigenvalues of  $G_\alpha$  equal 2, and so,

$$(3.11) \quad \phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} = (\phi_{\mathbf{q}}^\perp + \phi_{\mathbf{q}})^T G_\alpha (\phi_{\mathbf{q}}^\perp + \phi_{\mathbf{q}}) = (\phi_{\mathbf{q}}^\perp)^T G_\alpha \phi_{\mathbf{q}}^\perp \geq 2 |\phi_{\mathbf{q}}^\perp|^2.$$

*Step 2. Lower bound for  $|\phi_{\mathbf{q}}^\perp|$ .* We obtain an analogue of (3.9) for tensor product elements. Let  $\phi_{\mathbf{q},h}^\perp$  be the finite element function with coefficient vector  $\phi_{\mathbf{q}}^\perp$ . Then,

$$(3.12) \quad \|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}^2 = (\phi_{\mathbf{q}}^\perp)^T G_{\mathbf{k}}^T M_{\mathbf{k}} G_{\mathbf{k}} \phi_{\mathbf{q}}^\perp \leq \bar{\lambda}(M_{\mathbf{k}}) |G_{\mathbf{k}} \phi_{\mathbf{q}}^\perp|^2 \leq \bar{\lambda}(M_{\mathbf{k}}) \bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) |\phi_{\mathbf{q}}^\perp|^2.$$

Using (3.12) in conjunction with the fact that  $\bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) = 4$  for quads yields

$$(3.13) \quad |\phi_{\mathbf{q}}^\perp|^2 \geq \frac{1}{4\bar{\lambda}(M_{\mathbf{k}})} \|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}^2.$$

*Step 3. Lower bound for  $\|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}$ .* Bounding of the element forms for tensor product elements requires one more step that was not necessary for simplicial elements. This step estimates  $\|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}$  by a norm of the directional derivative  $\mathbf{t}_i \cdot \nabla \phi_h$  of the finite element function along the vector field aligned with the edges  $e_{\alpha^i}$ . Using (1.4) and the formula for  $\phi_{\mathbf{q}}^\perp$  yields

$$\nabla \phi_{\mathbf{q},h}^\perp = \sum_{e_\xi \in E(\mathbf{k})} \left( \mathbf{g}_\xi \phi_{\mathbf{k}} - \mathbf{g}_\xi \sum_{i=1}^2 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_\xi = \nabla \phi_h - \sum_{e_\xi \in E(\mathbf{k})} \mathbf{g}_\xi \left( \sum_{i=1}^2 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_\xi$$

Since  $\mathbf{q}_i \in \ker G_\alpha$ ,  $\mathbf{g}_{\alpha^j} \mathbf{q}_i = 0$  and summation over  $E(\mathbf{k})$  reduces to a sum over  $e_{\beta^j}$  only:

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - \sum_{j=1}^2 \mathbf{g}_{\beta^j} \left( \sum_{i=1}^2 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_{\beta^j}.$$

Using the expression (B.5) for the edge element basis functions we find that

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - \left[ \sum_{j=1}^2 \mathbf{g}_{\beta^j} \left( \sum_{i=1}^2 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) B_j(\mathbf{x}) \right] \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j} = \nabla \phi_h - r(\mathbf{x}) \frac{\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j}}{|\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j}|},$$

where  $r(\mathbf{x})$  is a scalar function corresponding to the expression in the square brackets scaled by  $|\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}|$ . On the other hand, from (B.1) it follows that

$$\frac{(\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j})^T}{|\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}|} \mathbf{t}_i = \frac{\mathbf{j}^T (\mathbf{J}_k^{-1}(\mathbf{x})\mathbf{J}_k(\mathbf{x})) \mathbf{i}}{|\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}| |\mathbf{J}_k(\mathbf{x})\mathbf{i}|} = \frac{\mathbf{j}^T \mathbf{i}}{|\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}| |\mathbf{J}_k(\mathbf{x})\mathbf{i}|} = 0 \quad \forall \mathbf{x} \in \mathbf{k}.$$

i.e.,  $\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}/|\mathbf{J}_k^{-T}(\mathbf{x})\mathbf{j}|$  is perpendicular to  $\mathbf{t}_i$ . Accordingly, we denote the former by  $\mathbf{t}_i^\perp$  and so,

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - r(\mathbf{x})\mathbf{t}_i^\perp.$$

Since  $\mathbf{t}_i$  and  $\mathbf{t}_i^\perp$  are orthonormal we have the orthogonal decomposition

$$(3.14) \quad \nabla \phi_h = (\mathbf{t}_i \cdot \nabla \phi_h) \mathbf{t}_i + (\mathbf{t}_i^\perp \cdot \nabla \phi_h) \mathbf{t}_i^\perp.$$

The second term is the orthogonal projection of the gradient onto  $\mathbf{t}_i^\perp$ . From the properties of orthogonal projections and (3.14) it follows that

$$(3.15) \quad |\nabla \phi_{\mathbf{q},h}^\perp| = |\nabla \phi_h - r(\mathbf{x})\mathbf{t}_i^\perp| \geq |\nabla \phi_h - (\mathbf{t}_i^\perp \cdot \nabla \phi_h) \mathbf{t}_i^\perp| = |(\mathbf{t}_i \cdot \nabla \phi_h) \mathbf{t}_i| \quad \forall \mathbf{x} \in \mathbf{k}.$$

Combining this result with the lower bounds in (3.11) and (3.13) yields

$$\phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} \geq \frac{1}{2\bar{\lambda}(M_{\mathbf{k}})} \|(\mathbf{t}_i \cdot \nabla \phi_h) \mathbf{t}_i\|_{0,\mathbf{k}}^2.$$

Repeating the same steps for the second term in (3.10) leads to

$$\phi_{\mathbf{k}}^T G_\beta \phi_{\mathbf{k}} \geq \frac{1}{2\bar{\lambda}(M_{\mathbf{k}})} \|(\mathbf{t}_j \cdot \nabla \phi_h) \mathbf{t}_j\|_{0,\mathbf{k}}^2.$$

Recalling that  $\mathbf{t}_i$  and  $\mathbf{t}_j$  are unit vectors for every  $\mathbf{x} \in \mathbf{k}$  gives the intermediate lower bound

$$\sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq \frac{1}{2\bar{\lambda}(M_{\mathbf{k}})} \left( |v_i| \|\mathbf{t}_i \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 + |v_j| \|\mathbf{t}_j \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 \right).$$

Since  $|v_i| \leq 1$ , we have that  $|v_i|^2 \leq |v_i|$  and the same for  $v_j$ . Using this, and the triangle inequality

$$\begin{aligned} |v_i| \|\mathbf{t}_i \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 + |v_j| \|\mathbf{t}_j \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 &\geq \|v_i \mathbf{t}_i \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 + \|v_j \mathbf{t}_j \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 \\ &\geq \|v_i \mathbf{t}_i \cdot \nabla \phi_h + v_j \mathbf{t}_j \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 = \|(\mathbf{t}_i v_i + \mathbf{t}_j v_j) \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 = \frac{1}{|\mathbf{u}|} \|\mathbf{u}_k \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2. \end{aligned}$$

As a result,

$$\sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq \frac{1}{2\bar{\lambda}(M_{\mathbf{k}})|\mathbf{u}|} \|\mathbf{u}_k \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2.$$

Together with (3.7) this inequality allows us to conclude that on quads

$$a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \rho_{\mathbf{k}} h_{\mathbf{k}} \|\mathbf{u}_k \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 \quad \text{with} \quad \rho_{\mathbf{k}} = \frac{1}{4} C_{\mathbf{k}} \kappa^{-1}(M_{\mathbf{k}}).$$

*Hexahedral elements.* Estimation of the element forms on hexahedra follows the same key steps. The analogue of (3.10) is given by

$$(3.16) \quad \sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq |v_i| \phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} + |v_j| \phi_{\mathbf{k}}^T G_\beta \phi_{\mathbf{k}} + |v_{\mathbf{k}}| \phi_{\mathbf{k}}^T G_\gamma \phi_{\mathbf{k}}.$$

where  $G_\alpha = \sum_{i=1}^4 \mathbf{g}_{\alpha^i}^T \mathbf{g}_{\alpha^i}$ ,  $G_\beta = \sum_{i=1}^4 \mathbf{g}_{\beta^i}^T \mathbf{g}_{\beta^i}$ , and  $G_\gamma = \sum_{i=1}^4 G_{\gamma^i}^T G_{\gamma^i}$  have four-dimensional nullspaces. We briefly discuss bounding the first term on the right hand side of (3.16).

*Step 1. Algebraic lower bound for  $\phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}}$ .* Let  $\{\mathbf{q}_i\}$  denote an orthonormal basis of  $\ker G_\alpha$ . For hexahedrons all nonzero eigenvalues of  $G_\alpha$  equal 2 and so, (3.11) continues to hold:

$$(3.17) \quad \phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} \geq 2 |\phi_{\mathbf{q}}^\perp|^2.$$

*Step 2. Lower bound for  $|\phi_{\mathbf{q}}^\perp|$ .* On hexahedra  $\bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) = 6$  and the analogue of (3.13) is

$$(3.18) \quad |\phi_{\mathbf{q}}^\perp|^2 \geq \frac{1}{6\bar{\lambda}(M_{\mathbf{k}})} \|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}^2.$$

*Step 3. Lower bound for  $\|\nabla \phi_{\mathbf{q},h}^\perp\|_{0,\mathbf{k}}$ .* As for quads, we have that

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - \sum_{e_\xi \in E(\mathbf{k})} \mathbf{g}_\xi \left( \sum_{i=1}^4 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_\xi.$$

Likewise, since  $\mathbf{q}_i \in \ker G_\alpha$ ,  $\mathbf{g}_{\alpha^j} \mathbf{q}_i = 0$ , summation over  $E(\mathbf{k})$  reduces to a sum over  $e_{\beta^j}$  and  $e_{\gamma^j}$

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - \sum_{j=1}^4 \mathbf{g}_{\beta^j} \left( \sum_{i=1}^4 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_{\beta^j} - \sum_{j=1}^4 \mathbf{g}_{\gamma^j} \left( \sum_{i=1}^4 (\mathbf{q}_i^T \phi_{\mathbf{k}}) \mathbf{q}_i \right) \vec{W}_{\gamma^j}$$

Using (B.6)

$$\nabla \phi_{\mathbf{q},h}^\perp = \nabla \phi_h - \left( r_{\mathbf{j}}(\mathbf{x}) \frac{\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j}}{|\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j}|} + r_{\mathbf{k}}(\mathbf{x}) \frac{\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{k}}{|\mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{k}|} \right)$$

For every  $\mathbf{x} \in \mathbf{k}$  the vector field in the parenthesis belongs in the orthogonal complement of the vector field  $\mathbf{t}_i$ . Using properties of orthogonal projections

$$|\nabla \phi_{\mathbf{q},h}^\perp| \geq |\nabla \phi_h - (\mathbf{t}_j \cdot \nabla \phi_h) \mathbf{t}_j - (\mathbf{t}_k \cdot \nabla \phi_h) \mathbf{t}_k| = |(\mathbf{t}_i \cdot \nabla \phi_h) \mathbf{t}_i| \quad \forall \mathbf{x} \in \mathbf{k}.$$

Repeating the same steps for all three terms leads to

$$\begin{aligned} \phi_{\mathbf{k}}^T G_\alpha \phi_{\mathbf{k}} &\geq \frac{1}{3\bar{\lambda}(M_{\mathbf{k}})} \|(\mathbf{t}_i \cdot \nabla \phi_h) \mathbf{t}_i\|_{0,\mathbf{k}}^2; & \phi_{\mathbf{k}}^T G_\beta \phi_{\mathbf{k}} &\geq \frac{1}{3\bar{\lambda}(M_{\mathbf{k}})} \|(\mathbf{t}_j \cdot \nabla \phi_h) \mathbf{t}_j\|_{0,\mathbf{k}}^2 \\ \text{and } \phi_{\mathbf{k}}^T G_\gamma \phi_{\mathbf{k}} &\geq \frac{1}{3\bar{\lambda}(M_{\mathbf{k}})} \|(\mathbf{t}_k \cdot \nabla \phi_h) \mathbf{t}_k\|_{0,\mathbf{k}}^2. \end{aligned}$$

Proceeding as in the case of quadrilaterals we find that

$$\sum_{e_\xi \in E(\mathbf{k})} z_\xi^2 |v_\xi| \geq \frac{1}{3\bar{\lambda}(M_{\mathbf{k}})} \|(v_i \mathbf{t}_i + v_j \mathbf{t}_j + v_{\mathbf{k}} \mathbf{t}_{\mathbf{k}}) \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 = \frac{1}{3\bar{\lambda}(M_{\mathbf{k}})|\mathbf{u}|} \|\mathbf{u}_{\mathbf{k}} \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2,$$

which produces an analogous lower bound for the element form:

$$a_{\mathbf{k}}^h(\phi_h, \phi_h) \geq \rho_{\mathbf{k}} h_{\mathbf{k}} \|\mathbf{u}_{\mathbf{k}} \cdot \nabla \phi_h\|_{0,\mathbf{k}}^2 \quad \text{with} \quad \rho_{\mathbf{k}} = \frac{1}{6} C_{\mathbf{k}} \kappa^{-1}(M_{\mathbf{k}}).$$

□

Theorem 3.3 implies that (2.14) remains stable for  $\varepsilon \ll h$ . In conjunction with the Lax-Milgram lemma, the theorem allows us to conclude that the stabilized method (2.14) has a unique solution.

**3.2. Error estimates.** We start with an upper bound on the artificial diffusion term  $Q(\cdot, \cdot)$ .

LEMMA 3.4. *There exists a constant  $C$ , independent of  $h$ , such that for all  $\phi_h \in N^h(\Omega)$*

$$(3.19) \quad Q(\phi_h, \phi_h) \leq C h |\mathbf{u}| \|\nabla \phi_h\|_0^2.$$

*Proof.* The stabilizing term  $Q(\cdot, \cdot)$  can be written as a sum of element forms  $Q_{\mathbf{k}}(\phi_h, \varphi_h) := Q(\phi_h, \varphi_h)|_{\mathbf{k}}$ . We use the techniques of Theorem 3.3 to bound these forms by  $\|\nabla \phi_h\|_{0,\mathbf{k}}^2$ . As in the proof of this theorem it suffices to consider only those  $\phi_h$  for which  $G_{\mathbf{k}} \phi_{\mathbf{k}} \neq 0$ . Using the upper bound for  $\theta_\alpha$  in (2.10) and the fact that  $\bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) \leq 6$  for all elements under consideration

$$Q_{\mathbf{k}}(\phi_h, \phi_h) = \phi_{\mathbf{k}}^T (G_{\mathbf{k}}^T D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}} G_{\mathbf{k}}) \phi_{\mathbf{k}}^T \leq \bar{\lambda}(M_{\mathbf{k}}) \bar{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) \frac{h_{\mathbf{k}} |\mathbf{u}|}{2} |\phi_{\mathbf{k}}|^2 \leq 3\bar{\lambda}(M_{\mathbf{k}}) h_{\mathbf{k}} |\mathbf{u}| |\phi_{\mathbf{k}}|^2.$$

On the other hand, since the smallest non-zero eigenvalue  $\underline{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) \geq 1$  for all element shapes

$$\|\nabla \phi_h\|_{0,\mathbf{k}}^2 = \phi_{\mathbf{k}}^T (G_{\mathbf{k}}^T M_{\mathbf{k}} G_{\mathbf{k}}) \phi_{\mathbf{k}}^T \geq \underline{\lambda}(M_{\mathbf{k}}) \underline{\lambda}(G_{\mathbf{k}}^T G_{\mathbf{k}}) |\phi_{\mathbf{k}}|^2 \geq \underline{\lambda}(M_{\mathbf{k}}) |\phi_{\mathbf{k}}|^2.$$

Therefore, for every element in  $K(\Omega)$

$$Q_{\mathbf{k}}(\phi_h, \phi_h) \leq C_{\mathbf{k}} h_{\mathbf{k}} |\mathbf{u}| \|\nabla \phi_h\|_{0,\mathbf{k}}^2 \quad \text{with} \quad C_{\mathbf{k}} = 3\kappa(M_{\mathbf{k}}),$$

and the Lemma follows with  $C = \max_{\mathbf{k}} C_{\mathbf{k}}$ . □

We proceed to estimate the discrete error of the solution of (2.14).

THEOREM 3.5. *Assume that the exact solution of (1.1) is in  $H^2(\Omega)$ ,  $\phi_h$  is the solution of (2.14), and  $\phi_{\mathcal{I}} := \mathcal{I}_N(\phi)$ . There exists a constant  $C$ , independent of  $h$  and  $\varepsilon$ , such that*

$$(3.20) \quad \|\phi_h - \phi_{\mathcal{I}}\| \leq \gamma_1 \sqrt{h} \|\phi\|_2.$$

*Proof.* We have the following “error orthogonality” relation

$$a(\phi, \varphi_h) = b(\varphi_h) = a^h(\phi_h, \varphi_h) \quad \forall \varphi_h \in N_0^h(\Omega).$$

This identity, the fact that  $a^h(\cdot, \cdot) = a(\cdot, \cdot) + Q(\cdot, \cdot)$ , and the Cauchy-Schwartz inequality imply that

$$(3.21) \quad \begin{aligned} \|\phi_h - \phi_{\mathcal{I}}\|^2 &= a^h(\phi_h - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) \pm a(\phi - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) \\ &= a(\phi - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) + [a(\phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) - a^h(\phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}})] \\ &= a(\phi - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) + Q(\phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}}) \\ &\leq a(\phi - \phi_{\mathcal{I}}, \phi - \phi_{\mathcal{I}})^{1/2} a(\phi_h - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}})^{1/2} + Q(\phi_{\mathcal{I}}, \phi_{\mathcal{I}})^{1/2} Q(\phi_h - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}})^{1/2}. \end{aligned}$$

Using the interpolation estimate (1.9) and the skew-symmetry of the advection term

$$a(\phi - \phi_{\mathcal{I}}, \phi - \phi_{\mathcal{I}})^{1/2} \leq \sqrt{\varepsilon} \|\nabla\phi - \nabla\phi_{\mathcal{I}}\|_0 \leq \sqrt{\varepsilon} h \|\phi\|_2,$$

while from identity  $a^h(\cdot, \cdot) = a(\cdot, \cdot) + Q(\cdot, \cdot)$  it follows that

$$a(\phi_h - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}})^{1/2} \leq \|\phi_h - \phi_{\mathcal{I}}\| \quad \text{and} \quad Q(\phi_h - \phi_{\mathcal{I}}, \phi_h - \phi_{\mathcal{I}})^{1/2} \leq \|\phi_h - \phi_{\mathcal{I}}\|.$$

Finally, using (3.19) in Lemma 3.4 and the boundedness (1.10) of  $\mathcal{I}_N$  we find that

$$Q(\phi_{\mathcal{I}}, \phi_{\mathcal{I}})^{1/2} \leq C \sqrt{h} |\mathbf{u}| \|\nabla\phi_{\mathcal{I}}\|_0 \leq C \sqrt{h} |\mathbf{u}| \|\phi\|_2.$$

Inserting these bounds into (3.21) and dividing through by the energy norm of  $\phi_h - \phi_{\mathcal{I}}$  yields

$$\|\phi_h - \phi_{\mathcal{I}}\| \leq \sqrt{\varepsilon} h \|\phi\|_2 + C \sqrt{h} |\mathbf{u}| \|\phi\|_2 = C \sqrt{h} \left( \sqrt{\varepsilon} h + |\mathbf{u}| \right) \|\phi\|_2,$$

which completes the proof.  $\square$

**COROLLARY 3.6.** *Under the assumptions of Theorem 3.5 there holds*

$$(3.22) \quad \sqrt{\varepsilon} \|\nabla\phi - \nabla\phi_h\|_0 \leq C \sqrt{h} \|\phi\|_2$$

with a constant  $C$  that does not depend on  $h$  and  $\varepsilon$ .

*Proof.* We use triangle inequality to split the error

$$\sqrt{\varepsilon} \|\nabla\phi - \nabla\phi_h\|_0 \leq \sqrt{\varepsilon} \|\nabla\phi - \nabla\phi_{\mathcal{I}}\|_0 + \sqrt{\varepsilon} \|\nabla\phi_{\mathcal{I}} - \nabla\phi_h\|_0.$$

The assertion follows by noting that

$$\varepsilon \|\nabla\phi_{\mathcal{I}} - \nabla\phi_h\|_0^2 = a(\phi_{\mathcal{I}} - \phi_h, \phi_{\mathcal{I}} - \phi_h) \leq a^h(\phi_{\mathcal{I}} - \phi_h, \phi_{\mathcal{I}} - \phi_h) = \|\phi_{\mathcal{I}} - \phi_h, \phi_{\mathcal{I}} - \phi_h\|^2$$

and applying the result of Theorem 3.5.  $\square$

Let  $e_h := \phi - \phi_h$ . To estimate the  $L^2$  error of the stabilized solution we assume that the adjoint problem: find  $\psi \in H_0^1(\Omega)$  such that

$$(3.23) \quad a(\varphi, \psi) = (e_h, \varphi), \quad \forall \varphi \in H_0^1(\Omega),$$

has full elliptic regularity, i.e., there is  $\gamma_\varepsilon > 0$  such that

$$(3.24) \quad \|\psi\|_2 \leq \gamma_\varepsilon \|e_h\|_0.$$

**THEOREM 3.7.** *Under the hypotheses of Theorem 3.5 there exists  $C(\varepsilon)$  such that*

$$(3.25) \quad \|\phi_h - \phi\|_0 \leq C(\varepsilon) h \|\phi\|_2.$$

*Proof.* Let  $\psi_{\mathcal{I}}$  be the interpolant of the solution of the adjoint problem and  $\phi_h$  the stabilized solution. Setting  $\varphi = e_h$  in (3.23) and using that  $a(\phi, \psi_{\mathcal{I}}) = a^h(\phi_h, \psi_{\mathcal{I}})$  gives

$$\begin{aligned} \|e_h\|_0^2 &= a(e_h, \psi) \pm a(e_h, \psi_{\mathcal{I}}) \\ &= a(\phi - \phi_h, \psi - \psi_{\mathcal{I}}) + a(\phi - \phi_h, \psi_{\mathcal{I}}) \\ &= a(\phi - \phi_h, \psi - \psi_{\mathcal{I}}) + (a^h(\phi_h, \psi_{\mathcal{I}}) - a(\phi_h, \psi_{\mathcal{I}})) \\ &= a(\phi - \phi_h, \psi - \psi_{\mathcal{I}}) - Q(\phi_h, \psi_{\mathcal{I}}). \end{aligned}$$

Using the Cauchy-Schwartz inequality for the first term gives the bound

$$a(\phi - \phi_h, \psi - \psi_{\mathcal{I}}) \leq \varepsilon \|\nabla e_h\|_0 \|\nabla \psi - \nabla \psi_{\mathcal{I}}\|_0 + |\mathbf{u}| \|e_h\|_0 \|\nabla \psi - \nabla \psi_{\mathcal{I}}\|_0$$

Using interpolation theory (1.9) and elliptic regularity (3.24)

$$\|\nabla \psi - \nabla \psi_{\mathcal{I}}\|_0 \leq Ch \|\psi\|_2 \leq Ch \|e_h\|_0,$$

while from Corollary 3.6  $\sqrt{\varepsilon} \|\nabla e_h\|_0 \leq C\sqrt{h} \|\phi\|_2$ . Therefore,

$$a(\phi - \phi_h, \psi - \psi_{\mathcal{I}}) \leq C \left( h^{3/2} \sqrt{\varepsilon} \|\phi\|_2 \|e_h\|_0 + h |\mathbf{u}| \|e_h\|_0^2 \right).$$

Using the Cauchy-Schwartz inequality, (3.19) in Lemma 3.4, (1.10) and elliptic regularity (3.24),

$$\begin{aligned} Q(\phi_h, \psi_{\mathcal{I}}) &\leq Q(\phi_h, \phi_h)^{1/2} Q(\psi_{\mathcal{I}}, \psi_{\mathcal{I}})^{1/2} \\ &\leq Ch |\mathbf{u}| \|\nabla \phi_h\|_0 \|\nabla \psi_{\mathcal{I}}\|_0 \leq Ch |\mathbf{u}| \|\nabla \phi_h\|_0 \|\psi\|_2 \leq Ch |\mathbf{u}| \|\nabla \phi_h\|_0 \|e_h\|_0. \end{aligned}$$

This gives the following intermediate upper bound for  $\|e_h\|_0$ :

$$\|e_h\|_0^2 \leq C \left( h^{3/2} \sqrt{\varepsilon} \|\phi\|_2 \|e_h\|_0 + h |\mathbf{u}| \|e_h\|_0^2 + h |\mathbf{u}| \|\nabla \phi_h\|_0 \|e_h\|_0 \right)$$

For  $h$  small enough,  $Ch |\mathbf{u}| < 1/2$  and  $Ch |\mathbf{u}| \|e_h\|_0^2$  can be absorbed into the left hand side. After dividing the result by  $\|e_h\|_0$  the intermediate bound assumes the form

$$\|e_h\|_0 \leq Ch \left( \sqrt{\varepsilon h} \|\phi\|_2 + |\mathbf{u}| \|\nabla \phi_h\|_0 \right).$$

Using Corollary 3.6, the regularity assumption  $\phi \in H^2(\Omega)$  and the triangle inequality,

$$\|\nabla \phi_h\|_0 \leq \|\nabla \phi_h - \nabla \phi\|_0 + \|\nabla \phi\|_0 \leq \frac{1}{\sqrt{\varepsilon}} (\sqrt{\varepsilon} \|\nabla \phi_h - \nabla \phi\|_0) + \|\phi\|_2 \leq \left( C\sqrt{\frac{h}{\varepsilon}} + 1 \right) \|\phi\|_2.$$

Combining all results together shows that for  $h$  small enough

$$\|e_h\|_0 \leq Ch \left( \sqrt{\varepsilon h} + |\mathbf{u}| \left( C\sqrt{\frac{h}{\varepsilon}} + 1 \right) \right) \|\phi\|_2,$$

which completes the proof.  $\square$

**4. Computational study.** This section provides a brief numerical illustration of the “symmetrized” stabilized method (2.14). We refer to [3] and [4] for thorough numerical studies of the finite element method (2.11) and the related control volume finite element method, respectively, both of which utilize the whole edge element flux  $F_E$ .

We solve (1.1) on the unit square  $\Omega = [0, 1]^2$  using conforming partitions  $K_h(\Omega)$  of  $\Omega$  into quadrilateral and triangular elements. Thus,  $N^h(\Omega)$  is either the isoparametric bilinear finite element space on quads, or the affine piecewise linear finite element on triangles. The domain boundary  $\Gamma = \Gamma_B \cup \Gamma_T \cup \Gamma_L \cup \Gamma_R$ , where  $\Gamma_B$ ,  $\Gamma_T$ ,  $\Gamma_L$  and  $\Gamma_R$  are the bottom, top, left and right sides of  $\Omega$ , respectively.

**4.1. Convergence rates.** The objective is to illustrate the theoretical error estimates (3.20), (3.22) and (3.25) in Section 3.2. To this end, we use the manufactured solution  $\phi = x^3 - y^2$  and the velocity field  $\mathbf{u} = (-\sin \pi/6, \cos \pi/6)$  from [8, Example 3.1.3, p.118]. Substitution of  $\phi$  and  $\mathbf{u}$  into (1.1) defines the boundary data and the forcing term.

Convergence rates of the discrete error  $\|\phi_{\mathcal{T}} - \phi_h\|$ , the  $H^1$ -seminorm error  $\|\nabla\phi - \nabla\phi_h\|_0$  and the  $L^2$ -norm error  $\|\phi - \phi_h\|$  are estimated by solving (2.14) on a sequence of uniform quadrilateral and triangular grids for two different values of  $\varepsilon$ . Triangular grids are obtained by splitting each element in the quadrilateral grids into two triangles. Tables 4.1–4.2 present solution errors in various norms and the corresponding convergence rate estimates.

TABLE 4.1  
Errors and convergence rates for (2.14) on uniform quads and triangles with  $\varepsilon = 0.001$ .

Grid	Quadrilaterals			Triangles		
	$\ \phi - \phi_h\ $	$\ \nabla\phi - \nabla\phi_h\ _0$	$\ \phi_{\mathcal{T}} - \phi_h\ $	$\ \phi - \phi_h\ $	$\ \nabla\phi - \nabla\phi_h\ _0$	$\ \phi_{\mathcal{T}} - \phi_h\ $
32	0.4260E-02	0.7533E-01	0.6333E-02	0.7707E-02	0.8804E-01	0.9588E-02
64	0.2073E-02	0.4794E-01	0.3052E-02	0.3854E-02	0.5712E-01	0.4633E-02
128	0.1061E-02	0.2764E-01	0.1425E-02	0.2029E-02	0.3804E-01	0.2338E-02
Rate	1.002	0.723	1.076	0.963	0.605	1.018

TABLE 4.2  
Errors and convergence rates for (2.14) on uniform quads and triangles with  $\varepsilon = 0.00001$ .

Grid	Quadrilaterals			Triangles		
	$\ \phi - \phi_h\ $	$\ \nabla\phi - \nabla\phi_h\ _0$	$\ \phi_{\mathcal{T}} - \phi_h\ $	$\ \phi - \phi_h\ $	$\ \nabla\phi - \nabla\phi_h\ _0$	$\ \phi_{\mathcal{T}} - \phi_h\ $
32	0.4739E-02	0.7949E-01	0.6835E-02	0.8594E-02	0.9514E-01	0.1054E-01
64	0.2518E-02	0.5497E-01	0.3552E-02	0.4616E-02	0.6664E-01	0.5485E-02
128	0.1299E-02	0.3842E-01	0.1809E-02	0.2402E-02	0.4684E-01	0.2796E-02
Rate	0.934	0.524	0.959	0.916	0.511	0.957

The results in the tables confirm that the solution of (2.14) is first-order accurate in  $L^2$ . These results also suggest that the dependence on  $\varepsilon$  in the theoretical bound (3.25) appears to be fairly benign. Indeed, decreasing the diffusion coefficient by a factor of 100 does not reduce significantly the  $L^2$  convergence rates, which drop only by 7% and 5%, on triangles and quads, respectively. The  $H^1$ -seminorm errors are also in line with the theoretical estimate (3.22). In particular, for  $\varepsilon = 0.00001$  the rate of convergence is in an almost perfect agreement with the theoretical prediction.

We note that numerically, the discrete error outperforms the theoretical rate of convergence in (3.20). This may indicate a possible superconvergent behavior of the discrete error on uniform grids that cannot be captured by the proof of Theorem 3.5, which addresses general unstructured grids. Since this superconvergent behavior does not seem to extend to the  $H^1$ -seminorm error, its further investigation is beyond the scope of this paper.

Finally, it is worth pointing out that convergence rates on quadrilaterals are slightly, but consistently, higher than on triangles. Theorem 3.3 provides a possible explanation of this observation. Specifically, stability bounds (3.3) and (3.4) suggest that the stabilized method (2.14) is *less diffusive* on quadrilateral grids. We will further examine this conjecture in Section 4.2.

**4.2. Qualitative properties.** Although asymptotic convergence rates are an important quantitative measure of the accuracy of numerical methods they don't always tell the whole story about their computational properties. In this section we use the challenging ‘‘Double Glazing’’ advection

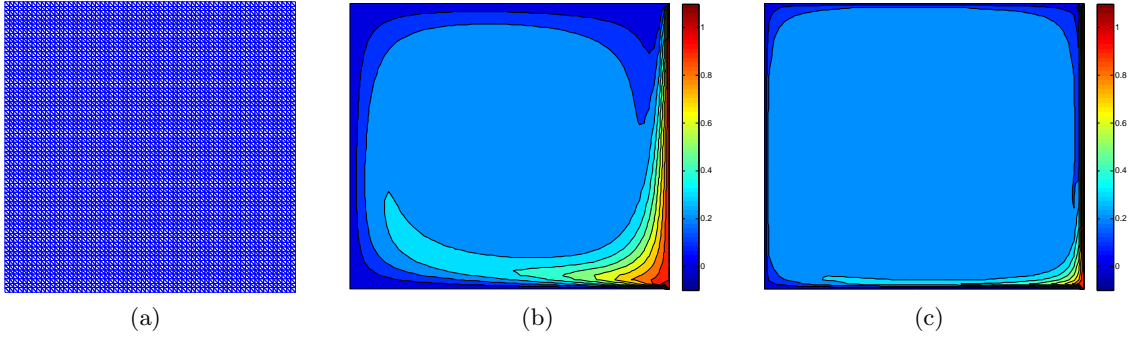


FIG. 4.1. Solution (b) of the Double Glazing test problem (4.1) with  $\varepsilon = 0.00001$  on a  $64 \times 64$  uniform triangular grid (a) and a  $64 \times 64$  uniform quadrilateral grid (c) by (2.14).

test [8, Example 3.1.4, p.119] to provide a complementary, qualitative perspective on (2.14). On the unit square the Double Glazing advection test is specified by

$$(4.1) \quad \mathbf{u} = \begin{pmatrix} 2(2y-1)(1-(2x-1)^2) \\ -2(2x-1)(1-(2y-1)^2) \end{pmatrix}; \quad f = 0; \quad \text{and} \quad g = \begin{cases} 1 & \text{on } \Gamma_R \\ 0 & \text{on } \Gamma_B \cup \Gamma_T \cup \Gamma_L \end{cases}.$$

Problem (4.1) models temperature distribution in a cavity with a “hot” external wall ( $\Gamma_R$ ). The discontinuities at the two corners of the hot wall create boundary layers near its corners.

*Triangles vs. quadrilaterals.* According to Theorem 3.3 on tensor product elements the stabilized form  $a^h(\cdot, \cdot)$  is bounded from below by the streamline derivative of the finite element solution along the effective elemental advective velocity  $\mathbf{u}_k$ . On the other hand, for simplicial elements this form is bounded from below by the gradient of the finite element solution. The difference in the stability bounds (3.3) and (3.4) prompts a conjecture that (2.14) will be less diffusive on tensor product elements. To test this conjecture we solve (4.1) with  $\varepsilon = 0.00001$  on uniform  $64 \times 64$  triangular and quadrilateral grids. Figure 4.1 presents the solution plots, which clearly show the more diffusive behavior of the finite element solution on the triangular grid.

*Edge flux stabilized vs. “symmetrized” stabilized formulations.* Direct stabilization by the edge element flux  $F_E(\phi_h)$  in (2.11) introduces the non-symmetric artificial diffusion kernel exposed in (2.12). The lack of symmetry in this kernel may be detrimental for the accuracy on unstructured grids. Figure 4.2 compares solution of the Double Glazing test problem with  $\varepsilon = 0.0001$  by (2.11) and the new “symmetrized” formulation (2.14) on a  $64 \times 64$  randomly perturbed quadrilateral grid. From the plots it is evident that the solution of (2.11) suffers from numerical noise, whereas the “symmetrized” formulation is not affected by the mesh structure. These observations can be further quantified by measuring the violation of global bounds by the two solutions. For the test problem  $\phi^{\min} = 0$  and  $\phi^{\max} = 1$  and so, the quantity  $\Delta = (\phi_h^{\max} - \phi_h^{\min}) - 1$  provides a measure of the total violation of global solution bounds by the finite element approximations. For the solution of (2.11)

$$\phi_h^{\max} = 1.0936, \quad \phi_h^{\min} = -0.0986 \quad \text{and} \quad \Delta = 0.19,$$

whereas for the solution of (2.14)

$$\phi_h^{\max} = 1.0792, \quad \phi_h^{\min} = -0.0052 \quad \text{and} \quad \Delta = 0.08$$

In other words, violation of global bounds in (2.11) is twice the violation in (2.14).



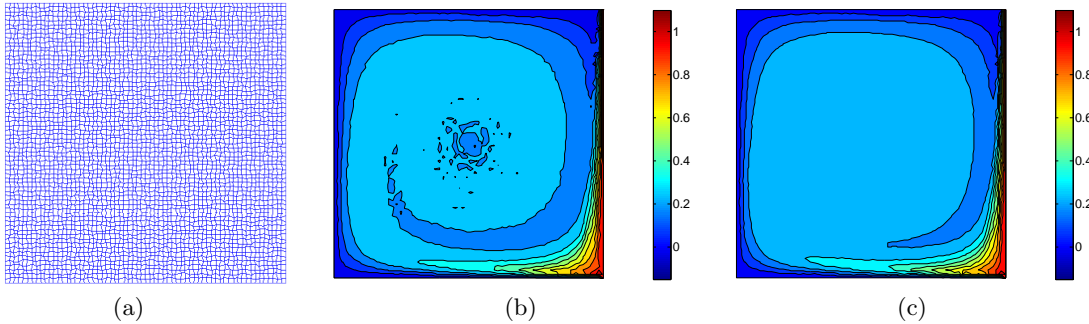


FIG. 4.2. Solution of the Double Glazing test problem (4.1) with  $\varepsilon = 0.0001$  on a  $64 \times 64$  randomly perturbed quadrilateral mesh (a) by the stabilized formulation (2.11) (b) and the “symmetrized” stabilized method (2.14) (c).

“Symmetrized” stabilized formulation vs. Artificial Diffusion and SUPG. The final qualitative study compares and contrasts (2.14) with the classical artificial diffusion method [14, p.181] and the Streamline Upwind Petrov-Galerkin (SUPG) method [12]. For this test we solve the Double Glazing problem with  $\varepsilon = 0.00001$  on a  $64 \times 64$  uniform quadrilateral grid. Figure 4.3 presents the results. From the solution plots in this figure it is immediately obvious that our formulation is significantly less diffusive than the artificial diffusion method. It also appears to be more robust than the SUPG solution, which develops spurious oscillations near the boundary layers, leading to significant violations of the global solution bounds. In particular, for the SUPG solution

$$\phi_h^{\max} = 1.19, \quad \phi_h^{\min} = -0.28 \quad \text{and} \quad \Delta = 0.47,$$

resulting in a 47% violation of the global solution bounds. In contrast, for the solution of (2.14)

$$\phi_h^{\max} = 1.00, \quad \phi_h^{\min} = -0.0097 \quad \text{and} \quad \Delta = 0.0097,$$

i.e., the total violation of the solution bounds is less than 1%.

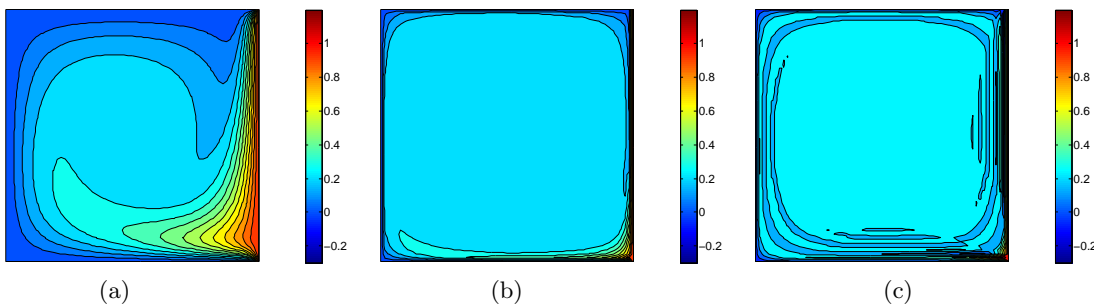


FIG. 4.3. Solution of the Double Glazing test problem (4.1) with  $\varepsilon = 0.00001$  on a  $64 \times 64$  uniform quadrilateral mesh by the Artificial Diffusion method (a) the “symmetrized” stabilized method (2.14) (b), and SUPG (c).

**Acknowledgements.** The authors acknowledge the support of the Advanced Scientific Computing Research program of the DoE Office of Science and the ASC program of the NNSA in carrying out this research. Discussions with our colleagues from the Charon team X. Gao, G. Hennigan, L. Musso, and T. Smith motivated many aspects of this work.

**Appendix A. Extensions to non-constant diffusion and velocity.** Throughout this section we assume that  $\varepsilon$  and  $\mathbf{u}$  are of class  $C^1(\Omega)$  and  $\nabla \cdot \mathbf{u} = 0$ . Whenever necessary, dependence on diffusivity and velocity will be indicated explicitly by including these fields in the list of arguments, e.g.,  $a^h(\phi_h, \varphi_h; \varepsilon, \mathbf{u})$ ,  $\theta_\alpha(\varepsilon, \bar{u}_\alpha)$ ,  $p_\alpha(\varepsilon, \bar{u}_\alpha)$ , and so on. To extend the results of Theorem 3.3 to non-constant  $\varepsilon$  and  $\mathbf{u}$  it suffices to show that on every  $\mathbf{k} \in K(\Omega)$  the element form  $a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u})$  can be bounded from below by a constant times  $a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}})$ , where  $\varepsilon_{\mathbf{k}}$  and  $\mathbf{u}_{\mathbf{k}}$  are constant approximations of  $\varepsilon$  and  $\mathbf{u}$  on element  $\mathbf{k}$ . Here, we set these approximations to

$$\varepsilon_{\mathbf{k}} := \max_{\mathbf{x} \in \mathbf{k}} \varepsilon(\mathbf{x}) \quad \text{and} \quad \mathbf{u}_{\mathbf{k}} = \frac{1}{|\mathbf{k}|} \int_{\mathbf{k}} \mathbf{u} dV,$$

respectively, i.e., the largest value of  $\varepsilon$  on  $\mathbf{k}$  and the average element velocity, respectively.

LEMMA A.1. *For sufficiently small mesh size  $h$  there exists a positive constant  $C$  such that*

$$(A.1) \quad a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u}) \geq C a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) \quad \forall \mathbf{k} \in K(\Omega).$$

*Proof.* For variable  $\varepsilon$  and  $\mathbf{u}$

$$a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u}) = \|\sqrt{\varepsilon} \nabla \phi_h\|_{0, \mathbf{k}}^2 + Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon, \mathbf{u}).$$

Regularity assumptions imply the existence of a positive constant  $\gamma$  such that  $\varepsilon_{\mathbf{k}} \leq (1 + \gamma) \varepsilon$  on every element  $\mathbf{k} \in K(\Omega)$ . This yields the following bound for the first term in  $a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u})$ :

$$(A.2) \quad \|\sqrt{\varepsilon} \nabla \phi_h\|_{0, \mathbf{k}}^2 \geq \frac{\varepsilon_{\mathbf{k}}}{1 + \gamma} \|\nabla \phi_h\|_{0, \mathbf{k}}^2.$$

We bound the second term from below in two steps, starting with an estimate of  $Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon, \mathbf{u})$  by  $Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}})$ . Proceeding as in the proof of Theorem 3.3, and using that  $\theta_\alpha(\varepsilon, \bar{u}_\alpha)$  is monotone decreasing with respect to its first argument, we find that

$$Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon, \mathbf{u}) \geq \underline{\lambda}(M_{\mathbf{k}}) \sum_{e_\xi \in E(\mathbf{k})} \theta_\xi(\varepsilon, \bar{u}_\xi) z_\xi^2 \geq \underline{\lambda}(M_{\mathbf{k}}) \sum_{e_\xi \in E(\mathbf{k})} \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) z_\xi^2,$$

where, as before,  $z_\xi$  is the element of  $\mathbf{z}_{\mathbf{k}} = \phi_{\mathbf{k}} G_{\mathbf{k}}$ . On the other hand,

$$Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) \leq \bar{\lambda}(M_{\mathbf{k}}) \sum_{e_\xi \in E(\mathbf{k})} \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) z_\xi^2.$$

Combining these two results yields the intermediate estimate

$$(A.3) \quad Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon, \mathbf{u}) \geq \kappa^{-1}(M_{\mathbf{k}}) Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}).$$

We proceed to bound  $Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}})$  from below in terms of  $Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}, \xi})$ . Let

$$D_{\mathbf{k}}(\mathbf{u}) := D_{\mathbf{k}}(\varepsilon_{\mathbf{k}}, \mathbf{u}) = \text{diag} \left( \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) \right); \quad D_{\mathbf{k}} := D_{\mathbf{k}}(\varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) = \text{diag} \left( \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, u_{\mathbf{k}, \xi}) \right)$$

and  $D_{\mathbf{k}}^\Delta = D_{\mathbf{k}}(\mathbf{u}) - D_{\mathbf{k}} = \text{diag} \left( \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) - \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, u_{\mathbf{k}, \xi}) \right),$

where  $\tilde{\theta}_\xi = \sqrt{\bar{\theta}_\xi}$  and  $\bar{u}_\xi$  and  $\bar{u}_{\mathbf{k},\xi}$  are the average tangential components of  $\mathbf{u}$  and  $\mathbf{u}_{\mathbf{k}}$ , respectively, along edge  $e_\xi$ . In terms of these matrices

$$\begin{aligned} Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}) &= \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} + D_{\mathbf{k}}^\Delta) M_{\mathbf{k}} (D_{\mathbf{k}} + D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} \\ &= \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}}) \mathbf{z}_{\mathbf{k}} + \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} + 2\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}}. \end{aligned}$$

Applying the Schwartz's and Young's inequalities to the last, mixed term gives

$$\begin{aligned} 2\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} &\leq 2(\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}})^{1/2} (\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}}) \mathbf{z}_{\mathbf{k}})^{1/2} \\ &\leq \delta \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} + \frac{1}{\delta} \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}} M_{\mathbf{k}} D_{\mathbf{k}}) \mathbf{z}_{\mathbf{k}} \end{aligned}$$

with a constant  $\delta > 1$  that will be determined later. Therefore,

$$(A.4) \quad Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}) \geq \left(1 - \frac{1}{\delta}\right) Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) - (\delta - 1) \mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}}.$$

We estimate the last term in (A.4) as follows. Note that

$$\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} \leq \bar{\lambda}(M_{\mathbf{k}}) \sum_{e_\xi \in E(\mathbf{k})} \left(\tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) - \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_{\mathbf{k},\xi})\right)^2 z_\xi^2$$

Recalling that  $\tilde{\theta}_\xi = \sqrt{\bar{\theta}_\xi}$  allows us to conclude that

$$\left(\tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) - \tilde{\theta}_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_{\mathbf{k},\xi})\right)^2 \leq |\theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) - \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_{\mathbf{k},\xi})|.$$

Expanding  $\theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi)$  in a Taylor series about  $\bar{u}_{\mathbf{k},\xi}$  shows that

$$|\theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_\xi) - \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{u}_{\mathbf{k},\xi})| \leq \left| \frac{\partial \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{v}_\xi)}{\partial \bar{u}_\xi} \right| |\bar{u}_\xi - \bar{u}_{\mathbf{k},\xi}|,$$

where  $\bar{v}_\xi = \lambda \bar{u}_\xi + (1 - \lambda) \bar{u}_{\mathbf{k},\xi}$  for some  $0 < \lambda < 1$ . Direct calculation reveals that

$$\left| \frac{\partial \theta_\xi(\varepsilon_{\mathbf{k}}, \bar{v}_\xi)}{\partial \bar{u}_\xi} \right| = \frac{h_\xi}{2} \left[ \coth p_\xi(\varepsilon_{\mathbf{k}}, \bar{v}_\xi) + p_\xi(\varepsilon_{\mathbf{k}}, \bar{v}_\xi) (1 - \coth^2 p_\xi(\varepsilon_{\mathbf{k}}, \bar{v}_\xi)) \right] \leq \frac{h_\xi}{2}$$

with the last inequality following from the fact that  $0 \leq \coth x + x(1 - \coth^2 x) \leq 1$ . Finally, consider the Taylor expansion  $\mathbf{u}(\mathbf{x}) = \mathbf{u}_{\mathbf{k}} + \nabla \mathbf{u}(\boldsymbol{\omega}) \cdot (\mathbf{x} - \bar{\mathbf{x}})$ , where  $\mathbf{u}(\bar{\mathbf{x}}) = \mathbf{u}_{\mathbf{k}}$  and  $\boldsymbol{\omega}$  is some point in  $\mathbf{k}$ . Using this expansion we find that

$$|\bar{u}_\xi - \bar{u}_{\mathbf{k},\xi}| = \frac{1}{h_\xi} \left| \int_{e_\xi} \mathbf{t}_\xi \cdot \nabla \mathbf{u}(\boldsymbol{\omega}) \cdot (\mathbf{x} - \bar{\mathbf{x}}) dl \right| \leq \frac{1}{h_\xi} \|\nabla \mathbf{u}\|_{0,\infty,\mathbf{k}} \int_{e_\xi} |\mathbf{x} - \bar{\mathbf{x}}| dl \leq h_\xi \|\nabla \mathbf{u}\|_{0,\infty,\mathbf{k}}.$$

After combining all of the above bounds we find that

$$\mathbf{z}_{\mathbf{k}}^T (D_{\mathbf{k}}^\Delta M_{\mathbf{k}} D_{\mathbf{k}}^\Delta) \mathbf{z}_{\mathbf{k}} \leq \bar{\lambda}(M_{\mathbf{k}}) \frac{h_{\mathbf{k}}^2}{2} \|\nabla \mathbf{u}\|_{0,\infty,\mathbf{k}} |\mathbf{z}_{\mathbf{k}}|^2.$$

Since  $\|\nabla\phi_h\|_{0,\mathbf{k}}^2 \geq \underline{\lambda}(M_{\mathbf{k}})|\mathbf{z}_{\mathbf{k}}|^2$  it follows that

$$\mathbf{z}_{\mathbf{k}}^T (D_{\hat{\mathbf{k}}}^\Delta M_{\mathbf{k}} D_{\hat{\mathbf{k}}}^\Delta) \mathbf{z}_{\mathbf{k}} \leq \kappa(M_{\mathbf{k}}) \frac{h_{\mathbf{k}}^2}{2} \|\nabla \mathbf{u}\|_{0,\infty,\mathbf{k}} \|\nabla \phi_h\|_{0,\mathbf{k}}^2.$$

Owing to the regularity assumptions, for a sufficiently small mesh size

$$\kappa(M_{\mathbf{k}}) \frac{h_{\mathbf{k}}^2}{2} \|\nabla \mathbf{u}\|_{0,\infty,\mathbf{k}} \leq \gamma \varepsilon_{\mathbf{k}}.$$

Using this inequality together with (A.4) allows us to conclude that

$$Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}) \geq \left(1 - \frac{1}{\delta}\right) Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) - (\delta - 1)\gamma \varepsilon_{\mathbf{k}} \|\nabla \phi_h\|_{0,\mathbf{k}}^2.$$

To complete the proof we combine the above inequality with (A.2) and (A.3):

$$\begin{aligned} a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u}) &\geq \frac{\varepsilon_{\mathbf{k}}}{1 + \gamma} \|\nabla \phi_h\|_{0,\mathbf{k}}^2 + \kappa^{-1}(M_{\mathbf{k}}) \left[ \left(1 - \frac{1}{\delta}\right) Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}) - (\delta - 1)\gamma \varepsilon_{\mathbf{k}} \|\nabla \phi_h\|_{0,\mathbf{k}}^2 \right] \\ &= \left( \frac{1}{1 + \gamma} - \kappa^{-1}(M_{\mathbf{k}})(\delta - 1)\gamma \right) \varepsilon_{\mathbf{k}} \|\nabla \phi_h\|_{0,\mathbf{k}}^2 + \kappa^{-1}(M_{\mathbf{k}}) \left(1 - \frac{1}{\delta}\right) Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}). \end{aligned}$$

Setting  $\delta = 1 + \kappa(M_{\mathbf{k}})/(\gamma + 1)^2$  gives

$$a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon, \mathbf{u}) \geq \frac{\varepsilon_{\mathbf{k}} \|\nabla \phi_h\|_{0,\mathbf{k}}^2}{(1 + \gamma)^2} + \frac{Q_{\mathbf{k}}(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}})}{\kappa(M_{\mathbf{k}}) + (1 + \gamma)^2} \geq C(\gamma) a_{\mathbf{k}}^h(\phi_h, \phi_h; \varepsilon_{\mathbf{k}}, \mathbf{u}_{\mathbf{k}}),$$

which proves the theorem with  $C(\gamma) = 1/(\kappa(M_{\mathbf{k}}) + (1 + \gamma)^2)$ .  $\square$

**Appendix B. Tensor product finite element spaces.** We label the edges of quadrilateral and hexahedral elements according to the versors of the reference coordinate system.

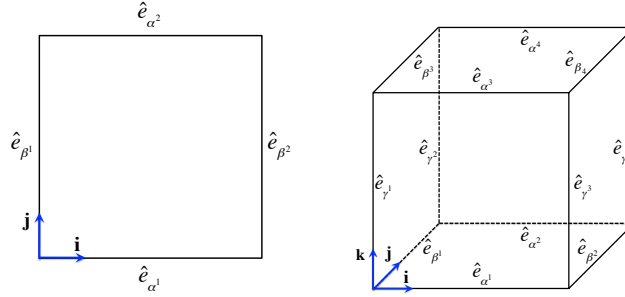


FIG. B.1. Edge numbering for reference quadrilateral and hexahedral elements.

A reference quadrilateral  $\hat{\mathbf{k}}$  has two pairs of edges  $\{\hat{\mathbf{e}}_{\alpha^1}, \hat{\mathbf{e}}_{\alpha^2}\}$  and  $\{\hat{\mathbf{e}}_{\beta^1}, \hat{\mathbf{e}}_{\beta^2}\}$ , parallel to  $\mathbf{i}$  and  $\mathbf{j}$ , respectively; see Fig. B.1. Let  $F_{\mathbf{k}}$  be the bilinear map between  $\hat{\mathbf{k}}$  and a quadrilateral  $\mathbf{k} \in K(\Omega)$ . The edges of  $\mathbf{k}$  are images of the reference edges under  $F_{\mathbf{k}}$  and form the pairs  $\{\mathbf{e}_{\alpha^1}, \mathbf{e}_{\alpha^2}\}$  and  $\{\mathbf{e}_{\beta^1}, \mathbf{e}_{\beta^2}\}$ , respectively. The unit tangents to the edges of  $\mathbf{k}$  form another two pairs labeled by  $\{\mathbf{t}_{\alpha^1}, \mathbf{t}_{\alpha^2}\}$  and  $\{\mathbf{t}_{\beta^1}, \mathbf{t}_{\beta^2}\}$ , respectively. Let  $J_{\mathbf{k}}(\mathbf{x})$  be the Jacobian of  $F_{\mathbf{k}}$  and define the vector fields

$$\mathbf{t}_{\mathbf{i}}(\mathbf{x}) = \frac{\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{i}}{|\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{i}|} \quad \text{and} \quad \mathbf{t}_{\mathbf{j}}(\mathbf{x}) = \frac{\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{j}}{|\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{j}|}.$$

It is easy to see that  $\{\mathbf{t}_{\alpha^1}, \mathbf{t}_{\alpha^2}\}$  are the restrictions of  $\mathbf{t}_i$  to  $\{\mathbf{e}_{\alpha^1}, \mathbf{e}_{\alpha^2}\}$ :  $\mathbf{t}_{\alpha^1} = \mathbf{t}_i(\mathbf{x})|_{\hat{\mathbf{e}}_{\alpha^1}}$  and  $\mathbf{t}_{\alpha^2} = \mathbf{t}_i(\mathbf{x})|_{\hat{\mathbf{e}}_{\alpha^2}}$ . Likewise,  $\{\mathbf{t}_{\beta^1}, \mathbf{t}_{\beta^2}\}$  are the restrictions of  $\mathbf{t}_j(\mathbf{x})$  to  $\{\mathbf{e}_{\beta^1}, \mathbf{e}_{\beta^2}\}$ .

A reference hexahedral  $\hat{\mathbf{k}}$  has three four-tuples of edges parallel to  $\mathbf{i}$ ,  $\mathbf{j}$  and  $\mathbf{k}$ , respectively, labeled as shown in Figure B.1. The edges of an element  $\mathbf{k} \in K(\Omega)$  are images of the reference edges under a trilinear map  $F_{\mathbf{k}}$ . They form three four-tuples, which we denote by  $\{\mathbf{e}_{\alpha^1}, \mathbf{e}_{\alpha^2}, \mathbf{e}_{\alpha^3}, \mathbf{e}_{\alpha^4}\}$ ,  $\{\mathbf{e}_{\beta^1}, \mathbf{e}_{\beta^2}, \mathbf{e}_{\beta^3}, \mathbf{e}_{\beta^4}\}$ , and  $\{\mathbf{e}_{\gamma^1}, \mathbf{e}_{\gamma^2}, \mathbf{e}_{\gamma^3}, \mathbf{e}_{\gamma^4}\}$ , respectively. Likewise, restrictions of the vector fields

$$(B.1) \quad \mathbf{t}_i(\mathbf{x}) = \frac{\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{i}}{|\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{i}|}; \quad \mathbf{t}_j(\mathbf{x}) = \frac{\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{j}}{|\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{j}|} \quad \text{and} \quad \mathbf{t}_k(\mathbf{x}) = \frac{\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{k}}{|\mathbf{J}_{\mathbf{k}}(\mathbf{x}) \mathbf{k}|}$$

to four-tuples of reference edges give the unit tangents of the associated four-tuples of edges on  $\mathbf{k}$ .

Stability analysis on tensor product elements requires a notion of an *effective advective velocity* field  $\mathbf{u}_{\mathbf{k}}$  for element  $\mathbf{k}$ . To define this field consider the normalized projections of the given velocity field  $\mathbf{u}$  onto the four-tuples of element edges associated with the versor directions

$$(B.2) \quad v_{\alpha^i} = \frac{\mathbf{u} \cdot \mathbf{t}_{\alpha^i}}{|\mathbf{u}|}; \quad v_{\beta^i} = \frac{\mathbf{u} \cdot \mathbf{t}_{\beta^i}}{|\mathbf{u}|}; \quad v_{\gamma^i} = \frac{\mathbf{u} \cdot \mathbf{t}_{\gamma^i}}{|\mathbf{u}|},$$

and let

$$(B.3) \quad v_i = \arg \min_i |v_{\alpha^i}|; \quad v_j = \arg \min_i |v_{\beta^i}|; \quad v_k = \arg \min_i |v_{\gamma^i}|.$$

The effective advective velocity on a tensor product element  $\mathbf{k}$  is then given by

$$(B.4) \quad \mathbf{u}_{\mathbf{k}} = \begin{cases} |\mathbf{u}|(v_i \mathbf{t}_i + v_j \mathbf{t}_j) & \text{for quadrilaterals} \\ |\mathbf{u}|(v_i \mathbf{t}_i + v_j \mathbf{t}_j + v_k \mathbf{t}_k) & \text{for hexahedrons} \end{cases}.$$

Finally, note that the edge element basis functions on a quadrilateral element  $\mathbf{k}$  can be written as

$$(B.5) \quad \vec{W}_{\alpha^i} = A_i(\mathbf{x}) \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{i} \quad \text{and} \quad \vec{W}_{\beta^i} = B_i(\mathbf{x}) \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j}; \quad i = 1, 2;$$

where  $A_i(\mathbf{x})$  and  $B_i(\mathbf{x})$  are scalar functions, whereas for hexahedral elements

$$(B.6) \quad \vec{W}_{\alpha^i} = A_i(\mathbf{x}) \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{i}; \quad \vec{W}_{\beta^i} = B_i(\mathbf{x}) \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{j} \quad \text{and} \quad \vec{W}_{\gamma^i} = C_i(\mathbf{x}) \mathbf{J}_{\mathbf{k}}^{-T}(\mathbf{x}) \mathbf{k} \quad i = 1, \dots, 4$$

for some other scalar functions  $A_i(\mathbf{x})$ ,  $B_i(\mathbf{x})$  and  $C_i(\mathbf{x})$ .

#### REFERENCES

- [1] D. N. Arnold, R. S. Falk, and R. Winther. Finite element exterior calculus, homological techniques, and applications. *Acta Numerica*, 15:1–155, 2006.
- [2] P. Bochev and M. Hyman. Principles of mimetic discretizations. In D. N. Arnold, P. Bochev, R. Lehoucq, R. Nicolaides, and M. Shashkov, editors, *Compatible Discretizations, Proceedings of IMA Hot Topics Workshop on Compatible Discretizations*, volume IMA 142, pages 89–120. Springer Verlag, 2006.
- [3] Pavel Bochev and Kara Peterson. A parameter-free stabilized finite element method for scalar advection-diffusion problems. *Central European Journal of Mathematics*, 11(8):1458–1477, 2013.
- [4] Pavel Bochev, Kara Peterson, and Xujiao Gao. A new control volume finite element method for the stable and accurate solution of the drift–diffusion equations on general unstructured grids. *Computer Methods in Applied Mechanics and Engineering*, 254(0):126 – 145, 2013.

- [5] Alexander N. Brooks and Thomas J.R. Hughes. Streamline upwind/Petrov-Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier-Stokes equations. *Computer Methods in Applied Mechanics and Engineering*, 32(1–3):199 – 259, 1982.
- [6] P. Ciarlet. *The Finite Element Method for Elliptic Problems*. SIAM Classics in Applied Mathematics. SIAM, Philadelphia, 2002.
- [7] Ramon Codina. Comparison of some finite element methods for solving the diffusion-convection-reaction equation. *Computer Methods in Applied Mechanics and Engineering*, 156(1-4):185 – 210, 1998.
- [8] H. C. Elman, D. J. Silvester, and A. J. Wathen. *Finite Elements and Fast Iterative Solvers with Applications in Incompressible Fluid Dynamics*. Numerical Mathematics and Scientific Computation. Oxford University Press, 2005.
- [9] A. Ern and J.-L. Guermond. *Theory and Practice of Finite Elements*. Number 159 in Applied Mathematical Sciences. Springer Verlag, New York, 2004.
- [10] I. Harari and T. J. R. Hughes. What are  $C$  and  $h$ ?: Inequalities for the analysis and design of finite element methods. *Comput. Meth. Appl. Mech. Eng.*, 97:157–192, 1992.
- [11] T. J. R. Hughes. Multiscale phenomena: Green’s function, the Dirichlet-to-Neumann map, subgrid scale models, bubbles and the origins of stabilized methods. *Comput. Meth. Appl. Mech. Eng.*, 127:387–401, 1995.
- [12] T. J. R. Hughes and A. Brooks. A theoretical framework for Petrov-Galerkin methods with discontinuous weighting functions: Application to the streamline-upwind procedure. In R. H. Gallagher et al, editor, *Finite Elements in Fluids*, volume 4, pages 47–65, New York, 1982. J. Wiley & Sons.
- [13] Thomas J.R. Hughes, Michel Mallet, and Mizukami Akira. A new finite element formulation for computational fluid dynamics: Ii. beyond supg. *Computer Methods in Applied Mechanics and Engineering*, 54(3):341 – 355, 1986.
- [14] C. Johnson. *Numerical Solution of Partial Differential Equations by the Finite Element Method*. Cambridge University Press, 1992.
- [15] P. Knobloch. On the definition of the SUPG parameter. *ETNA*, 32:76–89, 2008.
- [16] J. C. Nédélec. Mixed finite elements in  $R^3$ . *Numerische Mathematik*, 35:315–341, 1980. 10.1007/BF01396415.
- [17] D.L. Scharfetter and H.K. Gummel. Large-signal analysis of a silicon read diode oscillator. *IEEE Transactions on*, 16(1):64 – 77, jan 1969.
- [18] D. Z. Turner, K. B. Nakshatrala, and K. D. Hjelmstad. A stabilized formulation for the advection-diffusion equation using the generalized finite element method. *International Journal for Numerical Methods in Fluids*, 66(1):64–81, 2011.