

GUARANTEED AND ROBUST L^2 -NORM *A POSTERIORI* ERROR ESTIMATES FOR 1D LINEAR ADVECTION PROBLEMS

ALEXNDRE ERN^{1,2}, MARTIN VOHRALÍK^{2,1} AND MOHAMMAD ZAKERZADEH^{2,1,*}

Abstract. We propose a reconstruction-based *a posteriori* error estimate for linear advection problems in one space dimension. In our framework, a stable variational ultra-weak formulation is adopted, and the equivalence of the L^2 -norm of the error with the dual graph norm of the residual is established. This dual norm is showed to be localizable over vertex-based patch subdomains of the computational domain under the condition of the orthogonality of the residual to the piecewise affine hat functions. We show that this condition is valid for some well-known numerical methods including continuous/discontinuous Petrov–Galerkin and discontinuous Galerkin methods. Consequently, a well-posed local problem on each patch is identified, which leads to a global conforming reconstruction of the discrete solution. We prove that this reconstruction provides a guaranteed upper bound on the L^2 error. Moreover, up to a generic constant, it also gives local lower bounds on the L^2 error, where the constant only depends on the mesh shape-regularity. This, in particular, leads to robustness of our estimates with respect to the polynomial degree. All the above properties are verified in a series of numerical experiments, additionally leading to asymptotic exactness. Motivated by these results, we finally propose a heuristic extension of our methodology to any space dimension, achieved by solving local least-squares problems on vertex-based patches. Though not anymore guaranteed, the resulting error indicator is still numerically robust with respect to both advection velocity and polynomial degree in our collection of two-dimensional test cases including discontinuous solutions aligned and not aligned with the computational mesh.

Mathematics Subject Classification. 65N15, 65N30, 35F05.

Received April 21, 2019. Accepted June 10, 2020.

1. INTRODUCTION

This work deals with a linear advection equation of the form: find $u : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}$ such that

$$\mathbf{b} \cdot \nabla u = f, \quad \text{in } \Omega, \quad (1.1a)$$

$$u = 0, \quad \text{on } \partial_- \Omega. \quad (1.1b)$$

The velocity field $\mathbf{b} \in C^1(\bar{\Omega}; \mathbb{R}^d)$, $\mathbf{b} \neq 0$, is considered to be divergence-free and we take into account a general source term $f \in L^2(\Omega)$. The inflow, outflow, and characteristic parts of the boundary are denoted by $\partial_- \Omega$, $\partial_+ \Omega$,

Keywords and phrases. linear advection problem, discontinuous Galerkin method, Petrov–Galerkin method, *a posteriori* error estimate, local efficiency, advection robustness, polynomial-degree robustness.

¹ Université Paris-Est, CERMICS (ENPC), 77455 Marne-la-Vallée 2, France.

² Inria, 2 rue Simone Iff, 75589 Paris, France.

*Corresponding author: seyed-mohammad.zakerzadeh@inria.fr

and $\partial_0\Omega$, respectively, with the definitions

$$\partial_{\pm}\Omega := \{x \in \partial\Omega : \pm\mathbf{b}(x) \cdot \mathbf{n}(x) > 0\}, \quad \partial_0\Omega := \{x \in \partial\Omega : \mathbf{b}(x) \cdot \mathbf{n}(x) = 0\}.$$

In the main body of the paper, we focus on the one-dimensional case $d = 1$, where $\Omega \subset \mathbb{R}$ is a bounded interval; then \mathbf{b} is a constant scalar. We keep the notation in multi-dimensional form in order to be applicable when we discuss extensions of our results to the multi-dimensional case. For simplicity, we only consider a homogeneous boundary condition, but all the results can be extended to the non-homogeneous case, see Remarks 4.9 and 10.9 below.

The *a posteriori* error analysis for problem (1.1) admits a range of functional frameworks and consequently different norms in which the error can be measured. Our goal is to derive an L^2 -norm error estimate of the form

$$\|u - u_h\|_{L^2(\Omega)} \leq \eta, \tag{1.2}$$

where u is the weak solution of (1.1) in $L^2(\Omega)$, u_h is its numerical approximation, and η is an *a posteriori* error estimator that is *fully computable* from u_h by some *local procedure*. We seek to have a bound that is *guaranteed*, *i.e.*, featuring no unknown constant, in contrast to *reliability* where a bound up to a generic constant is sufficient. We develop a unified framework treating several classical numerical methods at once. Importantly, we also prove a converse estimate to (1.2) in the form

$$\eta \leq C\|u - u_h\|_{L^2(\Omega)} + \text{data oscillation}. \tag{1.3}$$

This is called *global efficiency* and yields equivalence between the incomputable error $\|u - u_h\|_{L^2(\Omega)}$ and the computable estimator η , up to the data oscillation term that vanishes for piecewise polynomial datum f and that is of higher order than the error for piecewise smooth datum f . Crucially, in our developments, the generic constant C in (1.3) only depends on the mesh shape regularity, requesting for $d = 1$ that any two neighboring elements be of comparable size. In particular, C is independent of the problem parameters \mathbf{b} and f as well as of the polynomial degree of the approximation k , yielding both data- and polynomial-degree-robustness. We actually also show local efficiency, *i.e.*, a localized version of (1.3), which is highly desirable on the practical side in view of adaptive mesh refinement. We observe that in one space dimension with constant velocity field \mathbf{b} , data-robustness boils down to a linear behavior of the error indicator with respect to the magnitude of the velocity field. Data robustness is thus expected to be true for any reasonable result from literature for this particular case.

To achieve the above-mentioned goals, we start with the ultra-weak variational formulation at the infinite-dimensional level, where the solution lies in the $L^2(\Omega)$ trial space and the test space is formed by functions in the graph space of the formal adjoint operator taking zero value at the outflow boundary ($H^1(\Omega)$ with zero value at the outflow in one space dimension). In this setting, we prove the equality of the L^2 -norm of the error with the dual graph norm (relying on $\|\mathbf{b} \cdot \nabla(\cdot)\|_{L^2(\Omega)}$) of the residual. In the one-dimensional case, we are able to prove that the global dual norm can be localized over vertex-based patches of elements under an orthogonality condition against the hat basis functions. Consequently, suitable discrete local problems posed over these patches are identified which lead to local reconstructions s_h^a combined into a global reconstruction s_h such that $\|u_h - s_h\|_{L^2(\Omega)}$ forms the main ingredient of the estimator η satisfying (1.2) and (1.3).

Let us recall some important contributions to *a posteriori* error estimation for problem (1.1). Bey and Oden [5] proposed an *a posteriori* error estimate for a discontinuous Galerkin (dG) formulation of the multi-dimensional advection–reaction problem. In this framework, two infinite-dimensional problems have to be solved on each mesh element; one to obtain the lower bound on the error and one for the upper bound, in two different and inequivalent weighted energy norms. This gives estimates similar to (1.2) and (1.3), but for two different estimators and in two different norms of the error. Additionally, one cannot solve analytically the infinite-dimensional elementwise problems, and, in practice, one needs to approximate them by some higher-order finite element approximation. Hence, neither simultaneous reliability and efficiency, nor robustness, are granted.

Süli in [31] applied the H^1 -stability result of Tartakoff [33] to the adjoint problem of (1.1) (with the presence of the reaction term and in the multi-dimensional case), and obtained a global reliable upper bound on the H^{-1} -norm of the error in terms of the L^2 -norm of the residual for a weak formulation of (1.1) with distinct trial and test spaces. He further turned this bound into a reliable H^{-1} -norm *a posteriori* error indicator for the streamline-diffusion finite element and the cell-vertex finite volume methods. However, neither the efficiency nor the robustness of this error indicator are theoretically discussed. Furthermore, in [31] by Süli and in [23] by Houston *et al.*, an *a posteriori* error analysis of the multi-dimensional advection–reaction problem in the graph space equipped with the full norm $\|\cdot\|_{L^2(\Omega)} + \|\mathbf{b}\cdot\nabla(\cdot)\|_{L^2(\Omega)}$ is provided. This functional setting yields the equivalence between the L^2 -norm of the error and the dual graph norm of the residual, up to some generic constants. Propositions 2.1 and 10.4 below are actually closely related to these results, upon replacing the full graph norm by only $\|\mathbf{b}\cdot\nabla(\cdot)\|_{L^2(\Omega)}$, which leads to a constant-free error–residual *equivalence*. The rigorous reliability and efficiency results of [23, 31] in this functional setting for the L^2 -norm of the error are restricted to the part of the L^2 -error generated inside each mesh element, by neglecting the advected L^2 -error from the upwind. In numerical experiments, the estimates of [23] behave well in different flow fields and adaptive meshes.

Becker *et al.* [4] derived reconstruction-based error estimators for the advection problem (1.1) in two space dimensions. An $\mathbf{H}(\text{div}, \Omega)$ -conforming reconstruction is proposed for the flux vector $\mathbf{b}u$ (instead of u in the present work) which is designed to produce a guaranteed upper bound on the error measured in some dual norm of the advection operator. A unified framework is built, covering the dG, nonconforming, and conforming finite element methods with stabilization terms. This dual norm is hard to evaluate even for a known exact solution, and, in practice, the authors replace it by the L^2 -norm, so that the guaranteed upper bound property is eventually lost. Proofs of efficiency or robustness are not given, but optimal convergence orders of the estimator are observed in numerical experiments. It is worth mentioning that, restricted to one space dimension, the dual norm of [4] reduces to the weak graph norm we employ. Our contribution in this respect consists in the proofs of (1.2) and (1.3), not given in [4] (where, recall, two space dimensions are treated).

In a recent result by Georgoulis *et al.* [20], the authors used the reconstruction proposed by Makridakis and Nochetto [25] for a dG approximation and provided a reliable upper bound on the error in the energy norm for one-dimensional advection–diffusion–reaction problems, as well as a reliable L^2 -norm estimate for the problem (1.1) in one space dimension. Though a proof of (1.3) is not given, efficiency is numerically observed. One might also note the earlier work of these authors [19], dedicated to the two-dimensional advection–reaction problem with a similar reconstruction. In that work, a reliable bound on the energy norm of the error is presented, though again without a theoretical elaboration on the efficiency and robustness.

Furthermore we mention the recent result of Dahmen and Stevenson [11] where the authors provide *a posteriori* error estimates for the discontinuous Petrov–Galerkin method tailored to the transport equations in multiple space dimensions. The equivalence of the errors of the bulk and skeleton quantities with the dual norm of the residual is established. This dual norm is later approximated by some equivalent yet computable indicator. The absorbed constants translate into a constant C in (1.3) which depends on the advective field \mathbf{b} and the polynomial degree of approximation, but one might obtain stronger results by employing the approach of [11] in the simplified settings of this paper, *i.e.*, one-dimensional pure advection with constant velocity, for a specific discontinuous Petrov–Galerkin scheme. This is not in the scope of the present paper. Another important result of [11] is the adaptive mesh refinement strategy to guarantee a fixed error reduction in the one-dimensional case.

Finally, we also mention that in the case of advection–diffusion(–reaction) problems, other approaches were previously considered to obtain robustness with respect to the advective field. Among them, Verfürth [34] proposed to augment the energy norm by a dual norm coming from the skew-symmetric part of the differential operator, and Sangalli [27, 28] used interpolated spaces and a fractional-order norm for the advective term. Extensions of these approaches can be found in [17, 29, 30]. However, the above results are not applicable when the diffusion parameter vanishes, *i.e.*, as the advection–diffusion problem reduces to (1.1), because the diffusive part of the operator is needed to evaluate the dual norm. We extend our approach to the advection–reaction

case in [35], where a guaranteed *a posteriori* error estimate that is locally efficient and robust with respect to the interplay of the advection and reaction phenomena is derived.

We treat problem (1.1) in one space dimension in Sections 2–9. Section 2 deals with the functional setting, whereby adopting the ultra-weak variational formulation. We prove, in particular, the equality of the L^2 -norm of the error and the dual norm of the residual. Section 3 introduces some numerical schemes for approximating (1.1). Section 4 presents a local potential reconstruction on the patch level and collects the main results. Section 5 discusses the localization of the dual norm of the residual over vertex-based patches, and Section 6 shows that this is possible for the schemes discussed in Section 3. Sections 7 and 8 then present the proofs for the upper and lower bounds as well as robustness in the form of (1.2)–(1.3). Section 9 then contains results of several numerical experiments to illustrate the developed theory. Finally, in Section 10, we consider the advection problem (1.1) in multiple space dimensions and derive a heuristic extension of our methodology to this case. Although we cannot prove here the guaranteed upper bound, (local) efficiency, and robustness, numerical experiments indicate appreciable properties of the derived estimates also in this case, including discontinuous solutions aligned and not aligned with the computational mesh.

2. ABSTRACT FRAMEWORK

We start with the presentation of the abstract framework.

2.1. Spaces

In the one-dimensional case, the constraint of \mathbf{b} being a non-zero divergence-free field is translated to \mathbf{b} being a constant nonzero scalar. Consequently, we are lead to work with the spaces

$$H_-^1(\Omega) = \{w \in H^1(\Omega), w = 0, \text{ on } \partial_- \Omega\}, \quad H_+^1(\Omega) = \{w \in H^1(\Omega), w = 0, \text{ on } \partial_+ \Omega\}. \quad (2.1)$$

The trace operator in these spaces is well-defined and the following *integration-by-parts* formula holds:

$$(v, \mathbf{b} \cdot \nabla w)_\Omega + (\mathbf{b} \cdot \nabla v, w)_\Omega = (\mathbf{b} \cdot \mathbf{n} v, w)_{\partial \Omega} \quad \forall v, w \in H^1(\Omega), \quad (2.2)$$

where the notation $(v, w)_D := \int_D v w$ is used for an open subdomain $D \subseteq \Omega$ or its boundary ∂D and for integrable functions v and w . Henceforth, $\|v\|_D$ denotes the norm $\|v\|_{L^2(D)} = \sqrt{(v, v)_D}$. We will drop the subscript when $D = \Omega$.

2.2. Poincaré inequalities

The Poincaré inequality states that

$$\|v - \bar{v}\|_D \leq h_D C_{P,D} \|\nabla v\|_D \quad \forall v \in H^1(D), \quad (2.3a)$$

with $C_{P,D} > 0$ a generic constant, in particular equal to $1/\pi$ for convex $D \subset \Omega$; here \bar{v} is the mean value of v over D defined as $\bar{v} = (v, 1)_D / |D|$ and h_D is the diameter of D . Similarly, another Poincaré (sometimes called Friedrichs) inequality states that

$$\|v\|_D \leq h_D C_{P,D,\Gamma_D} \|\nabla v\|_D, \quad \forall v \in \{H^1(D), v|_{\Gamma_D} = 0, |\Gamma_D| \neq 0\}, \quad (2.3b)$$

where $\Gamma_D \subset \partial D$; typically $C_{P,D,\Gamma_D} = 1$. Henceforth, we will use $C_{PF,D}$ as a general notation for both $C_{P,D}$ and C_{P,D,Γ_D} . It follows from the above that for a one-dimensional interval D , $C_{PF,D}$ can be taken as 1.

2.3. Ultra-weak variational formulation and residual

The variational framework hinges upon an appropriate choice of the trial and test spaces and their corresponding norms. In particular, it turns out natural to work on spaces well-suited to the non-symmetric structure of the problem. Here we consider Hilbert spaces (non-symmetric formulations in Banach spaces can be found in [8, 26]).

The (usual) weak formulation of (1.1) reads: find $u \in H^1_-(\Omega)$ such that

$$(\mathbf{b} \cdot \nabla u, v) = (f, v) \quad \forall v \in L^2(\Omega). \tag{2.4}$$

It is classically well-posed as one might confer with [18, 24], and Proposition 6 of [31], cf. also [14] and Remark 2.2 of [12]. Here, we rather adopt the so-called *ultra-weak* formulation of problem (1.1) where the bilinear form is obtained by casting the derivatives on the test function, using integration-by-parts. It reads: find $u \in L^2(\Omega)$ such that

$$-(u, \mathbf{b} \cdot \nabla v) = (f, v) \quad \forall v \in H^1_+(\Omega). \tag{2.5}$$

The well-posedness of (2.5) can be shown by inf-sup arguments (cf. [14], Thm. 2.6 and [12], Thm. 2.4).

Denote by $H^1_+(\Omega)'$ the dual space to $H^1_+(\Omega)$. For an arbitrary $u_h \in L^2(\Omega)$, the formulation (2.5) leads to the definition of the residual $\mathcal{R}(u_h)$, a bounded linear functional on $H^1_+(\Omega)'$, by

$$\langle \mathcal{R}(u_h), v \rangle := (f, v) + (u_h, \mathbf{b} \cdot \nabla v) \quad \forall v \in H^1_+(\Omega). \tag{2.6}$$

We define its velocity-scaled dual norm by

$$\|\mathcal{R}(u_h)\|_{\mathbf{b}; H^1_+(\Omega)'} := \sup_{v \in H^1_+(\Omega) \setminus \{0\}} \frac{\langle \mathcal{R}(u_h), v \rangle}{\|\mathbf{b} \cdot \nabla v\|}. \tag{2.7}$$

2.4. Error-residual equivalence

In this section, we present an important connection between the $L^2(\Omega)$ -norm of the error and the residual norm (2.7). To be self-contained, though this is not a new finding of this work, we include a proof of the following proposition:

Proposition 2.1 (Error-residual equivalence). *Let u be the ultra-weak solution of (2.5). Then*

$$\|u - u_h\| = \|\mathcal{R}(u_h)\|_{\mathbf{b}; H^1_+(\Omega)'} \quad \forall u_h \in L^2(\Omega).$$

Proof. The well-posedness of the weak formulation (2.4), for the velocity field $-\mathbf{b}$, implies that for all $v \in L^2(\Omega)$, there exists a unique $z \in H^1_+(\Omega)$ such that

$$-(\mathbf{b} \cdot \nabla z, w) = (v, w) \quad \forall w \in L^2(\Omega).$$

This clearly gives $\|\mathbf{b} \cdot \nabla z\| = \|v\|$. Hence, for any $w \in L^2(\Omega)$, we have

$$\|w\| = \sup_{v \in L^2(\Omega) \setminus \{0\}} \frac{(w, v)}{\|v\|} = \sup_{z \in H^1_+(\Omega) \setminus \{0\}} \frac{-(w, \mathbf{b} \cdot \nabla z)}{\|\mathbf{b} \cdot \nabla z\|},$$

and the claim follows by the choice $w = u - u_h$ and using the definitions (2.5) and (2.6). □

Compared to the similar equivalence provided in Theorem 3.3 of [23], Proposition 2.1 shows a form of *equality* which highlights the optimality of the chosen spaces and norms. This is advantageous for the sharpness of the *a posteriori* error estimation.

3. EXAMPLES OF NUMERICAL METHODS

Let $\mathcal{T}_h = \{K\}$ be a mesh of Ω , *i.e.*, a division of the one-dimensional domain Ω into non-overlapping intervals covering Ω , shape regular in the sense that two neighboring intervals are of comparable size, up to a constant $\kappa_{\mathcal{T}_h}$. Let us denote $h_K := \text{diam}(K)$ and $h := \max_{K \in \mathcal{T}_h} h_K$. We also denote by $\mathcal{E}_h := \cup_{K \in \mathcal{T}_h} \partial K$ the skeleton of the triangulation \mathcal{T}_h , coinciding with the set of mesh vertices \mathcal{V}_h in the present one-dimensional case. Moreover, we need to consider the decompositions $\mathcal{E}_h = \mathcal{E}_h^{\text{int}} \cup \mathcal{E}_h^{\text{bnd}}$ into internal and boundary faces and $\mathcal{V}_h = \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega} \cup \mathcal{V}_h^{\partial+\Omega}$ into internal, inflow, and outflow vertices, so that in the one-dimensional case $\mathcal{E}_h = \mathcal{V}_h$ and $\mathcal{E}_h^{\text{bnd}} = \mathcal{V}_h^{\partial-\Omega} \cup \mathcal{V}_h^{\partial+\Omega}$. Let $\mathcal{P}^k(\mathcal{T}_h)$ denote piecewise polynomial functions of at most degree k on the mesh \mathcal{T}_h . The following three numerical methods are classical examples of discretizations of (1.1). Please note that in Examples 3.1 and 3.3, we exclude the lowest polynomial degrees. We need to do so to comply with the orthogonality condition in Assumption 4.1, see Lemma 6.1 below.

The first finite element scheme is a finite-dimensional version of the weak formulation (2.4):

Example 3.1 (Continuous trial Petrov–Galerkin (PG1) finite element). Find $u_h \in X_h := H^1_-(\Omega) \cap \mathcal{P}^k(\mathcal{T}_h)$, $k \geq 2$, such that

$$(\mathbf{b} \cdot \nabla u_h, v_h) = (f, v_h) \quad \forall v_h \in Y_h := \mathcal{P}^{k-1}(\mathcal{T}_h). \tag{3.1}$$

The second finite element scheme stems from the ultra-weak formulation (2.5):

Example 3.2 (Discontinuous trial Petrov–Galerkin (PG2) finite element). Find $u_h \in X_h := \mathcal{P}^k(\mathcal{T}_h)$, $k \geq 0$, such that

$$-(u_h, \mathbf{b} \cdot \nabla v_h) = (f, v_h) \quad \forall v_h \in Y_h := H^1_+(\Omega) \cap \mathcal{P}^{k+1}(\mathcal{T}_h). \tag{3.2}$$

Finally, the dG method for problem (1.1) (letting ∇ also denote the broken (elementwise) gradient) reads:

Example 3.3 (dG finite element). Find $u_h \in X_h := \mathcal{P}^k(\mathcal{T}_h)$, $k \geq 1$, such that

$$\mathcal{B}_h(u_h, v_h) = (f, v_h) \quad \forall v_h \in Y_h := \mathcal{P}^k(\mathcal{T}_h), \tag{3.3a}$$

where

$$\begin{aligned} \mathcal{B}_h(u_h, v_h) := & - \sum_{K \in \mathcal{T}_h} (u_h, \mathbf{b} \cdot \nabla v_h)_K \\ & - \sum_{e \in \mathcal{E}_h^{\text{int}}} \mathbf{b} \cdot \mathbf{n} \{ \{ u_h \} \} [v_h] + \sum_{e \in \mathcal{E}_h^{\text{int}}} \frac{1}{2} | \mathbf{b} \cdot \mathbf{n} | [\{ u_h \}] [v_h] + \sum_{e \in \mathcal{E}_h^{\text{bnd}}} (\mathbf{b} \cdot \mathbf{n})^+ u_h v_h. \end{aligned} \tag{3.3b}$$

Here the notation u_h^- and u_h^+ stands for the trace value on a vertex from left and from right, respectively, the average is defined as $\{ \{ u_h \} \} := (u_h^- + u_h^+)/2$, and the jump is defined as $[\{ u_h \}] := u_h^+ - u_h^-$. In this formulation, the upwind dG flux is applied on the cell interfaces.

4. MAIN RESULTS

We first present here the heart of our approach, a local potential reconstruction on the patch level. We then collect and discuss our main results.

4.1. Patchwise potential reconstruction

Let \mathcal{V}_K be the set of vertices of a mesh element K and let $\mathcal{T}_\mathbf{a}$ denote the *patch* of all simplices which share the given vertex \mathbf{a} , $\mathcal{T}_\mathbf{a} := \{K, \mathbf{a} \in \mathcal{V}_K\}$. Let $\omega_\mathbf{a}$ be the corresponding open subdomain with $h_{\omega_\mathbf{a}} := \text{diam}(\omega_\mathbf{a})$. Then $\cup_{\mathbf{a} \in \mathcal{V}_h} \omega_\mathbf{a}$ forms an overlapping partition of Ω , with $\mathcal{N} = 2$ maximal overlap in one space dimension. For all $\mathbf{a} \in \mathcal{V}_h$, let $\psi_\mathbf{a} \in H^1(\Omega) \cap \mathcal{P}^1(\mathcal{T}_h)$ be the piecewise affine hat function, taking value 1 in vertex \mathbf{a} and 0 in all other vertices. The hat functions verify $\text{supp}(\psi_\mathbf{a}) = \overline{\omega_\mathbf{a}}$ and form a *partition of unity* as

$$\sum_{\mathbf{a} \in \mathcal{V}_h} \psi_\mathbf{a} = 1. \tag{4.1}$$

The following assumption on the $\psi_\mathbf{a}$ -orthogonality of the residual will be crucial to localize the error:

Assumption 4.1 ($\psi_\mathbf{a}$ -orthogonality). *The residual $\mathcal{R}(u_h) \in H^1_+(\Omega)'$ defined in (2.6) satisfies*

$$\langle \mathcal{R}(u_h), \psi_\mathbf{a} \rangle = (f, \psi_\mathbf{a})_{\omega_\mathbf{a}} + (u_h, \mathbf{b} \cdot \nabla \psi_\mathbf{a})_{\omega_\mathbf{a}} = 0 \quad \forall \mathbf{a} \in \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega}. \tag{4.2}$$

Having Assumption 4.1 satisfied, a local reconstruction technique which provides the key ingredient to evaluate our *a posteriori* error estimator is:

Definition 4.2 (Patchwise potential reconstruction). Let $u_h \in L^2(\Omega)$ satisfy Assumption 4.1. For all vertices $\mathbf{a} \in \mathcal{V}_h$, let $s_h^\mathbf{a} \in X_h^\mathbf{a}$ be the solution of the following advection–reaction problem on the patch $\omega_\mathbf{a}$

$$(\mathbf{b} \cdot \nabla(\psi_\mathbf{a} s_h^\mathbf{a}), v_h)_{\omega_\mathbf{a}} = (f \psi_\mathbf{a} + (\mathbf{b} \cdot \nabla \psi_\mathbf{a}) u_h, v_h)_{\omega_\mathbf{a}} \quad \forall v_h \in Y_h^\mathbf{a}, \tag{4.3}$$

with the finite-dimensional spaces

$$X_h^\mathbf{a} := \mathcal{P}^{k'}(\mathcal{T}_\mathbf{a}) \cap H^1(\omega_\mathbf{a}), \quad Y_h^\mathbf{a} := \mathcal{P}^{k'}(\mathcal{T}_\mathbf{a}),$$

and $k' \geq 0$. Define the global reconstruction s_h by

$$s_h := \sum_{\mathbf{a} \in \mathcal{V}_h} \psi_\mathbf{a} s_h^\mathbf{a}. \tag{4.4}$$

4.2. Main results

Our guaranteed upper bound on the L^2 -error can be presented as the following theorem:

Theorem 4.3 (Guaranteed *a posteriori* error estimate). *Let $u \in L^2(\Omega)$ be the ultra-weak solution of (2.5) and let $u_h \in L^2(\Omega)$ be arbitrary subject to the $\psi_\mathbf{a}$ -orthogonality in Assumption 4.1. Furthermore, consider s_h to be the reconstruction from Definition 4.2 with $k' \geq 0$. Then*

$$\|u - u_h\| \leq \eta := \left\{ \sum_{K \in \mathcal{T}_h} (\eta_{\text{NC},K} + \eta_{\text{Osc},K})^2 \right\}^{1/2},$$

where

$$\eta_{\text{NC},K} := \|u_h - s_h\|_K$$

and the data oscillation estimator is given as

$$\eta_{\text{Osc},K} := \frac{h_K}{\pi |\mathbf{b}|} \|(I - \Pi_{\mathcal{P}^{k'}(\mathcal{T}_h)})f\|_K, \tag{4.5}$$

with $\Pi_{\mathcal{P}^{k'}(\mathcal{T}_h)}$ the $L^2(\Omega)$ -orthogonal projection onto $\mathcal{P}^{k'}(\mathcal{T}_h)$.

The lower bound on the error and main theorem on local efficiency as well as robustness is presented in the following theorem:

Theorem 4.4 (Local efficiency and robustness). *Let $u \in L^2(\Omega)$ be the ultra-weak solution of (2.5) and let $u_h \in \mathcal{P}^k(\mathcal{T}_h)$, $k \geq 0$, be arbitrary subject to Assumption 4.1. Consider s_h as obtained by Definition 4.2 with $k' \geq k$ and $\eta_{\text{NC},K}$ as defined in Theorem 4.3. Then, for all the mesh elements $K \in \mathcal{T}_h$, the following holds:*

$$\eta_{\text{NC},K} \leq C_{\text{cont,PF}} \sum_{\alpha \in \mathcal{V}_K} \|u - u_h\|_{\omega_\alpha} + \sum_{\alpha \in \mathcal{V}_K} \frac{h_{\omega_\alpha}}{\pi |\mathbf{b}|} \|(I - \Pi_{\mathcal{P}^{k'}(\mathcal{T}_\alpha)})(f\psi_\alpha)\|_{\omega_\alpha}.$$

Here, $C_{\text{cont,PF}}$ is a generic positive constant that only depends on the mesh shape-regularity constant $\kappa_{\mathcal{T}_h}$ via

$$C_{\text{cont,PF}} := \max_{\alpha \in \mathcal{V}_h} (1 + C_{\text{PF},\omega_\alpha} h_{\omega_\alpha} \|\nabla \psi_\alpha\|_\infty). \quad (4.6)$$

Provided that $\psi_\alpha f$ is piecewise polynomial, one can obtain the global efficiency of the error indicator as an immediate consequence of Theorem 4.4 as:

Corollary 4.5 (Global efficiency and maximal overestimation). *Let the assumptions of Theorem 4.4 be verified and assume in addition that $\psi_\alpha f \in Y_h^\alpha$ for all the mesh vertices $\alpha \in \mathcal{V}_h$. Then*

$$\|u_h - s_h\| \leq 2C_{\text{cont,PF}} \|u - u_h\|.$$

4.3. Remarks

A few remarks are in order.

Remark 4.6 (Potential reconstruction and its local conservation). Lemma 5.4 below shows that

$$s_h \in \mathcal{P}^{k'+1}(\mathcal{T}_h) \cap H_-^1(\Omega), \quad (4.7)$$

i.e., it lies in a natural finite-dimensional functional space corresponding to the weak formulation (2.4). Moreover, the following orthogonality is satisfied

$$(f - \mathbf{b} \cdot \nabla s_h, v_h)_K = 0 \quad \forall v_h \in \mathcal{P}^{k'}(K), \quad \forall K \in \mathcal{T}_h. \quad (4.8)$$

Remark 4.7 (Lifting of the local residual). The potential reconstruction s_h^α of Definition 4.2 is such that the hat-function-weighted difference $\psi_\alpha(s_h^\alpha - u_h)$ is a lifting of the local hat-function-weighted residual $\langle \mathcal{R}(u_h), \psi_\alpha \cdot \rangle$ by a local advection problem. Indeed, let $v_h \in Y_h^\alpha \cap H^1(\omega_\alpha)$, $v_h(\mathbf{a}) = 0$ when $\mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}$. Then integration-by-parts, the property $s_h^\alpha|_{\partial-\Omega} = 0$ from Remark 4.6, and definition (2.6) of the residual give

$$\begin{aligned} (\psi_\alpha(u_h - s_h^\alpha), \mathbf{b} \cdot \nabla v_h)_{\omega_\alpha} &= (\psi_\alpha u_h, \mathbf{b} \cdot \nabla v_h)_{\omega_\alpha} + (\mathbf{b} \cdot \nabla(\psi_\alpha s_h^\alpha), v_h)_{\omega_\alpha} \\ &= (\psi_\alpha u_h, \mathbf{b} \cdot \nabla v_h)_{\omega_\alpha} + (f\psi_\alpha + (\mathbf{b} \cdot \nabla \psi_\alpha)u_h, v_h)_{\omega_\alpha} \\ &= (f, \psi_\alpha v_h)_{\omega_\alpha} + (u_h, \mathbf{b} \cdot \nabla(\psi_\alpha v_h))_{\omega_\alpha} = \langle \mathcal{R}(u_h), \psi_\alpha v_h \rangle. \end{aligned}$$

Remark 4.8 (Data oscillation). We call the estimator $\eta_{\text{Osc},K}$ in (4.5) “data oscillation” for the following reason: if u_h is piecewise polynomial of degree $k \geq 0$, the error $\|u - u_h\|$ may converge as $\mathcal{O}(h^{k+1})$. By choosing $k' \geq k$ one obtains, for sufficiently piecewise smooth data f , the higher convergence order $\mathcal{O}(h^{k'+2})$ for $\eta_{\text{Osc},K}$.

Remark 4.9 (Non-homogeneous boundary condition). For the sake of simplicity, we have just presented the case of a homogenous Dirichlet boundary condition. A non-homogeneous boundary condition $u = g$ on $\partial-\Omega$ is handled in the ultra-weak formulation (2.5) by subtracting $\mathbf{b} \cdot \mathbf{n}|_{\partial-\Omega}(g, v)_{\partial-\Omega}$ from the right-hand side, and in the definition of residual (2.6) as well. Employing Definition 4.2 leads to a reconstruction s_h which satisfies the boundary condition $s_h = g$ on $\partial-\Omega$, following the same arguments as those of Lemma 5.4. In the one-dimensional case, the boundary is a point, and its values can be captured in the finite-dimensional setting without any error.

Remark 4.10 (Specificity of the one-dimensional case). The solution of the one-dimensional problem (1.1) can actually be obtained by integration of the right-hand side. Importantly, no step above uses this fact. For this reason, most of the developments extend to multiple space dimensions, as we show in Section 10 below. Two specific points, though, do not seem to easily extend to multiple space dimensions, namely the reconstruction from Definition 4.2 and the use of the inverse operator in (8.2) in the proof of Theorem 4.4 in Section 8 below. There is a possible hope to overcome the latter obstacle with the multi-dimensional developments as for elliptic operators in [7, 15, 16]. The former obstacle, though, seems to be the true bottleneck that is not fully overcome in Definition 10.5 below.

Remark 4.11 (Extension to advection–reaction). In the advection–reaction problem, one replaces equation (1.1a) by $\mathbf{b} \cdot \nabla u + cu = f$, where $c \geq 0$ is a constant. Our results can be extended to this case, while providing a guaranteed *a posteriori* error estimate that is locally efficient and robust with respect to the interplay of the advection and reaction phenomena (mutual sizes of the constants \mathbf{b} and c), although the polynomial degree robustness is theoretically lost. Namely, in Definition 4.2, one merely replaces f by $f - cu_h$ in the right-hand side of (4.3). The analysis of the extension is, however, not straightforward, and is the subject of a stand-alone work [35].

5. CUT-OFF ESTIMATES, ERROR LOCALIZATION, AND WELL-POSEDNESS OF THE PATCHWISE PROBLEMS

In this section, we show that under Assumption 4.1, one can obtain a two-sided bound on the dual norm of the residual $\|\mathcal{R}(u_h)\|_{\mathbf{b}, H_+^1(\Omega)'}$ by identifying some (infinite-dimensional) problems on patches of elements around vertices. This identification allows us to localize the error. We then prove the well-posedness of the patchwise problems from Definition 4.2, as motivated by this localization.

5.1. Cut-off estimates

Similarly to (2.1), let $H_+^1(\omega_{\mathbf{a}})$ contain those functions from $H^1(\omega_{\mathbf{a}})$ whose trace is zero on the outflow boundary of $\omega_{\mathbf{a}}$. Define two patchwise spaces

$$H_{\#}^1(\omega_{\mathbf{a}}) := \begin{cases} H_0^1(\omega_{\mathbf{a}}), & \mathbf{a} \notin \mathcal{V}_h^{\partial-\Omega}, \\ H_+^1(\omega_{\mathbf{a}}), & \mathbf{a} \in \mathcal{V}_h^{\partial-\Omega}, \end{cases} \tag{5.1}$$

and

$$H_*^1(\omega_{\mathbf{a}}) := \begin{cases} \{H^1(\omega_{\mathbf{a}}) : (v, 1)_{\omega_{\mathbf{a}}} = 0\}, & \mathbf{a} \notin \mathcal{V}_h^{\partial+\Omega}, \\ H_+^1(\omega_{\mathbf{a}}), & \mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}. \end{cases} \tag{5.2}$$

In the sequel, we will use several times the following fact:

$$v \in H_*^1(\omega_{\mathbf{a}}) \implies \psi_{\mathbf{a}} v \in H_{\#}^1(\omega_{\mathbf{a}}). \tag{5.3}$$

Recall the constant $C_{\text{cont,PF}}$ from (4.6). The following important cut-off Poincaré estimate follows immediately from Theorem 3.1 of [10] or Section 3 of [7], cf. also Lemma 3.12 of [15], using that the present one-dimensional setting, \mathbf{b} is a constant scalar:

Lemma 5.1 (Local cut-off estimate). *For any mesh vertex $\mathbf{a} \in \mathcal{V}_h$, we have*

$$\|\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} v)\|_{\omega_{\mathbf{a}}} \leq C_{\text{cont,PF}} \|\mathbf{b} \cdot \nabla v\|_{\omega_{\mathbf{a}}} \quad \forall v \in H_*^1(\omega_{\mathbf{a}}).$$

5.2. Error localization

We have seen in Remark 4.7 that the reconstruction of Definition 4.2 is based on the $\psi_{\mathbf{a}}$ -orthogonality Assumption 4.1 and builds upon lifting the localized residual $\langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \cdot \rangle$. This is tightly connected with an equivalent, localized expression of the error/residual (recall Prop. 5.1). Define the restriction of $\mathcal{R}(u_h)$ from (2.6) to the space $H_{\#}^1(\omega_{\mathbf{a}})$ as

$$\|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'} := \sup_{v \in H_{\#}^1(\omega_{\mathbf{a}}) \setminus \{0\}} \frac{\langle \mathcal{R}(u_h), v \rangle}{\|\mathbf{b} \cdot \nabla v\|_{\omega_{\mathbf{a}}}}. \tag{5.4}$$

We then have:

Proposition 5.2 (Localizaion of residual dual norms with $\psi_{\mathbf{a}}$ -orthogonality). *Provided $\mathcal{R}(u_h)$ satisfies Assumption 4.1, we have*

$$\|\mathcal{R}(u_h)\|_{\mathbf{b}; H_+^1(\Omega)'}^2 \leq 2C_{\text{cont,PF}}^2 \sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'}^2. \tag{5.5a}$$

Independently of Assumption 4.1, the following always holds true:

$$\sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'}^2 \leq 2\|\mathcal{R}(u_h)\|_{\mathbf{b}; H_+^1(\Omega)'}^2. \tag{5.5b}$$

Proof. The proof proceeds along the lines in [6,7,10,15]. In particular, noting the partition of unity property (4.1) and the $\psi_{\mathbf{a}}$ -orthogonality of Assumption 4.1, one can use $v = \sum_{\mathbf{a} \in \mathcal{V}_h} \psi_{\mathbf{a}} v$ as the test function to obtain, for each $v \in H_+^1(\Omega)$,

$$\langle \mathcal{R}(u_h), v \rangle \stackrel{(4.1),(4.2)}{=} \sum_{\mathbf{a} \in \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega}} \langle \mathcal{R}(u_h), \psi_{\mathbf{a}}(v - \bar{v}_{\mathbf{a}}) \rangle + \sum_{\mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}} \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} v \rangle,$$

where $\bar{v}_{\mathbf{a}}$ is the mean value of v on $\omega_{\mathbf{a}}$. Let $w_{\mathbf{a}} := (v - \bar{v}_{\mathbf{a}})|_{\omega_{\mathbf{a}}}$ if $\mathbf{a} \in \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega}$ and $w_{\mathbf{a}} := v|_{\omega_{\mathbf{a}}}$ if $\mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}$. Then, $w_{\mathbf{a}} \in H_{\#}^1(\omega_{\mathbf{a}})$, so that $\psi_{\mathbf{a}} w_{\mathbf{a}} \in H_{\#}^1(\omega_{\mathbf{a}})$ by (5.3). Using the cut-off estimate of Lemma 5.1 for $v = w_{\mathbf{a}}$ and the definition (5.4), one can in particular obtain

$$\begin{aligned} \langle \mathcal{R}(u_h), v \rangle &= \sum_{\mathbf{a} \in \mathcal{V}_h} \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} w_{\mathbf{a}} \rangle \leq \sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'} \|\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} w_{\mathbf{a}})\|_{\omega_{\mathbf{a}}} \\ &\leq C_{\text{cont,PF}} \sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'} \|\mathbf{b} \cdot \nabla v\|_{\omega_{\mathbf{a}}} \\ &\stackrel{\text{C.S.}}{\leq} C_{\text{cont,PF}} 2^{1/2} \left(\sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'}^2 \right)^{1/2} \|\mathbf{b} \cdot \nabla v\|, \end{aligned}$$

which gives (5.5a).

To prove (5.5b), using the Riesz representation theorem one observes that there exists $\xi_{\mathbf{a}} \in H_{\#}^1(\omega_{\mathbf{a}})$ such that

$$(\mathbf{b} \cdot \nabla \xi_{\mathbf{a}}, \mathbf{b} \cdot \nabla v)_{\omega_{\mathbf{a}}} = \langle \mathcal{R}(u_h), v \rangle \quad \forall v \in H_{\#}^1(\omega_{\mathbf{a}}). \tag{5.6}$$

Consequently $\|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'} = \|\mathbf{b} \cdot \nabla \xi_{\mathbf{a}}\|_{\omega_{\mathbf{a}}}$. By extending $\xi_{\mathbf{a}}$ by zero outside of the patch $\omega_{\mathbf{a}}$ and defining $\sum_{\mathbf{a} \in \mathcal{V}_h} \xi_{\mathbf{a}} =: \xi \in H_+^1(\Omega)$, one has

$$\begin{aligned} \sum_{\mathbf{a} \in \mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})'}^2 &= \sum_{\mathbf{a} \in \mathcal{V}_h} (\mathbf{b} \cdot \nabla \xi_{\mathbf{a}}, \mathbf{b} \cdot \nabla \xi_{\mathbf{a}})_{\omega_{\mathbf{a}}} \stackrel{(5.6)}{=} \sum_{\mathbf{a} \in \mathcal{V}_h} \langle \mathcal{R}(u_h), \xi_{\mathbf{a}} \rangle \\ &= \langle \mathcal{R}(u_h), \xi \rangle \leq \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_+^1(\Omega)'} \|\mathbf{b} \cdot \nabla \xi\|. \end{aligned} \tag{5.7}$$

The application of the Cauchy–Schwarz inequality gives

$$\begin{aligned} \|\mathbf{b}\cdot\nabla\xi\|^2 &= \sum_{K\in\mathcal{T}_h} \|\mathbf{b}\cdot\nabla\xi\|_K^2 = \sum_{K\in\mathcal{T}_h} \left\| \sum_{\mathbf{a}\in\mathcal{V}_K} \mathbf{b}\cdot\nabla\xi_{\mathbf{a}} \right\|_K^2 \leq 2 \sum_{K\in\mathcal{T}_h} \sum_{\mathbf{a}\in\mathcal{V}_K} \|\mathbf{b}\cdot\nabla\xi_{\mathbf{a}}\|_K^2 \\ &= 2 \sum_{\mathbf{a}\in\mathcal{V}_h} \sum_{K\in\mathcal{T}_{\mathbf{a}}} \|\mathbf{b}\cdot\nabla\xi_{\mathbf{a}}\|_K^2 = 2 \sum_{\mathbf{a}\in\mathcal{V}_h} \|\mathbf{b}\cdot\nabla\xi_{\mathbf{a}}\|_{\omega_{\mathbf{a}}}^2 = 2 \sum_{\mathbf{a}\in\mathcal{V}_h} \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_{\#}^1(\omega_{\mathbf{a}})}^2, \end{aligned}$$

which proves (5.5b) in combination with (5.7). □

5.3. Well-posedness of the local problems

In order to use the reconstruction proposed in Definition 4.2, it is important to make sure of its well-posedness. We check it now.

A priori, the number of degrees of freedom in $X_h^{\mathbf{a}}$ and $Y_h^{\mathbf{a}}$ for an interior vertex $\mathbf{a} \in \mathcal{V}_h^{\text{int}}$ does not match; while there exist $2(k' + 1)$ linearly independent test functions in $Y_h^{\mathbf{a}}$, the trial space $X_h^{\mathbf{a}}$ has only $2k' + 1$ degrees of freedom. For any $\mathbf{a} \in \mathcal{V}_h^{\text{int}}$, though, the test function in (4.3) given by $v_h = 1$ on both elements $K \in \mathcal{T}_{\mathbf{a}}$ is actually superfluous. Indeed, on the one hand, we have

$$(\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}s_h^{\mathbf{a}}), 1)_{\omega_{\mathbf{a}}} = (\mathbf{b}\cdot\mathbf{n}, \psi_{\mathbf{a}}s_h^{\mathbf{a}})_{\partial\omega_{\mathbf{a}}} = 0, \tag{5.8}$$

according to the definition of $\psi_{\mathbf{a}}$. On the other hand, Assumption 4.1 guarantees that the right-hand side vanishes in such a case, hence

$$(f\psi_{\mathbf{a}} + (\mathbf{b}\cdot\nabla\psi_{\mathbf{a}})u_h, 1)_{\omega_{\mathbf{a}}} = \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle = 0.$$

Then, we can show that the solution of (4.3) uniquely exists and the proposed reconstruction is well-posed:

Lemma 5.3 (Well-posedness of Def. 4.2). *There exists a unique solution $s_h^{\mathbf{a}} \in X_h^{\mathbf{a}}$ of problem (4.3).*

Proof. For all $w_h \in X_h^{\mathbf{a}}$, $\psi_{\mathbf{a}}w_h \in H_0^1(\omega_{\mathbf{a}})$ for all $\mathbf{a} \in \mathcal{V}_h^{\text{int}}$, $\psi_{\mathbf{a}}w_h \in H_+^1(\omega_{\mathbf{a}})$ for $\mathbf{a} \in \mathcal{V}_h^{\partial-\Omega}$, and $\psi_{\mathbf{a}}w_h \in H_-^1(\omega_{\mathbf{a}})$ for $\mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}$; hence $\|\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}\cdot)\|_{\omega_{\mathbf{a}}}$ is a norm on $X_h^{\mathbf{a}}$. Noting that \mathbf{b} is constant and $\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}w_h) \in \mathcal{P}^{k'}(\mathcal{T}_{\mathbf{a}}) = Y_h^{\mathbf{a}}$, one can write the inf–sup condition of the bilinear form associated with the left-hand side of (4.3) as

$$\sup_{v_h \in Y_h^{\mathbf{a}} \setminus \{0\}} \frac{(\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}w_h), v_h)_{\omega_{\mathbf{a}}}}{\|v_h\|_{\omega_{\mathbf{a}}}} = \|\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}w_h)\|_{\omega_{\mathbf{a}}} \quad \forall w_h \in X_h^{\mathbf{a}},$$

with unit inf–sup constant. Since the dimensions of $X_h^{\mathbf{a}}$ and $Y_h^{\mathbf{a}}$ that count are the same, this injectivity implies the bijectivity of the operator. □

Lemma 5.4 (Properties of the reconstruction). *Definition 4.2 yields s_h satisfying (4.7) and (4.8).*

Proof. For (4.7) is clear from (4.4) that $s_h \in \mathcal{P}^{k'+1}(\mathcal{T}_h) \cap H^1(\Omega)$, and we only need to show that s_h satisfies the boundary condition requirement of the space $H_-^1(\Omega)$, i.e., $s_h|_{\partial-\Omega} = 0$. We check this by showing that $s_h^{\mathbf{a}}|_{\partial\omega_{\mathbf{a}} \cap \partial-\Omega} = 0$ for $\mathbf{a} \in \mathcal{V}_h^{\partial-\Omega}$. We see from (4.3) and Assumption 4.1 that

$$(\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}s_h^{\mathbf{a}}), 1)_{\omega_{\mathbf{a}}} = (f\psi_{\mathbf{a}} + (\mathbf{b}\cdot\nabla\psi_{\mathbf{a}})u_h, 1)_{\omega_{\mathbf{a}}} = \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle = 0,$$

so that the requested equality follows from integration-by-parts similarly to (5.8),

$$(\mathbf{b}\cdot\nabla(\psi_{\mathbf{a}}s_h^{\mathbf{a}}), 1)_{\omega_{\mathbf{a}}} = \mathbf{b}\cdot\mathbf{n}s_h^{\mathbf{a}}|_{\partial-\Omega},$$

and since $\mathbf{b} \cdot \mathbf{n} \neq 0$ on $\partial_- \Omega$ by definition.

To prove (4.8), first note that $\sum_{\mathbf{a} \in \mathcal{V}_K} \psi_{\mathbf{a}}|_K = 1$ and $\sum_{\mathbf{a} \in \mathcal{V}_K} (\mathbf{b} \cdot \nabla \psi_{\mathbf{a}}) u_h|_K = 0$. Thus, since $Y_h^{\mathbf{a}}|_K = \mathcal{P}^{k'}(K)$, extending the function $v_h \in \mathcal{P}^{k'}(K)$ by zero outside K , and using respectively definitions (4.4) of s_h and (4.3) of $s_h^{\mathbf{a}}$, one has

$$\begin{aligned} (f - \mathbf{b} \cdot \nabla s_h, v_h)_K &= \left(\sum_{\mathbf{a} \in \mathcal{V}_K} \{ \psi_{\mathbf{a}} f + (\mathbf{b} \cdot \nabla \psi_{\mathbf{a}}) u_h - \mathbf{b} \cdot \nabla (\psi_{\mathbf{a}} s_h^{\mathbf{a}}) \}, v_h \right)_K \\ &= \sum_{\mathbf{a} \in \mathcal{V}_K} (\psi_{\mathbf{a}} f + (\mathbf{b} \cdot \nabla \psi_{\mathbf{a}}) u_h - \mathbf{b} \cdot \nabla (\psi_{\mathbf{a}} s_h^{\mathbf{a}}), v_h)_{\omega_{\mathbf{a}}} = 0. \end{aligned}$$

□

6. $\psi_{\mathbf{a}}$ -ORTHOGONALITY OF THE RESIDUAL FOR THE METHODS OF SECTION 3

We now return to the three methods presented in Section 3 and show the validity of Assumption 4.1 for them:

Lemma 6.1 ($\psi_{\mathbf{a}}$ -orthogonality of the residual). *For methods PG1 of Example 3.1 with $k \geq 2$, PG2 of Example 3.2 with $k \geq 0$, and dG of Example 3.3 with $k \geq 1$, Assumption 4.1 holds true.*

Proof. Let $\mathbf{a} \in \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega}$. We verify the condition for each method:

- From definition (2.6), for the PG1 method (3.1), we have

$$\begin{aligned} \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle &= \sum_{K \in \mathcal{T}_{\mathbf{a}}} \{ (f, \psi_{\mathbf{a}})_K + (u_h, \mathbf{b} \cdot \nabla \psi_{\mathbf{a}})_K \} \\ &\stackrel{\text{I.B.P.}}{=} \sum_{K \in \mathcal{T}_{\mathbf{a}}} \{ (f, \psi_{\mathbf{a}})_K - (\mathbf{b} \cdot \nabla u_h, \psi_{\mathbf{a}})_K + (\mathbf{b} \cdot \mathbf{n} u_h, \psi_{\mathbf{a}})_{\partial K} \}. \end{aligned} \tag{6.1}$$

For all $\mathbf{a} \in \mathcal{V}_h^{\text{int}}$, the jump $[[\psi_{\mathbf{a}}]]$ vanishes at the vertex \mathbf{a} and $\psi_{\mathbf{a}} = 0$ on the boundary face of the patch. Hence, since u_h is also continuous in \mathbf{a} in the PG1 method, the last term in (6.1) disappears and one infers that

$$\langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle = (f, \psi_{\mathbf{a}}) - (\mathbf{b} \cdot \nabla u_h, \psi_{\mathbf{a}}) \stackrel{(3.1)}{=} 0,$$

since we assume $k \geq 2$, so that $\psi_{\mathbf{a}} \in Y_h$. The same result is valid for $\mathbf{a} \in \mathcal{V}_h^{\partial-\Omega}$ since $u_h = 0$ on the inflow as imposed in the definition of X_h .

- From definition (2.6) and employing the PG2 characterization (3.2), we obtain in a straightforward manner that

$$\langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle = 0$$

for all $k \geq 0$.

- For the dG method (3.3), noting that $(\mathbf{b} \cdot \mathbf{n})^+ = 0$ on the inflow and using the same arguments on the vanishing of the jump $[[\psi_{\mathbf{a}}]]$ and some $k \geq 1$ by assumption, we have for any vertex $\mathbf{a} \in \mathcal{V}_h^{\text{int}} \cup \mathcal{V}_h^{\partial-\Omega}$

$$\sum_{e \in \mathcal{E}_h^{\text{int}}} \left\{ \frac{1}{2} |\mathbf{b} \cdot \mathbf{n}| [[u_h]] - \mathbf{b} \cdot \mathbf{n} \{ \{ u_h \} \} \right\} [[\psi_{\mathbf{a}}]] + \sum_{e \in \mathcal{E}_h^{\text{bnd}}} (\mathbf{b} \cdot \mathbf{n})^+ u_h \psi_{\mathbf{a}} = 0.$$

Hence, also employing definition (2.6), we infer that

$$\langle \mathcal{R}(u_h), \psi_{\mathbf{a}} \rangle = \sum_{K \in \mathcal{T}_{\mathbf{a}}} \{ (f, \psi_{\mathbf{a}})_K + (u_h, \mathbf{b} \cdot \nabla \psi_{\mathbf{a}})_K \} = 0$$

for all $k \geq 1$ which implies $\psi_{\mathbf{a}} \in Y_h$.

□

7. PROOF OF GUARANTEED RELIABILITY (THM. 4.3)

We prove here Theorem 4.3. Since $s_h \in H^1_-(\Omega)$ by Lemma 5.4, for any $v \in H^1_+(\Omega)$ the integration-by-parts formula (2.2) implies that

$$(s_h, \mathbf{b} \cdot \nabla v) + (\mathbf{b} \cdot \nabla s_h, v) = (s_h, v_h \mathbf{b} \cdot \mathbf{n}) = 0. \tag{7.1}$$

By using the error–residual identity of Proposition 2.1, definitions (2.6) and (2.7), and the above equality, one can write

$$\|u - u_h\|_\Omega = \|\mathcal{R}(u_h)\|_{\mathbf{b}; H^1_+(\Omega)'} = \sup_{v \in H^1_+(\Omega) \setminus \{0\}} \frac{(f - \mathbf{b} \cdot \nabla s_h, v) + (u_h - s_h, \mathbf{b} \cdot \nabla v)}{\|\mathbf{b} \cdot \nabla v\|}.$$

Owing to (4.8), denoting by \bar{v}_K the mean value of v over the element K , we infer that

$$\begin{aligned} \|u - u_h\| &\stackrel{(4.8)}{=} \sup_{v \in H^1_+(\Omega) \setminus \{0\}} \frac{\sum_{K \in \mathcal{T}_h} \left[(u_h - s_h, \mathbf{b} \cdot \nabla v)_K + (f - \mathbf{b} \cdot \nabla s_h, v - \bar{v}_K)_K \right]}{\|\mathbf{b} \cdot \nabla v\|} \\ &\stackrel{(2.3a)}{\leq} \sup_{v \in H^1_+(\Omega) \setminus \{0\}} \frac{\sum_{K \in \mathcal{T}_h} \left[\|u_h - s_h\|_K \|\mathbf{b} \cdot \nabla v\|_K + \frac{h_K}{\pi |\mathbf{b}|} \|f - \mathbf{b} \cdot \nabla s_h\|_K \|\mathbf{b} \cdot \nabla v\|_K \right]}{\|\mathbf{b} \cdot \nabla v\|} \\ &\leq \left\{ \sum_{K \in \mathcal{T}_h} \left[\|u_h - s_h\|_K + \frac{h_K}{\pi |\mathbf{b}|} \|f - \mathbf{b} \cdot \nabla s_h\|_K \right]^2 \right\}^{1/2}. \end{aligned}$$

Noting that $\mathbf{b} \cdot \nabla s_h \in \mathcal{P}^{k'}(\mathcal{T}_h)$, it follows from (4.8) that $\mathbf{b} \cdot \nabla s_h = \Pi_{\mathcal{P}^{k'}(\mathcal{T}_h)} f$ so that

$$\|f - \mathbf{b} \cdot \nabla s_h\|_K = \|(I - \Pi_{\mathcal{P}^{k'}(\mathcal{T}_h)})f\|_K,$$

which completes the proof.

8. PROOF OF EFFICIENCY AND ROBUSTNESS

This section presents proofs of Theorem 4.4 and Corollary 4.5.

8.1. Proof of Theorem 4.4 (local efficiency and robustness)

Fix an element $K \in \mathcal{T}_h$. Noting that $\sum_{\mathbf{a} \in \mathcal{V}_K} \psi_{\mathbf{a}}|_K = 1$ and using definition (4.4), one has

$$\|u_h - s_h\|_K = \left\| \sum_{\mathbf{a} \in \mathcal{V}_K} \psi_{\mathbf{a}}(u_h - s_h^\mathbf{a}) \right\|_K \leq \sum_{\mathbf{a} \in \mathcal{V}_K} \|\psi_{\mathbf{a}}(u_h - s_h^\mathbf{a})\|_{\omega_{\mathbf{a}}}. \tag{8.1}$$

Recalling (5.2), we easily see that for any vertex $\mathbf{a} \in \mathcal{V}_h$, there is a unique $v^\mathbf{a} \in H^1_*(\omega_{\mathbf{a}})$ such that

$$\mathbf{b} \cdot \nabla v^\mathbf{a} = \psi_{\mathbf{a}}(u_h - s_h^\mathbf{a}) \tag{8.2}$$

in $\omega_{\mathbf{a}}$, and $v^\mathbf{a}$ is nonzero unless $\psi_{\mathbf{a}}u_h = \psi_{\mathbf{a}}s_h^\mathbf{a}$, in which case $\|\psi_{\mathbf{a}}(u_h - s_h^\mathbf{a})\|_{\omega_{\mathbf{a}}} = 0$. Moreover, first, $(\psi_{\mathbf{a}}s_h^\mathbf{a})(\mathbf{a}) = s_h(\mathbf{a}) = 0$ when $\mathbf{a} \in \mathcal{V}_h^{\partial-\Omega}$, using (4.7), and, second, $v^\mathbf{a}(\mathbf{a}) = 0$ when $\mathbf{a} \in \mathcal{V}_h^{\partial+\Omega}$, using (5.2). Thus, similarly to (7.1), for any $\mathbf{a} \in \mathcal{V}_h$, we have

$$(\psi_{\mathbf{a}}s_h^\mathbf{a}, \mathbf{b} \cdot \nabla v^\mathbf{a})_{\omega_{\mathbf{a}}} + (\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}}s_h^\mathbf{a}), v^\mathbf{a})_{\omega_{\mathbf{a}}} = 0.$$

From the two above identities, we infer that

$$\begin{aligned}
\|\psi_{\mathbf{a}}(u_h - s_h^{\mathbf{a}})\|_{\omega_{\mathbf{a}}} &= \frac{(\psi_{\mathbf{a}}(u_h - s_h^{\mathbf{a}}), \mathbf{b} \cdot \nabla v^{\mathbf{a}})_{\omega_{\mathbf{a}}}}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} \\
&= \frac{(\psi_{\mathbf{a}} u_h, \mathbf{b} \cdot \nabla v^{\mathbf{a}})_{\omega_{\mathbf{a}}} + (f \psi_{\mathbf{a}} + \mathbf{b} \cdot \nabla \psi_{\mathbf{a}} u_h, v^{\mathbf{a}})_{\omega_{\mathbf{a}}}}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} \\
&\quad + \frac{(\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} s_h^{\mathbf{a}}), v^{\mathbf{a}})_{\omega_{\mathbf{a}}} - (f \psi_{\mathbf{a}} + \mathbf{b} \cdot \nabla \psi_{\mathbf{a}} u_h, v^{\mathbf{a}})_{\omega_{\mathbf{a}}}}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} \\
&=: \text{I} + \text{II}.
\end{aligned} \tag{8.3}$$

For the term I, remark first that from (5.3) and from the definition of $H_{\#}^1(\omega_{\mathbf{a}})$ in (5.1), we have

$$\psi_{\mathbf{a}} v^{\mathbf{a}} \in H_{\#}^1(\omega_{\mathbf{a}}) \subseteq H_+^1(\Omega).$$

Second, recalling the residual definition (2.6) and the ultra-weak formulation (2.5), we have

$$(f \psi_{\mathbf{a}} + \mathbf{b} \cdot \nabla \psi_{\mathbf{a}} u_h, v^{\mathbf{a}})_{\omega_{\mathbf{a}}} + (u_h \psi_{\mathbf{a}}, \mathbf{b} \cdot \nabla v^{\mathbf{a}})_{\omega_{\mathbf{a}}} = \langle \mathcal{R}(u_h), \psi_{\mathbf{a}} v^{\mathbf{a}} \rangle = -(u - u_h, \mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} v^{\mathbf{a}})).$$

Consequently, employing the Cauchy–Schwarz inequality and Lemma 5.1, we infer that

$$\text{I} = \frac{-(u - u_h, \mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} v^{\mathbf{a}}))}{\|\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} v^{\mathbf{a}})\|_{\omega_{\mathbf{a}}}} \frac{\|\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} v^{\mathbf{a}})\|_{\omega_{\mathbf{a}}}}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} \leq C_{\text{cont,PF}} \|u - u_h\|_{\omega_{\mathbf{a}}}. \tag{8.4}$$

To bound the term II, we use the fact that $(\mathbf{b} \cdot \nabla \psi_{\mathbf{a}}) u_h \in Y_h^{\mathbf{a}}$ when $k' \geq k$ and that $\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} s_h^{\mathbf{a}}) \in Y_h^{\mathbf{a}}$, so that (4.3) actually holds pointwise, in the form

$$\mathbf{b} \cdot \nabla(\psi_{\mathbf{a}} s_h^{\mathbf{a}}) = \Pi_{\mathcal{P}^{k'}(\mathcal{T}_{\mathbf{a}})}(f \psi_{\mathbf{a}}) + (\mathbf{b} \cdot \nabla \psi_{\mathbf{a}}) u_h.$$

Hence, denoting $\bar{v}_K^{\mathbf{a}}$ the mean value of $v^{\mathbf{a}}$ over the element $K \in \mathcal{T}_{\mathbf{a}}$ and using (2.3a), we obtain

$$\begin{aligned}
\text{II} &= \frac{(\Pi_{\mathcal{P}^{k'}(\mathcal{T}_{\mathbf{a}})}(f \psi_{\mathbf{a}}) - f \psi_{\mathbf{a}}, v^{\mathbf{a}})_{\omega_{\mathbf{a}}}}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} = \frac{\sum_{K \in \mathcal{T}_{\mathbf{a}}} (\Pi_{\mathcal{P}^{k'}(\mathcal{T}_{\mathbf{a}})}(f \psi_{\mathbf{a}}) - f \psi_{\mathbf{a}}, v^{\mathbf{a}} - \bar{v}_K^{\mathbf{a}})_K}{\|\mathbf{b} \cdot \nabla v^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}} \\
&\leq \frac{\max_{K \in \mathcal{T}_{\mathbf{a}}} h_K}{\pi |\mathbf{b}|} \|\Pi_{\mathcal{P}^{k'}(\mathcal{T}_{\mathbf{a}})}(f \psi_{\mathbf{a}}) - f \psi_{\mathbf{a}}\|_{\omega_{\mathbf{a}}}.
\end{aligned}$$

The assertion follows by combining the bounds on I and II with (8.1).

8.2. Proof of Corollary 4.5 (global efficiency and maximal overestimation)

Proceeding as in the proof of Theorem 4.4, one has

$$\begin{aligned}
\|u_h - s_h\|^2 &= \sum_{K \in \mathcal{T}_h} \left\| \sum_{\mathbf{a} \in \mathcal{V}_K} \psi_{\mathbf{a}}(u_h - s_h^{\mathbf{a}}) \right\|_K^2 \leq 2 \sum_{K \in \mathcal{T}_h} \sum_{\mathbf{a} \in \mathcal{V}_K} \|\psi_{\mathbf{a}}(u_h - s_h^{\mathbf{a}})\|_K^2 \\
&= 2 \sum_{\mathbf{a} \in \mathcal{V}_h} \|\psi_{\mathbf{a}}(u_h - s_h^{\mathbf{a}})\|_{\omega_{\mathbf{a}}}^2 \stackrel{(8.3), (8.4)}{\leq} 2C_{\text{cont,PF}}^2 \sum_{\mathbf{a} \in \mathcal{V}_h} \|u - u_h\|_{\omega_{\mathbf{a}}}^2.
\end{aligned}$$

Another estimate for the overlapping of the patches yields $\sum_{\mathbf{a} \in \mathcal{V}_h} \|u - u_h\|_{\omega_{\mathbf{a}}}^2 \leq 2\|u - u_h\|^2$ and leads to the assertion.

TABLE 1. Effectivity indices I_{eff} for u_h obtained by the PG2 method (3.2) and dG method (3.3); $k = k' = 1$.

# Elements	# DOF	PG2	dG
4	8	1.234	1.126
16	32	1.058	1.032
64	128	1.014	1.008
256	512	1.004	1.002

9. NUMERICAL EXPERIMENTS

We provide in this section a numerical illustration of our results in one space dimension. In the first set of examples in Section 9.1, we consider a polynomial right-hand side function f and study the efficiency of the estimator. Then, in Section 9.2, we consider a more general case to investigate the effect of the increase of the polynomial degree on the quality of the estimators. Henceforth, we consider $\Omega = (0, 1)$ with the mesh $\mathcal{T}_h = \{K_i\}_1^n$ with $K_i = [x_{i-1}, x_i]$ of uniform size. In the experiments, the numerical solution $u_h \in \mathcal{P}^k(\mathcal{T}_h)$ will be computed by two methods:

- the PG2 method (3.2) with the finite-dimensional spaces as in Example 3.2, $k \geq 0$,
- the dG method (3.3) with the finite-dimensional spaces as in Example 3.3, $k \geq 1$.

The *effectivity index* is defined as $I_{\text{eff}} := \eta/\|u - u_h\|$, i.e., as the ratio of the estimated and the actual error from Theorem 4.3.

We always take $\mathbf{b} = 1$; recall that solely scaling \mathbf{b} in (1.1) by a factor implies the same scaling of the exact solution u , the numerical approximation u_h , the error $\|u - u_h\|$, the reconstruction s_h , and of the estimators in Theorem 4.3 by the inverse of this factor. Thus, in particular, the effectivity index is independent of \mathbf{b} .

9.1. Efficiency of the estimator

Here we consider the advection problem (1.1) with the piecewise quadratic right-hand side defined as

$$f(x) = x^2 + x + \sin(2\pi x_{i-1}), \quad \text{on } K_i, 1 \leq i \leq n,$$

whose exact solution can be easily computed by integration of the right-hand side. The numerical solutions u_h are obtained by both PG2 and dG methods with $k = 1, 2$.

If one sets $k' = 2$ in Definition 4.2, the oscillation estimators $\eta_{\text{osc},K}$ from (4.5) disappear, since $f \in \mathcal{P}^2(\mathcal{T}_h)$. In this case, actually, one has $s_h \in \mathcal{P}^3(\mathcal{T}_h) \cap H^1_-(\Omega)$, see (4.7). Moreover, owing to (4.8), $\mathbf{b} \cdot \nabla s_h = f$ pointwise. Hence, s_h in this setting coincides with the exact solution u , $\eta = \|u - u_h\|$, and $I_{\text{eff}} = 1$ (up to the machine precision).

To assess the behavior in the case where the reconstruction s_h does not coincide with the exact solution, we also test the choice $k' = 1$ in Definition 4.2 together with $k = 1$. The effectivity indices, for different mesh sizes, and for PG2 and dG methods have been reported in Table 1.

Moreover, we numerically observe asymptotic exactness with mesh refinement, for both tested schemes.

9.2. Robustness with respect to the polynomial degree

We now consider the advection problem (1.1) with a non-polynomial right-hand side $f(x) = \tan^{-1}(x)$, for different polynomial degrees $0 \leq k \leq 4$. The results are presented in Table 2 for the PG2 method and in Table 3 for the dG method. We always set $k' = k$. We use the notation $\eta_{\text{NC}} := (\sum_{K \in \mathcal{T}_h} \eta_{\text{NC},K}^2)^{1/2}$ and $\eta_{\text{osc}} := (\sum_{K \in \mathcal{T}_h} \eta_{\text{osc},K}^2)^{1/2}$. The mesh is refined uniformly until the error estimator $\eta \leq 10^{-14}$; we encountered some irregularities in I_{eff} beyond this point due to machine precision. We observe optimal convergence order of the estimators and the independence of I_{eff} from the polynomial degree, in accordance with the theory.

TABLE 2. Convergence of the error $\|u - u_h\|$, the error estimators η , η_{NC} , and η_{Osc} , and the effectivity indices I_{eff} for the PG2 method (3.2) with different polynomial degrees k .

$k = 0, k' = 0$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	4	3.562e-02	3.951e-02	3.574e-02	4.601e-03	1.11
16	16	8.934e-03	9.161e-03	8.936e-03	2.877e-04	1.03
64	64	2.234e-03	2.248e-03	2.234e-03	1.798e-05	1.01
256	256	5.585e-04	5.593e-05	5.585e-04	1.124e-06	1.00
1024	1024	1.396e-04	1.397e-05	1.396e-04	7.025e-08	1.00
$k = 1, k' = 1$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	8	1.868e-03	1.955e-03	1.867e-03	9.783e-05	1.05
16	32	1.167e-04	1.181e-04	1.167e-04	1.531e-06	1.02
64	128	7.294e-06	7.315e-06	7.294e-06	2.393e-08	1.00
256	512	4.559e-07	4.562e-07	4.559e-07	3.739e-10	1.00
1024	2048	2.849e-08	2.849e-08	2.849e-08	5.843e-12	1.00
$k = 2, k' = 2$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	12	2.600e-05	2.844e-05	2.598e-05	3.967e-06	1.09
16	48	4.066e-07	4.154e-07	4.066e-07	1.558e-08	1.02
64	192	6.354e-09	6.387e-09	6.354e-09	6.091e-11	1.01
256	768	9.928e-11	9.941e-11	9.928e-11	2.379e-13	1.00
1024	3072	1.552e-12	1.551e-12	1.551e-12	9.294e-16	1.00
$k = 3, k' = 3$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	16	7.859e-07	9.299e-07	7.852e-07	1.803e-07	1.18
16	64	3.085e-09	3.213e-09	3.085e-09	1.775e-10	1.04
64	256	1.205e-11	1.217e-11	1.205e-11	1.735e-13	1.01
256	1024	4.730e-14	4.730e-14	4.718e-14	1.694e-16	1.00
$k = 4, k' = 4$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	20	2.851e-08	3.517e-08	2.847e-08	8.486e-09	1.23
16	80	2.804e-11	2.948e-11	2.804e-11	2.095e-12	1.05
64	320	2.753e-14	2.776e-14	2.742e-14	5.118e-16	1.01

10. EXTENSION TO MULTIPLE SPACE DIMENSIONS

In this section, we investigate a possible extension of the ideas presented so far to the multi-dimensional case. We consider the advection equation (1.1) on a simply-connected Lipschitz polytope $\Omega \subset \mathbb{R}^d$ for $d \geq 2$. The velocity field $\mathbf{b} \in C^1(\bar{\Omega}; \mathbb{R}^d)$ is considered to be divergence-free. We also assume that \mathbf{b} is Ω -filling, i.e., its trajectories starting from the inflow boundary $\partial_- \Omega$ fill $\bar{\Omega}$ almost everywhere in a finite time. A sufficient condition for the validity of this property is given by [3] (see Lem. 10.1 below). One can find necessary and sufficient conditions in Lemma 2.3 of [13], see also [2, 8, 9, 12, 21].

10.1. Spaces

We start by introducing proper generalizations of (2.1). Let us define the operator related to (1.1) and its formal adjoint as

$$\mathcal{L}: v \mapsto \mathbf{b} \cdot \nabla v, \quad \mathcal{L}^*: v \mapsto -\nabla \cdot (\mathbf{b}v) = -\mathbf{b} \cdot \nabla v,$$

TABLE 3. Convergence of the error $\|u - u_h\|$, the error estimators η , η_{NC} , and η_{Osc} , and the effectivity indices I_{eff} for the dG method (3.3) with different polynomial degrees k .

$k = 1, k' = 1$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	8	3.021e-03	3.136e-03	3.048e-03	9.783e-05	1.04
16	32	1.901e-04	1.919e-03	1.906e-04	1.531e-06	1.01
64	128	1.190e-05	1.193e-05	1.191e-05	2.393e-08	1.00
256	512	7.444e-07	7.447e-07	7.445e-07	3.739e-10	1.00
1024	2048	4.653e-08	4.653e-08	4.653e-08	5.843e-12	1.00
$k = 2, k' = 2$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	12	4.045e-05	4.260e-05	4.210e-05	3.967e-06	1.05
16	48	6.307e-07	6.386e-07	6.299e-07	1.558e-08	1.01
64	192	9.847e-09	9.877e-09	9.844e-09	6.091e-11	1.00
256	768	1.538e-10	1.539e-10	1.538e-10	2.379e-13	1.00
1024	3072	2.403e-12	2.403e-12	2.403e-12	9.294e-16	1.00
$k = 3, k' = 3$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	16	1.169e-06	1.328e-06	1.186e-06	1.803e-07	1.14
16	64	4.647e-09	4.791e-09	4.664e-09	1.775e-10	1.03
64	256	1.821e-11	1.834e-11	1.822e-11	1.735e-13	1.01
256	1024	7.181e-14	7.184e-14	7.172e-14	1.694e-16	1.00
$k = 4, k' = 4$						
# Elements	# DOF	$\ u - u_h\ $	η	η_{NC}	η_{Osc}	I_{eff}
4	20	4.252e-08	4.895e-08	4.240e-08	8.486e-09	1.15
16	80	4.180e-11	4.323e-11	4.179e-11	2.095e-12	1.03
64	320	4.094e-14	4.117e-14	4.083e-14	5.118e-16	1.01

together with the following graph spaces

$$H(\mathcal{L}, \Omega) := \{v \in L^2(\Omega), \mathcal{L}v \in L^2(\Omega)\}, \quad H(\mathcal{L}^*, \Omega) := \{v \in L^2(\Omega), \mathcal{L}^*v \in L^2(\Omega)\}.$$

Then $\mathcal{L} : H(\mathcal{L}, \Omega) \rightarrow L^2(\Omega)$ and $\mathcal{L}^* : H(\mathcal{L}^*, \Omega) \rightarrow L^2(\Omega)$, and $H(\mathcal{L}, \Omega) = H(\mathcal{L}^*, \Omega)$. Moreover, one can define the following subspaces of the graph spaces with incorporated boundary conditions:

$$H_0(\mathcal{L}, \Omega) := \{v \in H(\mathcal{L}, \Omega), v = 0 \text{ on } \partial_- \Omega\},$$

$$H_0(\mathcal{L}^*, \Omega) := \{v \in H(\mathcal{L}^*, \Omega), v = 0 \text{ on } \partial_+ \Omega\}.$$

These definitions are consistent extensions from $d = 1$ in that the spaces $H(\mathcal{L}, \Omega), H(\mathcal{L}^*, \Omega), H_0(\mathcal{L}, \Omega), H_0(\mathcal{L}^*, \Omega)$ become respectively $H^1(\Omega), H^1(\Omega), H^1_-(\Omega)$, and $H^1_+(\Omega)$. One might confer with [31], page 131 and Theorems 2.1 and 2.2 of [23] for the justification of the trace operator which is discussed as an operator from $H(\mathcal{L}, \Omega)$ to $H^{-\frac{1}{2}}(\partial_- \Omega)$ (or from $H(\mathcal{L}^*, \Omega)$ to $H^{-\frac{1}{2}}(\partial_+ \Omega)$, respectively). The extension to $L^2(|\mathbf{b} \cdot \mathbf{n}|; \partial_- \Omega)$ is possible under slightly more restrictive conditions, see [31], page 133, Lemma 3.1 of [14], and more recently Proposition 2.3 of [12]. Moreover the following integration-by-parts formula holds true:

$$(v, \mathbf{b} \cdot \nabla w) + (\mathbf{b} \cdot \nabla v, w) = (\mathbf{b} \cdot \mathbf{n} v, w) \quad \forall v \in H(\mathcal{L}, \Omega), \quad \forall w \in H^1(\Omega). \tag{10.1}$$

The result (10.1) can be extended to $w \in H(\mathcal{L}^*, \Omega)$ if traces are meaningful in $L^2(|\mathbf{b} \cdot \mathbf{n}|; \partial_- \Omega)$.

10.2. Streamline Poincaré inequality

The following sufficient condition for the field \mathbf{b} to be Ω -filling is given in [3]:

Lemma 10.1 (Ω -filling sufficient condition). *Let $\mathbf{b} \in \mathcal{C}^1(\overline{\Omega}; \mathbb{R}^d)$ and assume that there is a fixed unit vector $\mathbf{k} \in \mathbb{R}^d$ and a real number $\alpha > 0$ such that*

$$\forall x \in \overline{\Omega}, \quad \mathbf{b} \cdot \mathbf{k} \geq \alpha. \quad (10.2)$$

Then \mathbf{b} is Ω -filling.

For Ω -filling \mathbf{b} , we can extend the inequality (2.3b) along the flow of \mathbf{b} , cf. [3]:

Lemma 10.2 (Streamline Poincaré inequality). *Let the field $\mathbf{b} \in \mathcal{C}^1(\overline{\Omega}; \mathbb{R}^d)$ be divergence-free and Ω -filling. Then there exists a streamline Poincaré constant $C_{\mathbf{P}, \mathbf{b}, \Omega}$ such that*

$$\|v\| \leq C_{\mathbf{P}, \mathbf{b}, \Omega} \|\mathbf{b} \cdot \nabla v\| \quad \forall v \in H_0(\mathcal{L}, \Omega). \quad (10.3)$$

The constant $C_{\mathbf{P}, \mathbf{b}, \Omega}$ is bounded by $C_{\mathbf{P}, \mathbf{b}, \Omega} \leq 2T$, where T is the longest time that trajectories of the field \mathbf{b} spend in the domain Ω . In particular, $T \leq \text{diam}(\Omega)/\alpha$ under assumption (10.2).

A similar result can also be obtained for a non divergence-free field, see [1]. In the case where the field \mathbf{b} is constant, one can easily set \mathbf{k} as the direction of the flow and $\alpha = |\mathbf{b}|$. A crucial consequence of Lemma 10.2 is that one can equip the spaces $H_0(\mathcal{L}, \Omega)$ and $H_0(\mathcal{L}^*, \Omega)$ with the norm $\|\mathbf{b} \cdot \nabla(\cdot)\|$.

Remark 10.3 (Functions with zero mean value). While, following from Lemma 10.2, the streamline Poincaré inequality holds true for functions with zero trace on the inflow of an arbitrary domain Ω , such a result is not valid for functions with zero mean value as a variant of the Poincaré inequality (2.3a) in multiple spatial dimensions. This leads to significant differences in the analysis of the multi-dimensional case compared to the one-dimensional one, and less complete results that we are able to present here.

10.3. Error–residual equivalence

We consider the multi-dimensional extension of the ultra-weak formulation (2.5): find $u \in L^2(\Omega)$ such that

$$-(u, \mathbf{b} \cdot \nabla v) = (f, v) \quad \forall v \in H_0(\mathcal{L}^*, \Omega). \quad (10.4)$$

Define the residual operator $\mathcal{R}(u_h) \in H_0(\mathcal{L}^*, \Omega)'$ and its dual norm as in (2.6) and (2.7), upon replacing $H_+^1(\Omega)$ by $H_0(\mathcal{L}^*, \Omega)$. One can extend the equivalence of Proposition 2.1 to the multi-dimensional case as follows:

Proposition 10.4 (Error–residual equivalence). *Let the field $\mathbf{b} \in \mathcal{C}^1(\overline{\Omega}; \mathbb{R}^d)$ be divergence-free and Ω -filling. Let $u \in L^2(\Omega)$ be the ultra-weak solution of (10.4). Then*

$$\|u - u_h\| = \|\mathcal{R}(u_h)\|_{\mathbf{b}, H_0(\mathcal{L}^*, \Omega)'} \quad \forall u_h \in L^2(\Omega).$$

Proof. We use the fact that for all $v \in L^2(\Omega)$, there exists a unique $z \in H_0(\mathcal{L}^*, \Omega)$ such that

$$-(\mathbf{b} \cdot \nabla z, w) = (v, w) \quad \forall w \in L^2(\Omega).$$

The rest of the proof goes along the lines of that of Proposition 2.1. □

10.4. Local problems and the error indicator

In this section, we propose a heuristic approach inspired by the rigorous discussions in the one-dimensional case. First, let us consider the following reconstruction, mimicking Definition 4.2. Here \mathcal{T}_h is a simplicial mesh of Ω , \mathcal{T}_a the patch of all simplices which share the given vertex $\mathbf{a} \in \mathcal{V}_h$, ω_a the corresponding open subdomain, and ψ_a the associated hat basis function.

Definition 10.5 (Patchwise problems). Let $u_h \in L^2(\Omega)$. For all vertices $\mathbf{a} \in \mathcal{V}_h$, let $s_h^\mathbf{a} \in X_h^\mathbf{a}$ be the solution of the following least-squares problem on the patch subdomain ω_a :

$$s_h^\mathbf{a} := \arg \min_{v_h \in X_h^\mathbf{a}} \left\{ \|\psi_a(u_h - v_h)\|_{\omega_a}^2 + C_{\text{opt}}^2 \|f\psi_a + (\mathbf{b} \cdot \nabla \psi_a) u_h - \mathbf{b} \cdot \nabla (\psi_a v_h)\|_{\omega_a}^2 \right\}. \tag{10.5}$$

For $k' \geq 0$, we take the finite-dimensional space $X_h^\mathbf{a} := \mathcal{P}^{k'}(\mathcal{T}_a) \cap H_0(\mathcal{L}, \omega_a)$ when the vertex \mathbf{a} lies in the closure of the inflow boundary $\partial_- \Omega$ and $X_h^\mathbf{a} := \mathcal{P}^{k'}(\mathcal{T}_a) \cap H(\mathcal{L}, \omega_a)$ otherwise. Here $C_{\text{opt}} > 0$ is a constant to be chosen. The global reconstruction s_h is defined by

$$s_h := \sum_{\mathbf{a} \in \mathcal{V}_h} \psi_a s_h^\mathbf{a}, \tag{10.6}$$

leading to $s_h \in \mathcal{P}^{k'+1}(\mathcal{T}_h) \cap H_0(\mathcal{L}, \Omega)$.

Remark 10.6 (Continuity of s_h). One might note that the reconstruction s_h of Definition 10.5, lying in the space $H_0(\mathcal{L}, \Omega)$, possibly allows capturing the discontinuity that may appear in the exact solution u across the streamlines. This is, however, only in reach if the triangulation is aligned with the streamlines. If this is not the case, the reconstruction s_h actually lies in the smoother space $H^1(\Omega)$.

In order to see the rationale behind the above reconstruction, one might note the following upper bound on the error exploiting Proposition 10.4, the integration-by-parts formula (10.1), the Cauchy–Schwarz inequality, and the streamline Poincaré inequality (10.3): for any $s_h \in H_0(\mathcal{L}, \Omega)$, we have

$$\begin{aligned} \|u - u_h\| &= \|\mathcal{R}(u_h)\|_{\mathbf{b}; H_0(\mathcal{L}^*, \Omega)'} = \sup_{v \in H_0(\mathcal{L}^*, \Omega) \setminus \{0\}} \frac{(f - \mathbf{b} \cdot \nabla s_h, v) + (u_h - s_h, \mathbf{b} \cdot \nabla v)}{\|\mathbf{b} \cdot \nabla v\|} \\ &\leq \|u_h - s_h\| + C_{\mathbf{P}, \mathbf{b}, \Omega} \|f - \mathbf{b} \cdot \nabla s_h\|. \end{aligned}$$

Furthermore, using that almost each point in Ω belongs to $(d + 1)$ patch subdomains ω_a and the partition of unity (4.1), the construction of s_h via (10.6) gives the following upper bound:

$$\|u - u_h\| \leq \left\{ 2(d + 1) \sum_{\mathbf{a} \in \mathcal{V}_h} \left[\|\psi_a(u_h - s_h^\mathbf{a})\|_{\omega_a}^2 + C_{\mathbf{P}, \mathbf{b}, \Omega}^2 \|f\psi_a + (\mathbf{b} \cdot \nabla \psi_a) u_h - \mathbf{b} \cdot \nabla (\psi_a s_h^\mathbf{a})\|_{\omega_a}^2 \right] \right\}^{1/2}.$$

In particular, the idea of adding $0 = \sum_{\mathbf{a} \in \mathcal{V}_h} \mathbf{b} \cdot \nabla \psi_a u_h$ is inspired by the analysis in the one-dimensional case. By comparison to Definition 10.5 one can see that the least-squares problems (10.5) minimize contributions to the upper bound on the error, and a theoretically-motivated choice for C_{opt} would be $C_{\text{opt}} = C_{\mathbf{P}, \mathbf{b}, \Omega}$. This in particular leads to the guaranteed estimate as the following theorem:

Theorem 10.7 (Guaranteed *a posteriori* error estimate). *Let $u \in L^2(\Omega)$ be the ultra-weak solution of (10.4) and let $u_h \in L^2(\Omega)$ be arbitrary. Furthermore, consider s_h to be the reconstruction from Definition 10.5 with $k' \geq 0$ and arbitrary C_{opt} . Then*

$$\|u - u_h\| \leq \eta := \left\{ \sum_{K \in \mathcal{T}_h} \eta_{\text{NC}, K}^2 \right\}^{1/2} + \left\{ \sum_{K \in \mathcal{T}_h} \eta_{\text{R}, K}^2 \right\}^{1/2}, \tag{10.7}$$

where

$$\eta_{\text{NC}, K} := \|u_h - s_h\|_K, \quad \eta_{\text{R}, K} := C_{\mathbf{P}, \mathbf{b}, \Omega} \|f - \mathbf{b} \cdot \nabla s_h\|_K.$$

Numerical experiments show that the estimate (10.7) is rather sharp when the solution u is discontinuous and the discontinuity line of u is not aligned with the triangulation. When, however, (i) the solution u is smooth; (ii) u is discontinuous and the discontinuity line of u is aligned with the triangulation, the estimators $\eta_{R,K}$ do not converge with the right order so that the corresponding effectivity indices increase with mesh refinement. This apparently comes from the special structure of the minimization term which cannot be approximated up to the projection error (see Rem. 10.6), in contrast to the one-dimensional case, where (4.8) holds true. Congruently, the lack of the Poincaré inequality in the streamline form (see Rem. 10.3) implies the loss of the scaling by the mesh element diameters h_K in the second term in (10.7), compare with η_{Osc} given by (4.5) in one space dimension.

The following remark provides a heuristic rectification for this under assumption (10.2):

Remark 10.8 (Heuristic modification). In both cases (i) or (ii) mentioned above, we heuristically replace $C_{P,b,\Omega}$ (which typically scales as $2 \text{diam}(\Omega)/\alpha$, see Lem. 10.2) in the estimator $\eta_{R,K}$ of (10.7) by local terms $C'h_K/\alpha$, where C' is a user-dependent constant and h_K the diameter of the mesh element K . Then the modification of the guaranteed estimator η from (10.7) is the non-guaranteed error indicator

$$\eta^{\text{mod}} := \left\{ \sum_{K \in \mathcal{T}_h} (\eta_{\text{NC},K}^2 + (\eta_{R,K}^{\text{mod}})^2) \right\}^{1/2}, \quad (10.8a)$$

$$\eta_{\text{NC},K} := \|u_h - s_h\|_K, \quad \eta_{R,K}^{\text{mod}} := \frac{C'h_K}{\alpha} \|f - \mathbf{b} \cdot \nabla s_h\|_K, \quad (10.8b)$$

where s_h is the reconstruction from Definition 10.5 with $k' \geq 0$ and arbitrary C_{opt} .

Overall, one has in cases (i) or (ii) two free parameters to choose, C_{opt} for the local problems in (10.5) and C' in (10.8b). We set below $C_{\text{opt}} = 2 \text{diam}(\Omega)/\alpha$, as suggested by Lemma 10.2, and $C' = 2$. Numerically, our results are actually not sensitive to the choice of the parameter C_{opt} .

Remark 10.9 (Non-homogeneous boundary condition). The treatment of the non-homogeneous boundary condition $u = g$ on $\partial_- \Omega$ for $g \in H^{\frac{1}{2}}(\partial \Omega)$ is similar to Remark 4.9 in the one-dimensional setting. In the reconstruction (10.5), one, however, needs to impose the boundary condition in the definition of the space X_h^a strongly by a piecewise polynomial projection of the datum g . Then an additional technicality comes from the difference between this projection and g , which then appears as a second data oscillation term in the estimator.

10.5. Numerical experiments

In this section, we provide some numerical tests in two space dimensions. We consider $\Omega = (0, 1)^2$ and uniformly refined structured triangulations aligned with the slope 45° . We only test here the dG method (3.3), since it is the only method among those considered in Section 3 which is well-defined in multiple space dimensions. The implementation is done in the framework of FreeFEM++ [22] and based on the scripts for the reconstruction-based *a posteriori* estimation by [32].

Below, we will consider three different test cases: In Section 10.5.1, we show an example where both the exact solution u and the reconstruction s_h are in $H^1(\Omega)$. Section 10.5.2 deals with a case with discontinuous solution $H_0(\mathcal{L}, \Omega) \setminus H^1(\Omega)$ aligned with the triangulation and discontinuous reconstruction $s_h \in H_0(\mathcal{L}, \Omega) \setminus H^1(\Omega)$. Finally, Section 10.5.3 discusses the case of a discontinuous solution $u \in H_0(\mathcal{L}, \Omega) \setminus H^1(\Omega)$ not aligned with the triangulation, where the reconstruction becomes continuous, $s_h \in H^1(\Omega)$. In Sections 10.5.1 and 10.5.2 the heuristic indicator η^{mod} of (10.8) is used, whereas in Section 10.5.3, we will show that the guaranteed error indicator η of (10.7) actually works well. In the former case, $\eta_{\text{R}}^{\text{mod}} := (\sum_{K \in \mathcal{T}_h} (\eta_{R,K}^{\text{mod}})^2)^{1/2}$ from (10.8b) with $C' = 2$ is employed, whereas in the latter case, $\eta_{\text{R}} := \{\sum_{K \in \mathcal{T}_h} \eta_{R,K}^2\}^{1/2}$; we always set $\eta_{\text{NC}} := (\sum_{K \in \mathcal{T}_h} \eta_{\text{NC},K}^2)^{1/2}$ and rely on Definition 10.5 with the choice $k' = k+1$ and typically $C_{\text{opt}} = 2 \text{diam}(\Omega)/\alpha$. We set $I_{\text{eff}}^{\text{mod}} := \eta^{\text{mod}}/\|u - u_h\|$ and $I_{\text{eff}} := \eta/\|u - u_h\|$, where only $I_{\text{eff}} \geq 1$ is guaranteed.

TABLE 4. Smooth solution (10.9); error $\|u - u_h\|$, error estimators η^{mod} , η_{NC} , and $\eta_{\text{R}}^{\text{mod}}$, and effectivity indices $I_{\text{eff}}^{\text{mod}}$ and I_{eff} for the dG method (3.3); $\mathbf{b} = (1, 1)^t$ and different polynomial degrees k .

$k = 1, k' = 2$							
# Elements	# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$	I_{eff}
8	24	1.097e-01	2.284e-01	9.365e-02	2.083e-01	2.08	2.67
32	96	2.963e-02	4.894e-02	2.584e-02	4.156e-02	1.65	4.03
128	384	7.553e-03	1.101e-02	6.786e-03	8.666e-03	1.45	6.54
512	1536	1.897e-03	2.630e-03	1.727e-03	1.983e-03	1.38	11.8
2048	6144	4.749e-04	6.456e-04	4.347e-04	4.773e-04	1.35	22.7
8192	24576	1.187e-04	1.601e-04	1.088e-04	1.173e-04	1.34	44.7
$k = 2, k' = 3$							
# Elements	# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$	I_{eff}
8	48	1.882e-02	5.317e-02	2.271e-02	4.807e-02	2.82	3.81
32	192	2.476e-03	4.896e-03	3.106e-03	3.785e-03	1.97	4.50
128	768	3.135e-04	5.742e-04	3.972e-04	4.147e-04	1.83	7.58
512	3072	3.929e-05	7.076e-05	4.995e-05	5.012e-05	1.80	14.4
2048	12288	4.934e-06	8.817e-06	6.253e-06	6.216e-06	1.78	28.5
8192	49152	6.270e-07	1.107e-06	7.822e-07	7.843e-07	1.76	56.6

TABLE 5. Smooth solution (10.9); error $\|u - u_h\|$, error estimators η^{mod} , η_{NC} , and $\eta_{\text{R}}^{\text{mod}}$, and effectivity indices $I_{\text{eff}}^{\text{mod}}$ for the dG method (3.3); different velocity fields \mathbf{b} and $k = 1$.

$k = 1, k' = 2, \mathbf{b} = (100, 100)^t$						
# Elements	# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$
8	24	1.097e-01	2.284e-01	9.365e-02	2.083e-01	2.08
32	96	2.963e-02	4.894e-02	2.584e-02	4.156e-02	1.65
128	384	7.553e-03	1.101e-02	6.786e-03	8.666e-03	1.45
512	1536	1.897e-03	2.630e-03	1.727e-03	1.983e-03	1.38
2048	6144	4.749e-04	6.456e-04	4.347e-04	4.773e-04	1.35
8192	24576	1.187e-04	1.601e-04	1.088e-04	1.173e-04	1.34
$k = 1, k' = 2, \mathbf{b} = (10, 1)^t$						
# Elements	# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$
8	24	1.009e-01	2.361e-01	8.299e-02	2.216e-01	2.33
32	96	2.896e-02	5.140e-02	2.057e-02	4.714e-02	1.77
128	384	7.965e-03	1.188e-02	5.325e-03	1.062e-02	1.49
512	1536	2.069e-03	3.014e-03	1.370e-03	2.684e-03	1.45
2048	6144	5.241e-04	7.636e-04	3.459e-04	6.807e-04	1.45
8192	24576	1.316e-04	1.918e-04	8.667e-05	1.711e-04	1.45
$k = 1, k' = 2, \mathbf{b} = (y, x + 1)^t (\alpha = 1)$						
# Elements	# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$
8	24	1.134e-01	2.435e-01	9.582e-02	2.239e-01	2.14
32	96	3.152e-02	5.787e-02	2.513e-02	5.212e-02	1.83
128	384	8.007e-03	1.393e-02	6.478e-03	1.233e-02	1.74
512	1536	2.013e-03	3.409e-03	1.636e-03	2.991e-03	1.69
2048	6144	5.053e-04	8.443e-04	4.103e-04	7.379e-04	1.67
8192	24576	1.267e-04	2.101e-04	1.027e-04	1.833e-04	1.65

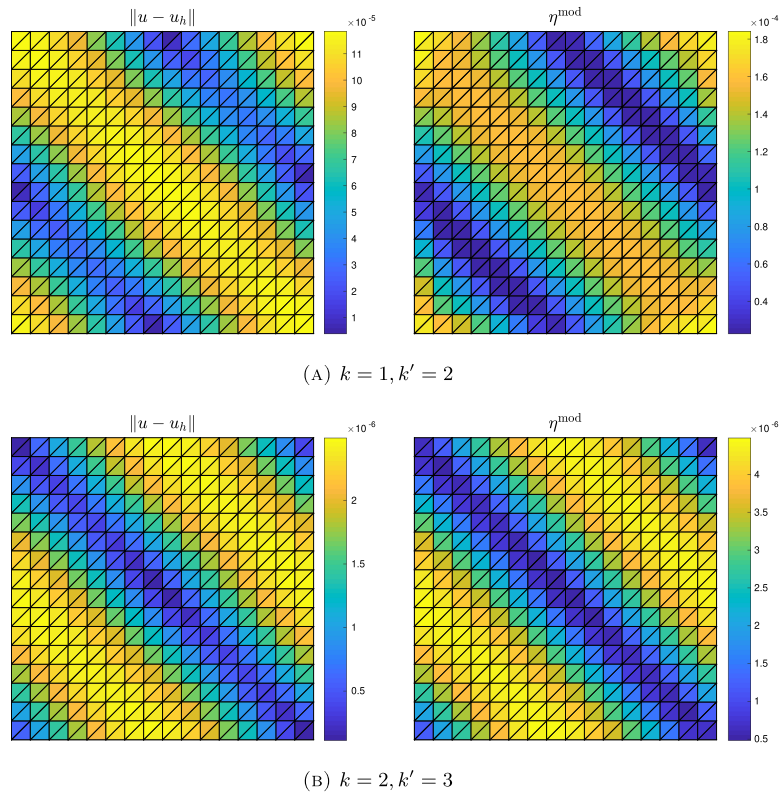


FIGURE 1. Smooth solution (10.9); distribution of the errors $\|u - u_h\|_K$ (left) and of the local error estimators η_K^{mod} (right) for the dG method (3.3) with 512 elements; $\mathbf{b} = (1, 1)^t$ and different polynomial degrees k .

TABLE 6. Discontinuous solution (10.10) with aligned triangulation; error $\|u - u_h\|$, error estimators η^{mod} , η_{NC} , and $\eta_{\text{R}}^{\text{mod}}$, and effectivity indices $I_{\text{eff}}^{\text{mod}}$ and I_{eff} for the dG method (3.3); $\mathbf{b} = (1, 1)^t$ and different polynomial degrees k .

$k = 1, k' = 2$						
# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$	I_{eff}
24	7.75e-02	1.61e-01	6.62e-02	1.47e-01	2.08	2.67
96	2.09e-02	3.46e-02	1.82e-02	2.94e-02	1.65	4.04
384	5.34e-03	7.78e-03	4.79e-03	6.12e-03	1.46	6.55
1536	1.34e-03	1.86e-03	1.22e-03	1.40e-03	1.38	11.8
6144	3.35e-04	4.56e-04	3.07e-04	3.37e-04	1.36	22.7
24576	8.39e-05	1.13e-04	7.70e-05	8.29e-05	1.35	44.7
$k = 2, k' = 3$						
# DOF	$\ u - u_h\ $	η^{mod}	η_{NC}	$\eta_{\text{R}}^{\text{mod}}$	$I_{\text{eff}}^{\text{mod}}$	I_{eff}
48	1.33e-02	3.75e-02	1.61e-02	3.39e-02	2.82	3.81
192	1.75e-03	3.46e-03	2.19e-03	2.67e-03	1.97	4.50
768	2.21e-04	4.06e-04	2.81e-04	2.93e-04	1.83	7.58
3072	2.77e-05	5.00e-05	3.53e-05	3.54e-05	1.80	14.4
12288	3.48e-06	6.23e-06	4.42e-06	4.39e-06	1.78	28.5
49152	4.43e-07	7.83e-07	5.53e-07	5.54e-07	1.76	56.6

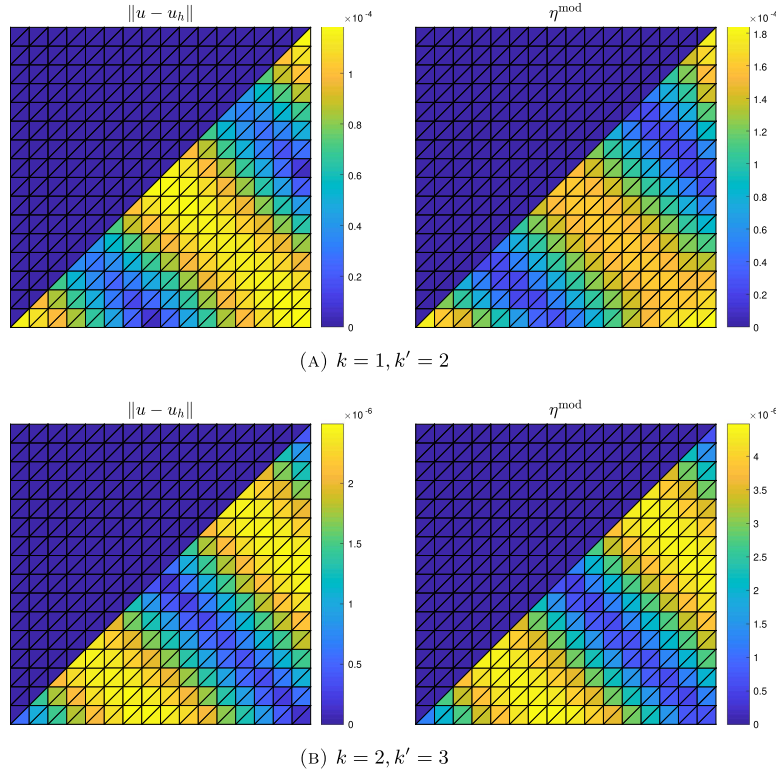


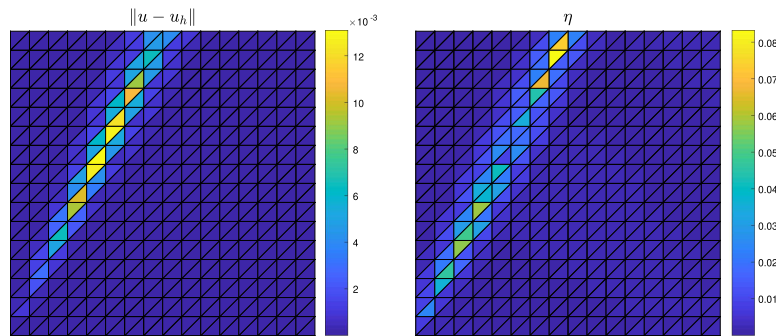
FIGURE 2. Discontinuous solution (10.10) with aligned triangulation; distribution of the errors $\|u - u_h\|_K$ (left) and of the local error estimators η_K^{mod} (right) for the dG method (3.3) with 512 elements; $\mathbf{b} = (1, 1)^t$ and different polynomial degrees k .

TABLE 7. Discontinuous solution (10.11a) with non-aligned triangulation; error $\|u - u_h\|$, error estimators η , η_{NC} , and η_{R} with convergence rates, and effectivity indices I_{eff} for the dG method (3.3); different polynomial degrees k .

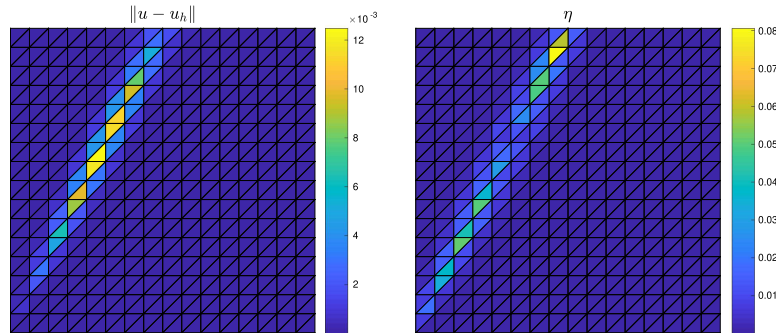
$k = 1, k' = 2$						
# DOF	$\ u - u_h\ $	η	η_{NC}	η_{R}	I_{eff}	
24	1.41e-01	5.70e-01	7.60e-02	5.65e-01	4.03	
96	8.36e-02 (0.76)	4.02e-01 (0.50)	3.11e-02 (1.29)	4.01e-01 (0.50)	4.80	
384	5.34e-02 (0.65)	2.89e-01 (0.48)	1.17e-02 (1.41)	2.89e-01 (0.47)	5.42	
1536	4.08e-02 (0.39)	2.31e-01 (0.32)	5.51e-03 (1.09)	2.31e-01 (0.32)	5.67	
6144	3.16e-02 (0.37)	1.93e-01 (0.26)	2.93e-03 (0.91)	1.94e-01 (0.26)	6.13	
24576	2.45e-02 (0.37)	1.70e-01 (0.18)	1.62e-03 (0.86)	1.71e-01 (0.18)	6.97	
$k = 2, k' = 3$						
# DOF	$\ u - u_h\ $	η	η_{NC}	η_{R}	I_{eff}	
48	1.28e-01	4.17e-01	4.31e-02	4.15e-01	3.24	
192	7.08e-02 (0.85)	2.82e-01 (0.56)	1.12e-02 (1.94)	2.82e-01 (0.54)	3.99	
768	4.75e-02 (0.58)	2.29e-01 (0.30)	5.59e-03 (1.00)	2.29e-01 (0.30)	4.83	
3072	3.50e-02 (0.44)	1.84e-01 (0.32)	2.83e-03 (0.98)	1.84e-01 (0.31)	5.26	
12288	2.54e-02 (0.46)	1.45e-01 (0.33)	1.50e-03 (0.92)	1.45e-01 (0.33)	5.73	
49152	1.85e-02 (0.46)	1.20e-01 (0.28)	8.41e-04 (0.83)	1.20e-01 (0.27)	6.47	

TABLE 8. Discontinuous solution (10.11b) with non-aligned triangulation; error $\|u - u_h\|$, error estimators η , η_{NC} , and η_R with convergence rates, and effectivity indices I_{eff} for the dG method (3.3); different polynomial degrees k .

$k = 1, k' = 2$					
# DOF	$\ u - u_h\ $	η	η_{NC}	η_R	I_{eff}
24	1.70e-01	6.14e-01	7.30e-02	6.09e-01	3.60
96	9.31e-02 (0.87)	4.42e-01 (0.47)	2.99e-02 (1.29)	4.41e-01 (0.47)	4.75
384	6.01e-02 (0.63)	3.24e-01 (0.45)	1.16e-02 (1.37)	3.24e-01 (0.44)	5.39
1536	4.62e-02 (0.38)	2.67e-01 (0.28)	5.31e-03 (1.13)	2.68e-01 (0.27)	5.79
6144	3.57e-02 (0.37)	2.36e-01 (0.18)	2.79e-03 (0.93)	2.37e-01 (0.18)	6.61
24576	2.78e-02 (0.36)	2.29e-01 (0.04)	1.54e-03 (0.86)	2.29e-01 (0.05)	8.26
$k = 2, k' = 3$					
# DOF	$\ u - u_h\ $	η	η_{NC}	η_R	I_{eff}
48	9.83e-02	4.31e-01	3.72e-02	4.29e-01	4.38
192	5.72e-02 (0.78)	2.85e-01 (0.59)	1.06e-02 (1.81)	2.85e-01 (0.59)	4.98
768	4.64e-02 (0.30)	2.34e-01 (0.29)	5.14e-03 (1.04)	2.34e-01 (0.28)	5.03
3072	3.31e-02 (0.48)	1.90e-01 (0.29)	2.78e-03 (0.89)	1.90e-01 (0.30)	5.75
12288	2.59e-02 (0.35)	1.72e-01 (0.14)	1.55e-03 (0.84)	1.72e-01 (0.14)	6.63
49152	1.92e-02 (0.43)	1.58e-01 (0.12)	8.44e-04 (0.88)	1.58e-01 (0.12)	8.27



(A) $k = 1, k' = 2$



(B) $k = 2, k' = 3$

FIGURE 3. Discontinuous solution (10.11a) with non-aligned triangulation; distribution of the errors $\|u - u_h\|_K$ (left) and of the local error estimators η_K (right) for the dG method (3.3) with 512 elements; different polynomial degrees k .

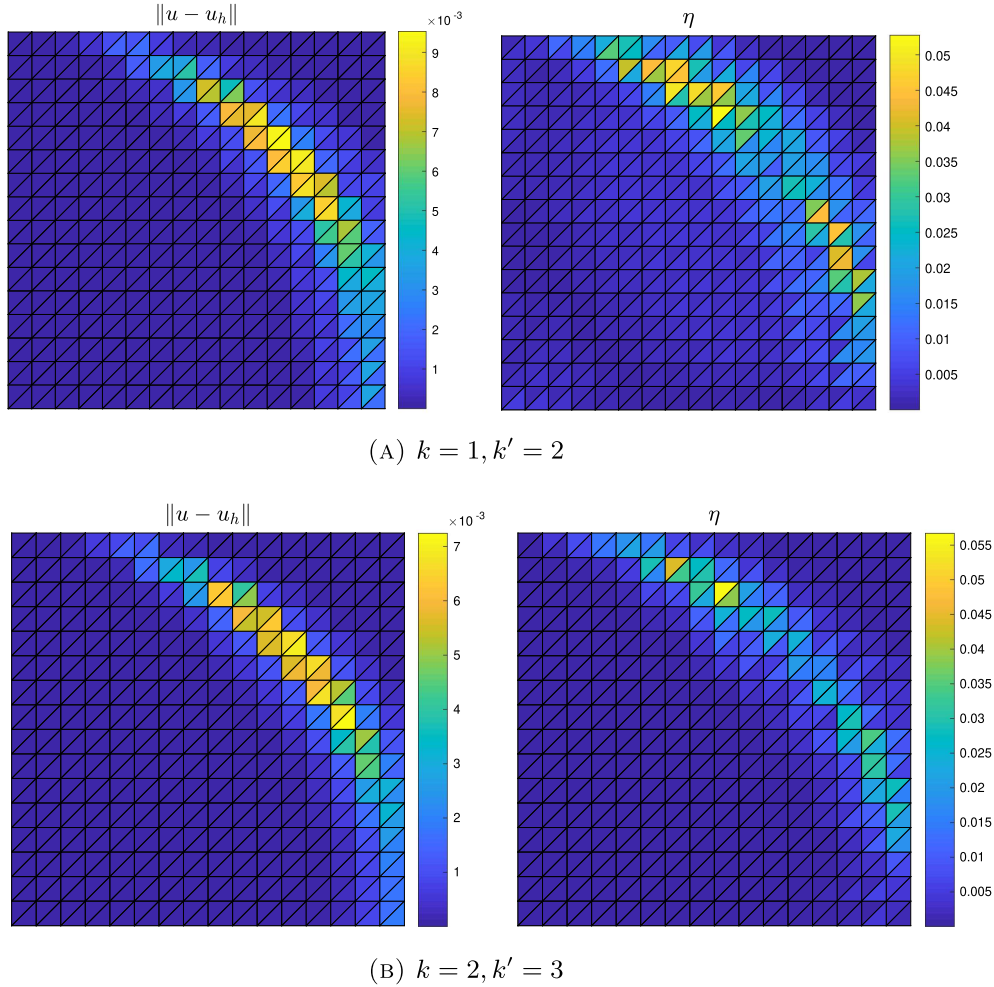


FIGURE 4. Discontinuous solution (10.11b) with non-aligned triangulation; distribution of the errors $\|u - u_h\|_K$ (left) and of the local error estimators η_K (right) for the dG method (3.3) with 512 elements; different polynomial degrees k .

10.5.1. Smooth solution

We apply the right-hand side f such that the solution of (1.1) is

$$u(x, y) = \sin(\pi x) \sin(\pi y), \tag{10.9}$$

for different velocity fields \mathbf{b} . The results are presented in Tables 4 and 5 for various polynomial degrees k . We use $\alpha = |\mathbf{b}|$ except for the non-constant velocity where $\alpha = 1$ is taken. The error indicator η^{mod} of (10.8) performs well in actually providing an upper bound on the error and simultaneously not overestimating it excessively. Moreover, the efficiency results numerically appear to be robust with respect to both the velocity field \mathbf{b} and the polynomial degree k . As in Section 9.1, both u and u_h , but actually also s_h constructed following Definition 10.5, turn out to be insensitive to the scaling of \mathbf{b} by a constant, so that the estimators in (10.8) do not change either. In Figure 1, the distributions of the errors $\|u - u_h\|_K$ and of the error estimators $\eta_K^{\text{mod}} := (\eta_{\text{NC},K}^2 + (\eta_{\text{R},K}^{\text{mod}})^2)^{1/2}$

are presented. These distributions show a very close behavior, which suggests that the presented indicators should be suitable for adaptive mesh/polynomial degree refinement.

10.5.2. Discontinuous solution with aligned triangulation

In this example, we consider a discontinuous exact solution for (1.1). For the velocity field $\mathbf{b} = (1, 1)^t$ with $\alpha = |\mathbf{b}|$, we set

$$u(x, y) = \begin{cases} 0, & x < y, \\ \sin(\pi x) \sin(\pi y), & x > y, \end{cases} \quad (10.10)$$

and prescribe accordingly the right-hand side f . As the triangulation is set to be aligned with this discontinuity, the reconstruction s_h is continuous everywhere but not at the discontinuity line of the exact solution. The results are presented in Table 6 for different polynomial degrees k . They show robustness with respect to the polynomial degree of approximation. In Figure 2, the distributions of the error and of the error estimators η_K^{mod} are presented, showing again a very close behavior.

10.5.3. Discontinuous solution with non-aligned triangulation

In this section, we finally consider a discontinuous exact solution whose discontinuity is not aligned with the triangulation. We consider the following two examples:

$$u(x, y) = \begin{cases} 0, & 2x < y, \\ \sin(\pi x) \sin(\pi y), & 2x > y, \end{cases} \quad \mathbf{b} = (1, 2)^t, \quad (10.11a)$$

which gives rise to straight streamlines which are not aligned with the triangulation, and

$$u(x, y) = \begin{cases} 0, & x^2 + y^2 > 1, \\ \sin(\pi x) \sin(\pi y), & x^2 + y^2 < 1, \end{cases} \quad \mathbf{b} = (y, -x)^t, \quad (10.11b)$$

with a circular rotation around the origin that cannot be captured by triangular elements. We define accordingly the right-hand side functions f . In the spirit of Remark 10.6 and following Definition 10.5, we obtain $s_h \in H^1(\Omega)$, whereas the exact solution has a discontinuity and lies in $H_0(\mathcal{L}, \Omega) \setminus H^1(\Omega)$. The velocity field \mathbf{b} of (10.11b) does not satisfy the sufficient condition (10.2) to be Ω -filling. Nevertheless, one can verify that it is in fact Ω -filling with $T = 1/4$, so that we take $C_{\text{opt}} = 1/2$. In the first case, $\alpha = |\mathbf{b}|$, and we take $C_{\text{opt}} = 2\text{diam}(\Omega)/\alpha$.

The results are presented in Tables 7 and 8. One first observes that the rate of convergence of η_{NC} can be much worse compared to $\|u - u_h\|$, originating from the fact that s_h is a less accurate reconstruction of u . Despite the fact that the effectively indices are larger, they still remain rather independent of the mesh refinement and the polynomial degree of approximation. In Figures 3 and 4, the distributions of the error and of the error estimators $\eta_K := (\eta_{\text{NC}, K}^2 + (\eta_{\text{R}, K})^2)^{1/2}$ are presented, showing again a very close behavior.

11. CONCLUSIONS

In this work, we proposed a local reconstruction for numerical approximations of the one-dimensional linear advection equation, easily and independently obtained on each vertex patch. The reconstruction is proved to be well-posed and leads to a guaranteed upper bound of the L^2 -norm error between the actual solution u and the approximation u_h . This error estimator is also proved to be locally efficient with the local efficiency constant only depending on mesh shape-regularity. These results hold in a unified framework that only requires the residual of u_h to satisfy an orthogonality condition with respect to the hat basis functions. Numerical illustrations support the theory and additionally suggest asymptotic exactness. Motivated by these results, a heuristic extension to any space dimension is presented, with numerical experiments in 2D being rather encouraging.

Acknowledgements. This work was funded by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement No. 647134 GATIPOR).

REFERENCES

- [1] O. Axelsson, J. Karátson and B. Kovács, Robust preconditioning estimates for convection-dominated elliptic problems via a streamline Poincaré-Friedrichs inequality. *SIAM J. Numer. Anal.* **52** (2014) 2957–2976.
- [2] B. Ayuso and L.D. Marini, Discontinuous Galerkin methods for advection-diffusion-reaction problems. *SIAM J. Numer. Anal.* **47** (2009) 1391–1420.
- [3] P. Azérad and J. Pousin, Inégalité de Poincaré courbe pour le traitement variationnel de l'équation de transport. *C. R. Acad. Sci. Paris Sér. I Math.* **322** (1996) 721–727.
- [4] R. Becker, D. Capatina and R. Luce, Reconstruction-based *a posteriori* error estimators for the transport equation. In: Numerical Mathematics and Advanced Applications 2011. Springer, Berlin-Heidelberg (2013) 13–21.
- [5] K.S. Bey and J.T. Oden, *hp*-version discontinuous Galerkin methods for hyperbolic conservation laws. *Comput. Methods Appl. Mech. Eng.* **133** (1996) 259–286.
- [6] J. Blechta, J. Málek and M. Vohralík, Localization of the $W^{-1,q}$ norm for local *a posteriori* efficiency. *IMA J. Numer. Anal.* **40** (2019) 914–950.
- [7] D. Braess, V. Pillwein and J. Schöberl, Equilibrated residual error estimates are *p*-robust. *Comput. Methods Appl. Mech. Eng.* **198** (2009) 1189–1197.
- [8] P. Cantin, Well-posedness of the scalar and the vector advection-reaction problems in Banach graph spaces. *C. R. Math. Acad. Sci. Paris* **355** (2017) 892–902.
- [9] P. Cantin and A. Ern, An edge-based scheme on polyhedral meshes for vector advection-reaction equations. *ESAIM: M2AN* **51** (2017) 1561–1581.
- [10] C. Carstensen and S.A. Funken, Fully reliable localized error control in the FEM. *SIAM J. Sci. Comput.* **21** (1999) 1465–1484.
- [11] W. Dahmen and R.P. Stevenson, Adaptive strategies for transport equations. *Comput. Methods Appl. Math.* **19** (2019) 431–464.
- [12] W. Dahmen, C. Huang, C. Schwab and G. Welper, Adaptive Petrov–Galerkin methods for first order transport equations. *SIAM J. Numer. Anal.* **50** (2012) 2420–2445.
- [13] A. Devinatz, R. Ellis and A. Friedman, The asymptotic behavior of the first real eigenvalue of second order elliptic operators with a small parameter in the highest derivatives. II. *Indiana Univ. Math. J.* **23** (1973–1974) 991–1011.
- [14] A. Ern and J.-L. Guermond, Discontinuous Galerkin methods for Friedrichs' systems. I. General theory. *SIAM J. Numer. Anal.* **44** (2006) 753–778.
- [15] A. Ern and M. Vohralík, Polynomial-degree-robust *a posteriori* estimates in a unified setting for conforming, nonconforming, discontinuous Galerkin, and mixed discretizations. *SIAM J. Numer. Anal.* **53** (2015) 1058–1081.
- [16] A. Ern and M. Vohralík, Stable broken H^1 and $H(\text{div})$ polynomial extensions for polynomial-degree-robust potential and flux reconstruction in three space dimensions. *Math. Comput.* **89** (2020) 551–594.
- [17] A. Ern, A.F. Stephansen and M. Vohralík, Guaranteed and robust discontinuous Galerkin *a posteriori* error estimates for convection-diffusion-reaction problems. *J. Comput. Appl. Math.* **234** (2010) 114–130.
- [18] K.O. Friedrichs, Symmetric positive linear differential equations. *Comm. Pure Appl. Math.* **11** (1958) 333–418.
- [19] E.H. Georgoulis, E. Hall and C. Makridakis, Error control for discontinuous Galerkin methods for first order hyperbolic problems. In: Vol. 157 of Recent Developments in Discontinuous Galerkin Finite Element Methods for Partial Differential Equations. *IMA Vol. Math. Appl.* Springer, Cham (2014) 195–207.
- [20] E.H. Georgoulis, E. Hall and C. Makridakis, An *a posteriori* error bound for discontinuous Galerkin approximations of convection-diffusion problems. *IMA J. Numer. Anal.* **39** (2019) 34–60.
- [21] J.-L. Guermond, A finite element technique for solving first-order PDEs in l^p . *SIAM J. Numer. Anal.* **42** (2004) 714–737.
- [22] F. Hecht, New development in FreeFEM++. *J. Numer. Math.* **20** (2012) 251–265.
- [23] P. Houston, J.A. Mackenzie, E. Süli, and G. Warnecke, *A posteriori* error analysis for numerical approximations of Friedrichs systems. *Numer. Math.* **82** (1999) 433–470.
- [24] P.D. Lax and R.S. Phillips, Local boundary conditions for dissipative symmetric linear differential operators. *Comm. Pure Appl. Math.* **13** (1960) 427–455.
- [25] C. Makridakis and R.H. Nochetto, *A posteriori* error analysis for higher order dissipative methods for evolution problems. *Numer. Math.* **104** (2006) 489–514.
- [26] I. Muga, M.J. Tyler and K. van der Zee, The discrete-dual minimal-residual method (DDMRes) for weak advection-reaction problems in Banach spaces. Preprint [arXiv:1808.04542](https://arxiv.org/abs/1808.04542) (2018).
- [27] G. Sangalli, Analysis of the advection-diffusion operator using fractional order norms. *Numer. Math.* **97** (2004) 779–796.
- [28] G. Sangalli, A uniform analysis of nonsymmetric and coercive linear operators. *SIAM J. Math. Anal.* **36** (2005) 2033–2048.
- [29] G. Sangalli, Robust *a posteriori* estimator for advection-diffusion-reaction problems. *Math. Comput.* **77** (2008) 41–70.
- [30] D. Schötzau and L. Zhu, A robust *a posteriori* error estimator for discontinuous Galerkin methods for convection-diffusion equations. *Appl. Numer. Math.* **59** (2009) 2236–2255.
- [31] E. Süli, *A posteriori* error analysis and adaptivity for finite element approximations of hyperbolic problems. In: An Introduction to Recent Developments in Theory and Numerics for Conservation Laws (Freiburg/Littenweiler, 1997) Vol. 5 of *Lect. Notes Comput. Sci. Eng.* Springer, Berlin-Heidelberg (1999) 123–194.

- [32] Z. Tang, <https://who.rocq.inria.fr/Zuqi.Tang/freefem++.html> (2015).
- [33] D.S. Tartakoff, Regularity of solutions to boundary value problems for first order systems. *Indiana Univ. Math. J.* **21** (1972) 1113–1129.
- [34] R. Verfürth, Robust a posteriori error estimates for stationary convection–diffusion equations. *SIAM J. Numer. Anal.* **43** (2005) 1766–1782.
- [35] M. Vohralík and M. Zakerzadeh, Guaranteed and robust L^2 -norm a posteriori error estimates for 1D linear advection–reaction problems. In preparation (2020).