# Robust error bounds for finite element approximation of reaction-diffusion problems with non-constant reaction coefficient in arbitrary space dimension

*Mark Ainsworth*

*Tomáš Vejchodský*

# ROBUST ERROR BOUNDS FOR FINITE ELEMENT APPROXIMATION OF REACTION-DIFFUSION PROBLEMS WITH NON-CONSTANT REACTION COEFFICIENT IN ARBITRARY SPACE DIMENSION

MARK AINSWORTH AND TOMÁŠ VEJCHODSKÝ

ABSTRACT. We present a fully computable a posteriori error estimator for piecewise linear finite element approximations of reaction-diffusion problems with mixed boundary conditions and piecewise constant reaction coefficient formulated in arbitrary dimension. The estimator provides a guaranteed upper bound on the energy norm of the error and it is robust for all values of the reaction coefficient, including the singularly perturbed case. The approach is based on robustly equilibrated boundary flux functions [1] and on subsequent robust and explicit flux reconstruction. This paper simplifies and extends the applicability of the previous result [2] in three aspects: (i) arbitrary dimension, (ii) mixed boundary conditions, and (iii) non-constant reaction coefficient. It is the first robust upper bound on the error with these properties. An auxiliary result that is of independent interest is the derivation of new explicit constants for two types of trace inequalities on simplices.

## 1. INTRODUCTION

Consider a linear reaction-diffusion problem in a domain $\Omega \subset \mathbb{R}^d$ with mixed boundary conditions:

$$-\Delta u + \kappa^2 u = f \quad \text{in } \Omega; \qquad u = 0 \quad \text{on } \Gamma_{\mathrm{D}}; \qquad \partial u / \partial \boldsymbol{n} = g_{\mathrm{N}} \quad \text{on } \Gamma_{\mathrm{N}}, \qquad (1)$$

where $\boldsymbol{n}$ stands for the unit outward normal vector to the boundary $\partial \Omega$. The dimension $d \geq 2$ is chosen arbitrarily. For simplicity we assume $\Omega$ to be a polytope. The portions $\Gamma_{\mathrm{D}}$ and $\Gamma_{\mathrm{N}}$ of the boundary $\partial \Omega$ are open, disjoint, and satisfy $\overline{\Gamma}_{\mathrm{D}} \cup \overline{\Gamma}_{\mathrm{N}} = \partial \Omega$. The reaction coefficient $\kappa \geq 0$ is considered to be piecewise constant. In order to guarantee unique solvability of (1), we consider $\kappa > 0$ in a subdomain of $\Omega$ of a positive measure or a positive measure of $\Gamma_{\mathrm{D}}$. We use the

finite element method to approximate the exact solution $u$ by a piecewise affine function $u_h$ with respect to a simplicial partition $\mathcal{T}_h$ of $\Omega$.

In this paper we derive a computable a posteriori error estimate based on robust flux equilibration and explicit flux reconstruction. This error estimate $\eta$ provides a guaranteed and fully computable upper bound on the energy norm of the error $\|u - u_h\|$ and it is robust with respect to both $\kappa$ and the mesh-size $h$.

A posteriori error estimates are useful for adaptive algorithms, where they play two roles. Firstly, they indicate where the computational mesh should be refined or coarsened. Secondly, they provide quantitative information about the size of the error for reliable stopping criterion. Unfortunately, many existing estimators do not provide actual numerical bounds that can be used as a stopping criterion.

Adaptive algorithms are convergent [3] provided the error estimates are *locally efficient* and *reliable*. If $\eta_K$ stand for local error indicators on elements $K \in \mathcal{T}_h$ and $\eta^2 = \sum_{K \in \mathcal{T}_h} \eta_K^2$ is the global error estimator, then the indicators $\eta_K$ are said to be *locally efficient* if there exists a constant $c > 0$ such that

$$c\eta_K \leq \|u - u_h\|_{\widetilde{K}},$$

where $\|u - u_h\|_{\widetilde{K}}$ stands for the energy norm restricted to a patch $\widetilde{K}$ of elements consisting of $K$ and neighbouring elements sharing at least one vertex with $K$. Similarly, the error estimator $\eta$ is *reliable* if there exists a constant $C > 0$ such that

$$\|u - u_h\| \leq C\eta.$$

The error estimate $\eta$ is *robust* if the constants $c$ and $C$ are independent of $\kappa$ and mesh-size $h$. The error estimate $\eta$ is a *guaranteed upper bound* if $\|u - u_h\| \leq \eta$, i.e. the reliability constant $C$ is equal to one. Finally, the error bound $\eta$ is *fully computable* if it can be evaluated in terms of the approximation $u_h$ and given data without the need for generic (unknown) constants.

A robust, reliable, locally efficient explicit a posteriori error estimate for problem (1) was first derived by Verfürth in [4]. This estimate, however, does not provide guaranteed upper bound on the error. An estimator which does provide an upper bound along with robust local efficiency was obtained by Ainsworth and Babuška in [5], but this upper bound depends on an exact solution of a Neumann problem and as noted in [5] is not fully computable. Subsequently in [2] we were able to develop fully computable error bounds in the two dimensional setting by a complementarity technique combined with robustly equilibrated fluxes and explicit flux reconstruction. Here, we develop a simpler flux reconstruction that is suitable for any dimension $d \geq 2$ and is applicable to the case of piecewise constant coefficient $\kappa$ including the situation where $\kappa$ can be very large in some parts of the domain and vanishingly small in others. Furthermore, we extend the previous results by considering nonhomogeneous Neumann boundary conditions. In order to achieve these goals, we develop some new techniques and tools for the analysis that are of wider applicability than the problem addressed here.

The question of robust a posteriori error estimates for singularly perturbed problems is studied by other authors as well. In [6], an error estimate that is robust with respect to anisotropic meshes is obtained, but unfortunately does not provide guaranteed upper bound on the error. A robust, locally efficient and fully computable guaranteed upper bound was obtained in [7] for the finite volume method and $d = 2$ and 3. Recently, a robust estimator for the error in the maximum norm was obtained in [8] for the case $d = 1$.

The basic idea behind our work can be traced back to the method of the hypercircle [9] and later to [10, 11, 12]. This approach has been adopted by Repin [13] and his group for a wide class of problems in conjunction with the solution of a global minimization problem to compute the error bound. We avoid any global computations and instead develop local algorithms for guaranteed and fully computable error bounds based on *flux equilibration* [5, 14, 15, 16, 17, 18, 19, 20] etc. In the present work we will utilize the robust flux equilibration from [5].

The rest of the paper is organized as follows. Section 2 defines the finite element approximation and corresponding assumptions. The core of the paper lies in Section 3, where we present new trace inequalities on simplices, and develop two new flux reconstructions both of which are used to derive the a posteriori error and our main result. Finally, Section 4 provides an illustrative numerical example and Section 5 draws the conclusions.

## 2. Model Problem and Its Approximate Solution

The weak formulation of (1) reads: find $u \in V = \{v \in H^1(\Omega) : v = 0 \text{ on } \Gamma_{\mathrm{D}}\}$ such that

$$\mathcal{B}(u, v) = \mathcal{F}(v) \quad \forall v \in V, \tag{2}$$

where $\mathcal{B}$ and $\mathcal{F}$ are bilinear and linear forms, respectively, defined on $V$ by

$$\mathcal{B}(u, v) = \int_\Omega (\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} v + \kappa^2 uv) \, \mathrm{d}\boldsymbol{x}; \quad \mathcal{F}(v) = \int_\Omega fv \, \mathrm{d}\boldsymbol{x} + \int_{\Gamma_{\mathrm{N}}} g_{\mathrm{N}} v \, \mathrm{d}\boldsymbol{s}.$$

In order to discretize problem (1) we consider a family of partitions $\mathcal{G} = \{\mathcal{T}_h\}$ of the domain $\Omega$. Each partition $\mathcal{T}_h$ consists of simplices (elements), their union is $\overline{\Omega}$, their interiors are pairwise disjoint, and every facet of each simplex lies either in $\partial\Omega$ or it is completely shared by exactly two neighbouring simplices. We assume that all partitions in $\mathcal{G}$ are compatible the coefficient $\kappa$ meaning that $\kappa$ is a constant $\kappa_K$ in any element $K$ of $\mathcal{T}_h$ for all $\mathcal{T}_h \in \mathcal{G}$.

We denote by $h_K$, $\boldsymbol{x}_K$, $\rho_K$, and $\boldsymbol{n}_K$ the diameter, the *incentre*, the *inradius* of simplex $K$, and the unit outward-facing normal vector to the boundary $\partial K$, respectively. The family of partitions $\mathcal{G}$ is assumed to be regular, i.e. there exists a constant $C > 0$ such that

$$\sup_{\mathcal{T}_h \in \mathcal{G}} \max_{K \in \mathcal{T}_h} \frac{h_K}{\rho_K} \leq C, \tag{3}$$

but is not requested to be quasi-uniform, thereby permitting the use of locally refined meshes. Throughout the paper we use symbol $C$ for a generic constant whose value is independent of $\kappa$ and any mesh-size and whose actual numerical value can differ in different occurrences. Furthermore, we define

$$\widetilde{K} = \mathrm{int}\left\{\bigcup \overline{K'} : \overline{K'} \cap \overline{K} \neq \emptyset\right\} \tag{4}$$

to be the patch consisting of $K$ and elements sharing at least one common point with $K$.

The regularity assumption implies several facts that we will use in the subsequent analysis. Firstly, the number of elements in any patch is uniformly bounded over the family $\mathcal{G}$ as is the number of patches containing a particular element. Secondly, within each patch $\widetilde{K}$, a local quasi-uniformity condition $ch_K \leq h_{K'} \leq Ch_K$ holds for all elements $K' \subset \widetilde{K}$ with uniform constants $c > 0$ and $C > 0$ over the family $\mathcal{G}$. Thirdly, the elements are shape regular meaning that there exists a positive constant $\mathcal{C}_0$ such that

$$\frac{1}{\mathcal{C}_0}\rho_K \leq \rho_{K'} \leq \mathcal{C}_0\rho_K \tag{5}$$

for all elements $K' \subset \widetilde{K}$, all $K \in \mathcal{T}_h$, and all $\mathcal{T}_h \in \mathcal{G}$.

The coefficient $\kappa$ is assumed to be piecewise constant, and such that for some constant $C > 0$ the following conditions hold for all triangulations $\mathcal{T}_h \in \mathcal{G}$ and all elements $K \in \mathcal{T}_h$:

$$\text{if} \quad \kappa_K \neq 0 \quad \text{then} \quad \kappa_{K'} \leq C\kappa_K \quad \text{for all } K' \subset \widetilde{K}; \tag{6}$$

$$\text{if} \quad \kappa_K = 0 \quad \text{then} \quad \kappa_{K'} \leq C \quad \text{for all } K' \subset \widetilde{K}. \tag{7}$$

These assumptions rule out the case of arbitrarily high jumps in values $\kappa_K$ between neighbouring elements.

One consequence of assumptions (6)–(7) together with the quasi-uniformity of $h_K$ is existence of a constant $C > 0$ such that for all $\mathcal{T}_h \in \mathcal{G}$, all $K \in \mathcal{T}_h$, and all elements $K' \subset \widetilde{K}$, we have

$$C^{-1}\min\{h_{K'}, \kappa_{K'}^{-1}\} \leq \min\{h_K, \kappa_K^{-1}\} \leq C\min\{h_{K'}, \kappa_{K'}^{-1}\}. \tag{8}$$

The quantity $\min\{h_K, \kappa_K^{-1}\}$ appears extensively throughout the paper, and for the avoidance of doubt, we note explicitly that

$$\min\{h_K, \kappa_K^{-1}\} = h_K \quad \text{if } \kappa_K = 0. \tag{9}$$

Let $X_h$ be the space of continuous and piecewise affine functions with respect to the partition $\mathcal{T}_h$, and let the subspace $V_h = X_h \cap H_0^1(\Omega)$. The finite element approximation $u_h \in V_h$ of (1) is then given by

$$\mathcal{B}(u_h, v_h) = \mathcal{F}(v_h) \quad \forall v_h \in V_h. \tag{10}$$

Finally, we use local counterparts of the bilinear and linear forms defined by

$$\mathcal{B}_K(u,v) = \int_K (\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} v + \kappa_K^2 uv)\,\mathrm{d}\boldsymbol{x}; \quad \mathcal{F}_K(v) = \int_K fv\,\mathrm{d}\boldsymbol{x} + \int_{\Gamma_{\mathrm{N}} \cap \partial K} g_{\mathrm{N}} v\,\mathrm{d}\boldsymbol{s}.$$

The associated global and local energy norms $\|\!|\cdot|\!\|$ and $\|\!|\cdot|\!\|_K$ are defined by $\|\!|v|\!\|^2 = \mathcal{B}(v,v)$ and $\|\!|v|\!\|_K^2 = \mathcal{B}_K(v,v)$, respectively. Analogously, we use $\|\cdot\|$ and $\|\cdot\|_K$ for the $L^2(\Omega)$ and $L^2(K)$ norms, respectively.

## 3. A Posteriori Error Estimator

3.1. **Trace inequalities on simplices.** The derivation of complementarity based error estimates for problems with nonhomogeneous Neumann boundary conditions requires certain types of trace inequalities. Moreover, the constants appearing in these inequalities are present in the final error bounds. We derive two new trace inequalities for simplices together with explicit formulas for the corresponding constants.

**Lemma 1.** *Let $K$ be a $d$-dimensional non-degenerate simplex and let $\gamma$ be one of its facets. Let $h_K$ be the diameter of $K$ and $\kappa_K \geq 0$ a constant. Let $v \in H^1(K)$ and let $\bar{v}_\gamma$ denote the average value of $v$ on $\gamma$. Then*

$$\|v\|_\gamma \leq C_{\mathrm{T}} \|\!|v|\!\|_K \quad \text{for } \kappa_K > 0, \tag{11}$$

$$\|v - \bar{v}_\gamma\|_\gamma \leq \overline{C}_{\mathrm{T}} \|\!|v|\!\|_K, \tag{12}$$

*hold with constants $C_{\mathrm{T}}, \overline{C}_{\mathrm{T}} > 0$ given by*

$$C_{\mathrm{T}}^2 = \frac{|\gamma|}{|K|} \frac{1}{d+1} \frac{1}{\kappa_K} \sqrt{(2h_K)^2 + (d/\kappa_K)^2},$$

$$\overline{C}_{\mathrm{T}}^2 = \frac{|\gamma|}{|K|} \frac{1}{d+1} \min\{h_K/\pi, \kappa_K^{-1}\} \left(2h_K + d\min\{h_K/\pi, \kappa_K^{-1}\}\right).$$

*Proof.* Let $\boldsymbol{x}_0$ be the vertex of $K$ opposite to the facet $\gamma$. Define $\boldsymbol{\varphi}(\boldsymbol{x}) = \boldsymbol{x} - \boldsymbol{x}_0$ for $\boldsymbol{x} \in K$. Note that $\boldsymbol{n}_K \cdot \boldsymbol{\varphi} = 0$ on $\partial K \setminus \gamma$ and $\boldsymbol{n} \cdot \boldsymbol{\varphi} = \tilde{\varrho}_K$ on $\gamma$, where $\boldsymbol{n}_K$ is the unit outward normal to $\partial K$ and $\tilde{\varrho}_K$ is the distance between $\gamma$ and $\boldsymbol{x}_0$, i.e. the altitude of $K$. In particular, $\tilde{\varrho}_K = (d+1)|K|/|\gamma|$.

Let $v \in H^1(K)$ then

$$\frac{d+1}{|\gamma|}|K| \|v\|_\gamma^2 = \int_\gamma v^2 \boldsymbol{n}_K \cdot \boldsymbol{\varphi}\,\mathrm{d}\boldsymbol{s} = \int_{\partial K} v^2 \boldsymbol{n}_K \cdot \boldsymbol{\varphi}\,\mathrm{d}\boldsymbol{s} = \int_K \mathrm{div}(v^2 \boldsymbol{\varphi})\,\mathrm{d}\boldsymbol{x}$$

$$= 2\int_K v\boldsymbol{\varphi} \cdot \boldsymbol{\nabla} v\,\mathrm{d}\boldsymbol{x} + \int_K v^2 \,\mathrm{div}\,\boldsymbol{\varphi}\,\mathrm{d}\boldsymbol{x} \leq \|v\|_K \left(2h_K \|\boldsymbol{\nabla} v\|_K + d\|v\|_K\right). \tag{13}$$

Using $\|v\|_K \leq \kappa^{-1}\|\!|v|\!\|_K$ and $2h_K\|\boldsymbol{\nabla} v\|_K + d\|v\|_K \leq \left((2h_K)^2 + (d/\kappa_K)^2\right)^{1/2}\|\!|v|\!\|_K$ in (13), we obtain (11).

Now, consider $\bar{v}_\gamma = |\gamma|^{-1} \int_\gamma v \, \mathrm{d}s$ and $\bar{v}_K = |K|^{-1} \int_K v \, \mathrm{d}x$. Applying estimate (13) to $v - \bar{v}_K$ yields

$$\|v - \bar{v}_\gamma\|_\gamma^2 \leq \|v - \bar{v}_K\|_\gamma^2 \leq \frac{|\gamma|}{|K|} \frac{1}{d+1} \|v - \bar{v}_K\|_K \left(2h_K \|\boldsymbol{\nabla} v\|_K + d \|v - \bar{v}_K\|_K\right).$$
(14)

The norm $\|v - \bar{v}_K\|_K$ can be bounded in either of the two ways:

$$\|v - \bar{v}_K\|_K \leq \|v\|_K \leq \kappa_K^{-1} \|v\|_K \quad \text{and} \quad \|v - \bar{v}_K\|_K \leq \frac{h_K}{\pi} \|\boldsymbol{\nabla} v\|_K \leq \frac{h_K}{\pi} \|v\|_K,$$

where we use Poincaré inequality [21]. Thus, $\|v - \bar{v}_K\|_K \leq \min\{h_K/\pi, \kappa_K^{-1}\} \|v\|_K$. Using this estimate and inequality $\|\boldsymbol{\nabla} v\|_K \leq \|v\|_K$ in (14), we derive (12). □

The constants $C_{\mathrm{T}}$ and $\overline{C}_{\mathrm{T}}$ from Lemma 1 have the correct asymptotic behaviour with respect to $h_K$ and $\kappa_K$, but they are not optimal in terms of absolute values. Optimal values for trace constants are not known in general, but their two-sided bounds can be computed numerically for quite general domains [22].

3.2. **General framework.** We define $\Pi_K : L^2(K) \to \mathbb{P}^1(K)$ to be the $L^2(K)$-orthogonal projector to the space of affine functions defined over an element $K \in \mathcal{T}_h$. Similarly, for a facet $\gamma \subset \partial K$ we define $\Pi_\gamma : L^2(\gamma) \to \mathbb{P}^1(\gamma)$ to be the $L^2(\gamma)$-orthogonal projector to the space of affine functions defined over the facet $\gamma \subset \partial K$. The following generalization of the corresponding result in [2] forms the basis of our approach:

**Lemma 2.** *Let $u \in V$ be the weak solution (2) and $u_h \in V$ be an arbitrary function. Further let $\boldsymbol{\tau} \in \boldsymbol{H}(\mathrm{div}, \Omega)$ be such that $\Pi_K f + \mathrm{div}\,\boldsymbol{\tau} = 0$ in those elements $K \in \mathcal{T}_h$ where $\kappa_K = 0$ and $\boldsymbol{\tau} \cdot \boldsymbol{n} = \Pi_\gamma g_{\mathrm{N}}$ on facets $\gamma \subset \Gamma_{\mathrm{N}} \cap \partial K$. Then*

$$\|u - u_h\|^2 \leq \sum_{K \in \mathcal{T}_h} \left[ \eta_K(\boldsymbol{\tau}) + \mathrm{osc}_K(f) + \sum_{\gamma \subset \Gamma_{\mathrm{N}} \cap \partial K} \mathrm{osc}_\gamma(g_{\mathrm{N}}) \right]^2$$

*where $\eta_K(\boldsymbol{\tau}) \geq 0$, $\mathrm{osc}_K(f)$, and $\mathrm{osc}_\gamma(g_{\mathrm{N}})$ are defined by*

$$\eta_K^2(\boldsymbol{\tau}) = \begin{cases} \|\boldsymbol{\tau} - \boldsymbol{\nabla} u_h\|_K^2 + \kappa_K^{-2} \|\Pi_K f - \kappa_K^2 u_h + \mathrm{div}\,\boldsymbol{\tau}\|_K^2 & \text{if } \kappa_K > 0, \\ \|\boldsymbol{\tau} - \boldsymbol{\nabla} u_h\|_K^2 & \text{if } \kappa_K = 0, \end{cases}$$
(15)

$$\mathrm{osc}_K(f) = \min\left\{ \frac{h_K}{\pi}, \frac{1}{\kappa_K} \right\} \|f - \Pi_K f\|_K,$$

$$\mathrm{osc}_\gamma(g_{\mathrm{N}}) = \min\{C_{\mathrm{T}}, \overline{C}_{\mathrm{T}}\} \|g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}\|_\gamma.$$

*Proof.* Let $v \in V$ be arbitrary. Using the weak formulation (2) for $u$, the fact that the global forms $\mathcal{B}$ and $\mathcal{F}$ are sums of the local forms $\mathcal{B}_K$ and $\mathcal{F}_K$, and the

divergence theorem, we obtain the following identity

$$
\begin{aligned}
\mathcal{B}(u - u_h, v) = \sum_{K \in \mathcal{T}_h} \mathcal{F}_K(v) - \mathcal{B}_K(u_h, v) = \sum_{K \in \mathcal{T}_h} \Bigg[ & \int_K (\boldsymbol{\tau} - \boldsymbol{\nabla} u_h) \cdot \boldsymbol{\nabla} v \, \mathrm{d}\boldsymbol{x} \\
+ \int_K (\Pi_K f - \kappa_K^2 u_h + \operatorname{div} \boldsymbol{\tau}) v \, \mathrm{d}\boldsymbol{x} & + \sum_{\gamma \subset \Gamma_{\mathrm{N}} \cap \partial K} \int_\gamma (\Pi_\gamma g_{\mathrm{N}} - \boldsymbol{\tau} \cdot \boldsymbol{n}) v \, \mathrm{d}\boldsymbol{s} \\
+ \int_K (f - \Pi_K f) v \, \mathrm{d}\boldsymbol{x} & + \sum_{\gamma \subset \Gamma_{\mathrm{N}} \cap \partial K} \int_\gamma (g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}) v \, \mathrm{d}\boldsymbol{s} \Bigg]. \quad (16)
\end{aligned}
$$

Now we estimate the five integrals on the right-hand side of (16). The sum of the first two integrals is clearly bounded by $\eta_K(\boldsymbol{\tau}) \|v\|_K$ for both $\kappa_K > 0$ and $\kappa_K = 0$. The third integral on the right-hand side of (16) vanishes since $\boldsymbol{\tau} \cdot \boldsymbol{n} = \Pi_\gamma g_{\mathrm{N}}$ on $\Gamma_{\mathrm{N}}$.

The fourth integral can be estimated as

$$
\int_K (f - \Pi_K f) v \, \mathrm{d}\boldsymbol{x} \leq \min \left\{ \frac{h_K}{\pi}, \frac{1}{\kappa_K} \right\} \|f - \Pi_K f\|_K \|v\|_K = \operatorname{osc}_K(f) \|v\|_K,
$$

where the constant $h_K/\pi$ comes from the Poincaré inequality [21] and $1/\kappa_K$ comes from the inequality $\|v\|_K \leq \kappa_K^{-1} \|v\|_K$, see [2, p. 228] for details. The last integral in (16) can be bounded in the following two ways:

$$
\int_\gamma (g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}) v \, \mathrm{d}\boldsymbol{s} \leq \|g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}\|_\gamma \|v\|_\gamma,
$$

$$
\int_\gamma (g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}) v \, \mathrm{d}\boldsymbol{s} = \int_\gamma (g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}})(v - \bar{v}_\gamma) \, \mathrm{d}\boldsymbol{s} \leq \|g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}\|_\gamma \|v - \bar{v}_\gamma\|_\gamma,
$$

where $\bar{v}_\gamma = |\gamma|^{-1} \int_\gamma v \, \mathrm{d}\boldsymbol{s}$. Employing trace inequalities (11) and (12) we end up with the estimate

$$
\int_\gamma (g_{\mathrm{N}} - \Pi_\gamma g_{\mathrm{N}}) v \, \mathrm{d}\boldsymbol{s} \leq \operatorname{osc}_\gamma(g_{\mathrm{N}}) \|v\|_K.
$$

Hence,

$$
\mathcal{B}(u - u_h, v) \leq \sum_{K \in \mathcal{T}_h} \left[ \eta_K(\boldsymbol{\tau}) + \operatorname{osc}_K(f) + \sum_{\gamma \subset \Gamma_{\mathrm{N}} \cap \partial K} \operatorname{osc}_\gamma(g_{\mathrm{N}}) \right] \|v\|_K.
$$

The Cauchy-Schwarz inequality and substitution $v = u - u_h$ finishes the proof. $\quad\square$

The vector field $\boldsymbol{\tau} \in \boldsymbol{H}(\operatorname{div}, \Omega)$ is referred to as a *flux reconstruction* and its specific choice is crucial for the efficiency and robustness of the resulting error estimators. We reconstruct the flux $\boldsymbol{\tau} \in \boldsymbol{H}(\operatorname{div}, \Omega)$ in two steps. Firstly, we find

boundary fluxes $g_K$ satisfying the following conditions:

$$g_K|_\gamma \in \mathbb{P}^1(\gamma) \quad \text{for all facets } \gamma \subset \partial K, \tag{17}$$

$$g_K = \Pi_\gamma g_{\mathrm{N}} \quad \text{for all facets } \gamma \subset \Gamma_{\mathrm{N}} \cap \partial K, \tag{18}$$

$$g_K + g_{K'} = 0 \qquad \text{on facets } \gamma = \partial K \cap \partial K'. \tag{19}$$

Secondly, we locally reconstruct vector fields $\boldsymbol{\tau}_K \in \boldsymbol{H}(\mathrm{div}, K)$ satisfying boundary conditions $\boldsymbol{\tau}_K \cdot \boldsymbol{n}_K = g_K$ on $\partial K$. The values of $\boldsymbol{\tau}$ in the interior of the elements will be presented in detail below. Irrespective of this, the resulting vector field $\boldsymbol{\tau}$ is defined elementwise by $\boldsymbol{\tau}|_K = \boldsymbol{\tau}_K$ for all $K \in \mathcal{T}_h$, so that $\boldsymbol{\tau} \in \boldsymbol{H}(\mathrm{div}, \Omega)$ due to the consistency condition (19).

Conditions (17)–(19) do not determine a unique set of fluxes. The specific choice of fluxes satisfying (17)–(19) will be crucial to the robustness of the associated estimator. We say that boundary fluxes $g_K$ are equilibrated with respect to linear functions if the following condition holds:

$$\int_K f\theta \, \mathrm{d}\boldsymbol{x} - \mathcal{B}_K(u_h, \theta) + \int_{\partial K} g_K \theta \, \mathrm{d}\boldsymbol{s} = 0 \quad \forall \theta \in \mathbb{P}^1(K). \tag{20}$$

or, equally well,

$$\int_K (\boldsymbol{\tau} - \boldsymbol{\nabla}u_h) \cdot \boldsymbol{\nabla}\theta \, \mathrm{d}\boldsymbol{x} + \int_K (f - \kappa_K^2 u_h + \mathrm{div}\,\boldsymbol{\tau})\theta \, \mathrm{d}\boldsymbol{x} = 0 \quad \forall \theta \in \mathbb{P}^1(K).$$

Fluxes satisfying the equilibration (20) yield accurate error bounds for small $\kappa$, but are not robust for large values of $\kappa$ [5]. Therefore, we will follow the approach from [5] for large values of $\kappa$.

### 3.3. Robust equilibration of boundary fluxes.

A detailed algorithm for the construction of boundary fluxes $g_K$ satisfying conditions (17)–(20) can be found in [1]. A modification of this approach that is robust for large values of $\kappa$ is described in [5] and [2] and we will briefly recall it here.

The idea is to replace the affine functions in (20) by their approximate minimum energy extensions. Clearly, it suffices to satisfy condition (20) for the barycentric coordinates $\theta_n$, $n = 1, 2, \ldots, d+1$, in $K$. The approximate minimum energy extensions $\theta_n^*$ of $\theta_n$ are defined in [5] for $d = 1$, 2, and 3 dimensions. Here, we define them for general $d$-dimensional simplices.

Consider a simplex $K$ with vertices $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{d+1}$ and facets $\gamma_1, \ldots, \gamma_{d+1}$ opposite to these vertices. The standard basis functions $\theta_n$ are determined by the conditions $\theta_n(\boldsymbol{x}_m) = \delta_{nm}$, $n, m \in \{1, 2, \ldots, d+1\}$. For each $n = 1, 2, \ldots, d+1$, define approximate minimum energy extension $\theta_n^*$ as follows. If $\kappa_K \rho_K \leq 1$ then $\theta_n^* = \theta_n$. If $\kappa_K \rho_K > 1$ then define a point $\boldsymbol{x}_P$ by its barycentric coordinates $\lambda_i(\boldsymbol{x}_P) = \delta$ for $i \neq n$ and $\lambda_n(\boldsymbol{x}_P) = 1 - d\delta$ with $\delta = \min\{1, 1/(\kappa_K \rho_K)\}/d$, and consider a submesh in $K$ created by simplices $K_i = \overline{\gamma_i \boldsymbol{x}_P}$, $i = 1, 2, \ldots, d+1$. The approximate minimum energy extension $\theta_n^*$ is then defined as a piecewise affine function with respect to this submesh such that $\theta_n^*(\boldsymbol{x}_n) = 1$, $\theta_n^*(\boldsymbol{x}_P) = 0$,
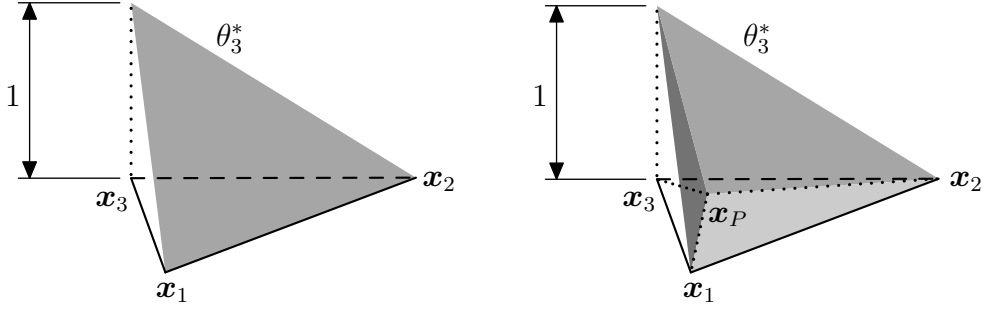
FIGURE 1. A graph of the approximate minimum energy extension $\theta_n^*$ for $\kappa_K \rho_K \leq 1$ (left) and for $\kappa_K \rho_K > 1$ (right).

and $\theta_n^*(\boldsymbol{x}_i) = 0$ for all $i \neq n$. A two-dimensional illustration of functions $\theta_n^*$ is provided in Figure 1.

It is easy to verify that $\theta_n^*$, $n \in \{1, 2, \ldots, d+1\}$, satisfy

- $\theta_n^* = \theta_n$ on the boundary $\partial K$;
- if $\kappa_K \rho_K \leq 1$ then $\theta_n^* = \theta_n$ on $K$;
- $C_1 h_K^{d-1} \min\{h_K, 1/\kappa_K\} \leq \|\theta_n^*\|_K^2 \leq C_2 h_K^{d-1} \min\{h_K, 1/\kappa_K\}$;
- $C_1 h_K^{d-1} \min\{h_K, 1/\kappa_K\}^{-1} \leq \|\boldsymbol{\nabla}\theta_n^*\|_K^2 \leq C_2 h_K^{d-1} \min\{h_K, 1/\kappa_K\}^{-1}$.

These key features of the approximate minimum energy extensions $\theta_n^*$ were identified already in [5] and are crucial for the robustness of the resulting fluxes.

The robust flux reconstruction is obtained in [5] by replacing functions $\theta_n$ by $\theta_n^*$ in (20) and finding the least squares minimizer of the system. This approach must be modified to deal with case of variable $\kappa$ considered here.

In particular we require

$$\mathcal{F}_K(\theta_n) - \mathcal{B}_K(u_h, \theta_n) + \int_{\partial K \backslash \Gamma_{\mathrm{N}}} g_K \theta_n \, \mathrm{d}\boldsymbol{s} = 0 \tag{21}$$

for all elements $K \in \mathcal{T}_h$, where $\kappa_K \rho_K \leq 1$, and for all $n = 1, 2, \ldots, d+1$. For elements $K \in \mathcal{T}_h$, where $\kappa_K \rho_K > 1$, we impose a similar condition in a least-squares sense:

$$\mathcal{F}_K(\theta_n^*) - \mathcal{B}_K(u_h, \theta_n^*) + \int_{\partial K \backslash \Gamma_{\mathrm{N}}} g_K \theta_n^* \, \mathrm{d}\boldsymbol{s} \approx 0 \quad \forall n = 1, 2, \ldots, d+1. \tag{22}$$

The resulting constrained least-squares problem (21)–(22) can be transformed into a series of small constrained least-squares problems on patches of elements corresponding to vertices of $\mathcal{T}_h$ as follows.

We define the set of vertices $\mathcal{N}(\gamma)$ of a facet $\gamma$ of a simplex $K$ and functions $\psi_\gamma^m \in \mathbb{P}^1(\gamma)$ satisfying $\int_\gamma \psi_\gamma^m \theta_n \, \mathrm{d}\boldsymbol{s} = \delta_{mn}$. Further, we consider a fixed orientation $\sigma_{K,\gamma}$ of facets $\gamma$ of simplices $K \in \mathcal{T}_h$. The orientation $\sigma_{K,\gamma}$ is either 1 or $-1$ and satisfies

$$\sigma_{K,\gamma} + \sigma_{K',\gamma} = 0 \quad \text{on } \gamma = \partial K \cap \partial K'.$$

Finally, we introduce the average and the jump flux across a common facet of two neighbouring simplices $K$ and $K'$ as

$$\left\langle \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right\rangle = \frac{1}{2} \boldsymbol{n}_K \cdot (\boldsymbol{\nabla} u_h|_K + \boldsymbol{\nabla} u_h|_{K'}) \quad \text{and} \quad \left[ \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right] = \boldsymbol{n}_K \cdot (\boldsymbol{\nabla} u_h|_K - \boldsymbol{\nabla} u_h|_{K'}).$$

On the boundary $\partial\Omega$ we set $\langle \partial u_h / \partial \boldsymbol{n}_K \rangle = \partial u_h / \partial \boldsymbol{n}_K$ and $[\partial u_h / \partial \boldsymbol{n}_K] = 0$. The boundary flux $g_K$ on a facet $\gamma$ of a simplex $K$ is then defined in the form

$$g_K = \left\langle \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right\rangle + \sigma_{K,\gamma} \sum_{m \in \mathcal{N}(\gamma)} \alpha_\gamma^m \psi_\gamma^m. \tag{23}$$

Notice that this construction of $g_K$ immediately guarantees the consistency condition (19). Furthermore, if $\gamma \subset \Gamma_\mathrm{N}$ then the coefficients $\alpha_\gamma^m$ are uniquely determined by (18).

The substitution (23) transforms the global constrained least-squares problem (21)–(22) into small local constrained least-squares problems

$$\sum_{\gamma:\gamma\subset\partial K\backslash\Gamma_\mathrm{N}, \gamma\ni\boldsymbol{x}_n} \sigma_{K,\gamma}\alpha_\gamma^n = -D_K(\theta_n) \quad \forall K \in \omega(\boldsymbol{x}_n), \text{where } \kappa_K\rho_K \le 1, \tag{24}$$

$$\sum_{\gamma:\gamma\subset\partial K\backslash\Gamma_\mathrm{N}, \gamma\ni\boldsymbol{x}_n} \sigma_{K,\gamma}\alpha_\gamma^n \approx -D_K(\theta_n^*) \quad \forall K \in \omega(\boldsymbol{x}_n), \text{where } \kappa_K\rho_K > 1, \tag{25}$$

for all vertices $\boldsymbol{x}_n$ in the partition $\mathcal{T}_h$, where we define $\omega(\boldsymbol{x}_n) = \{K \in \mathcal{T}_h : \boldsymbol{x}_n \in \overline{K}\}$ to be the set of elements sharing the vertex $\boldsymbol{x}_n$ and

$$D_K(\theta) = \mathcal{F}_K(\theta) - \mathcal{B}_K(u_h, \theta) + \int_{\partial K\backslash\Gamma_\mathrm{N}} \left\langle \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right\rangle \theta \, \mathrm{d}\boldsymbol{s}.$$

The summations in (24)–(25) are performed over those facets $\gamma \subset \partial K \setminus \Gamma_\mathrm{N}$ that contain the vertex $\boldsymbol{x}_n$.

Observe that the system (24)–(25) is always solvable. It was shown in [5] and [1] that if $\kappa_K\rho_K \le 1$ for all elements $K \in \omega(\boldsymbol{x}_n)$ then the linear system (24) is always solvable. Trivially, removing some (or all) equality constraints and replacing them by the requirement of least-squares fit (25), preserves the solvability of the remaining system of linear constraints (24).

To summarize, we require the satisfaction of the exact equilibration condition (24) for those elements where the coefficient $\kappa_K$ is small. For the other elements we mimic the equilibration by the least-squares fit (25). The resulting constrained least-squares problem (24)–(25) is always solvable and its solution depends continuously on the data.

3.4. **Auxiliary results.** In this section we recall several estimates from [5] and extend them to include the case of piecewise constant $\kappa$ and to Neumann boundary conditions. Lemma 5(2) from [5] says that if $\gamma$ is an interior facet (i.e. shared

by two elements) then

$$\left\| \left[ \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right] \right\|_\gamma \le C \left[ \min\{h_\gamma, \kappa_K^{-1}\}^{-\frac{1}{2}} \|u - u_h\|_{\widetilde{\gamma}} + \min\{h_\gamma, \kappa_K^{-1}\}^{\frac{1}{2}} \|f - \Pi f\|_{\widetilde{\gamma}} \right], \quad (26)$$

where $\widetilde{\gamma}$ is the pair of elements sharing the facet $\gamma$ and $\Pi f$ is defined piecewise by $(\Pi f)|_K = \Pi_K f$ for all $K \in \mathcal{T}_h$. Notice that thanks to the assumption (8) estimate (26) holds for any $K \in \widetilde{\gamma}$.

Further, Lemma 6 from [5] provides the estimate

$$\left\| g_K - \left\langle \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right\rangle \right\|_\gamma \le C \left[ \min\{h_\gamma, \kappa_K^{-1}\}^{-\frac{1}{2}} \|u - u_h\|_{\widetilde{K}} + \min\{h_\gamma, \kappa_K^{-1}\}^{\frac{1}{2}} \|f - \Pi f\|_{\widetilde{K}} \right],$$
$$(27)$$

for all $K \in \mathcal{T}_h$, where $\gamma$ is a facet of $K$ that is either interior or it lies on the Dirichlet boundary $\Gamma_\mathrm{D}$. The patch of elements $\widetilde{K}$ was defined in (4) and condition (8) is needed to generalize the proof from [5] to the case of piecewise constant $\kappa$. The final estimate on page 343 in [5] states that

$$\|\Pi_K r_h\|_K \le C \left[ \min\{h_K, \kappa_K^{-1}\} \|u - u_h\|_K + \|f - \Pi_K f\|_K \right], \quad (28)$$

for all elements $K$ in $\mathcal{T}_h$. Here, $r_h = f - \kappa_K^2 u_h + \Delta u_h$ stands for the residual on $K$. This bound is local and independent of values of $\kappa$ in the other elements and therefore applies to the case considered here.

We emphasize that estimates (26)–(28) are proved in [5] for the case of pure Dirichlet boundary conditions and constant coefficient $\kappa$. However, their proofs remain valid even in the presence of Neumann boundary conditions and due to the condition (8) also for piecewise constant $\kappa$. Nevertheless, estimate (27) is not valid for a facet on the Neumann boundary. We use a slight modification of the proof of Lemma 5(2) from [5] and derive the estimate

$$\|g_K - \boldsymbol{\nabla} u_h|_K \cdot \boldsymbol{n}_K\|_\gamma \le C \left[ \min\{h_K, \kappa_K^{-1}\}^{-\frac{1}{2}} \|u - u_h\|_K \right.$$
$$\left. + \min\{h_K, \kappa_K^{-1}\}^{\frac{1}{2}} \|f - \Pi_K f\|_K + \|g_\mathrm{N} - \Pi_\gamma g_\mathrm{N}\|_\gamma \right] \quad (29)$$

for those facets $\gamma$ of $K$ located on the Neumann boundary $\Gamma_\mathrm{N}$. Thus, defining $R = g_K - \boldsymbol{\nabla} u_h|_K \cdot \boldsymbol{n}_K$, using (26), (27), (29) and the fact that

$$R = g_K - \left\langle \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right\rangle - \frac{1}{2} \left[ \frac{\partial u_h}{\partial \boldsymbol{n}_K} \right]$$

we easily derive the estimate

$$\|R\|_{\partial K} \le C \left( \min\{h_K, \kappa_K^{-1}\}^{-\frac{1}{2}} \|u - u_h\|_{\widetilde{K}} \right.$$
$$\left. + \min\{h_K, \kappa_K^{-1}\}^{\frac{1}{2}} \|f - \Pi f\|_{\widetilde{K}} + \left\| g_\mathrm{N} - \Pi_\gamma^K g_\mathrm{N} \right\|_{\Gamma_\mathrm{N} \cap \partial K} \right) \quad (30)$$

for all $K \in \mathcal{T}_h$. In view of notation (9) and thanks to the assumption (7) the above estimates hold even if $\kappa_K = 0$.

3.5. **Flux reconstruction #1.** For each simplex $K \in \mathcal{T}_h$ on which $\kappa_K \rho_K \leq 1$, we use a reconstruction of the form

$$\boldsymbol{\tau}_K^{(1)} = \boldsymbol{\nabla} u_h|_K + \boldsymbol{\tau}_K^{\mathrm{L}} + \boldsymbol{\tau}_K^{\mathrm{Q}}. \tag{31}$$

The vector field $\boldsymbol{\tau}_K^{\mathrm{L}}$ is defined as

$$\boldsymbol{\tau}_K^{\mathrm{L}} = -\sum_{n=1}^{d+1} \lambda_n \sum_{\substack{m=1 \\ m \neq n}}^{d+1} R_{|\gamma_m}(\boldsymbol{x}_n) \, |\boldsymbol{\nabla}\lambda_m| \, \boldsymbol{t}_{nm}, \tag{32}$$

where $\boldsymbol{x}_n$, $n = 1, 2, \ldots, d+1$, stand for vertices of $K$, $\gamma_n$ are the facets opposite to the vertices $\boldsymbol{x}_n$, $\lambda_n$ are the corresponding barycentric coordinates, $\boldsymbol{t}_{mn} = \boldsymbol{x}_n - \boldsymbol{x}_m$ denote the edge-vectors from $\boldsymbol{x}_m$ to $\boldsymbol{x}_n$, and function $R = g_K - \boldsymbol{\nabla} u_h|_K \cdot \boldsymbol{n}_K$ is affine on each facet of $K$. The quadratic vector field $\boldsymbol{\tau}_K^{\mathrm{Q}}$ is given by

$$\boldsymbol{\tau}_K^{\mathrm{Q}} = \frac{1}{d+1} \sum_{n=1}^{d+1} \sum_{\substack{m=2 \\ m>n}}^{d+1} \lambda_m \lambda_n \boldsymbol{t}_{mn} \boldsymbol{t}_{mn}^T \boldsymbol{\nabla} r(\overline{\boldsymbol{x}}_K) \tag{33}$$

where $r = \Pi_K f - \kappa_K^2 u_h$ is affine on $K$, and $\overline{\boldsymbol{x}}_K$ denotes the centroid of simplex $K$.

It can be easily shown that $\boldsymbol{\tau}_K^{\mathrm{L}} \cdot \boldsymbol{n}_K = R$ and $\boldsymbol{\tau}_K^{\mathrm{Q}} \cdot \boldsymbol{n}_K = 0$ on each facet of $K$. Indeed, if we denote the outward normal unit vector to the facet $\gamma_k$ by $\boldsymbol{n}_k$ then the following identity holds

$$\boldsymbol{\tau}_K^{\mathrm{L}} \cdot \boldsymbol{n}_k|_{\gamma_k} = -\sum_{\substack{n=1 \\ n \neq k}}^{d+1} \lambda_n R_{|\gamma_k}(\boldsymbol{x}_n) \, |\boldsymbol{\nabla}\lambda_k| \, \boldsymbol{t}_{nk} \cdot \boldsymbol{n}_k$$

$$= \sum_{\substack{n=1 \\ n \neq k}}^{d+1} \lambda_n R_{|\gamma_k}(\boldsymbol{x}_n) \, \boldsymbol{t}_{nk} \cdot \boldsymbol{\nabla}\lambda_k = \sum_{\substack{n=1 \\ n \neq k}}^{d+1} \lambda_n R_{|\gamma_k}(\boldsymbol{x}_n) = R_{|\gamma_k},$$

where we use the facts that: $\lambda_k|_{\gamma_k} = 0$; if $n \neq k$ then $\boldsymbol{t}_{nm} \cdot \boldsymbol{n}_k = 0$ for $m \neq k$; $\boldsymbol{n}_k = -\boldsymbol{\nabla}\lambda_k/|\boldsymbol{\nabla}\lambda_k|$; and $\boldsymbol{t}_{nk} \cdot \boldsymbol{\nabla}\lambda_k = 1$. Similarly, we show that

$$\boldsymbol{\tau}_K^{\mathrm{Q}} \cdot \boldsymbol{n}_k|_{\gamma_k} = \frac{1}{d+1} \sum_{n=1}^{d+1} \sum_{\substack{m=2 \\ m>n}}^{d+1} \lambda_m|_{\gamma_k} \lambda_n|_{\gamma_k} (\boldsymbol{t}_{mn} \cdot \boldsymbol{\nabla} r(\overline{\boldsymbol{x}}_K))(\boldsymbol{t}_{mn} \cdot \boldsymbol{n}_k) = 0,$$

because if $n = k$ then $\lambda_n|_{\gamma_k} = 0$ and if $n \neq k$ then $m \neq k$ and $\boldsymbol{t}_{mn} \cdot \boldsymbol{n}_k = 0$.

**Lemma 3.** *Let $K \in \mathcal{T}_h$ then the vector field $\boldsymbol{\tau}_K^{(1)}$ defined by (31), (32), and (33) satisfies $\boldsymbol{\tau}_K^{(1)} \cdot \boldsymbol{n}_K = g_K$ on all facets of $K$ and, if $\kappa_K \rho_K \leq 1$, then*

$$\Pi_K f - \kappa_K^2 u_h + \operatorname{div} \boldsymbol{\tau}_K^{(1)} = 0 \quad \text{in } K.$$

*Proof.* The first assertion is a consequence of the foregoing arguments. Suppose $\kappa_K \rho_K \leq 1$, then since $\boldsymbol{\tau}_K^{\mathrm{L}}$ has constant divergence over the element $K$ and $\boldsymbol{\nabla} u_h|_K$ has vanishing divergence over $K$, we have

$$\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{L}} = \frac{1}{|K|} \int_{\partial K} \boldsymbol{\tau}_K^{\mathrm{L}} \cdot \boldsymbol{n}_K \, \mathrm{d}\boldsymbol{s} = \frac{1}{|K|} \int_{\partial K} g_K \, \mathrm{d}\boldsymbol{s}.$$

Since $\kappa_K \rho_K \leq 1$, the exact equilibration condition (21) is satisfied in $K$ for $n = 1, 2, \ldots, d+1$ and, consequently, using the fact that $\sum_{n=1}^{d+1} \theta_n = 1$ on $K$ in (21) we end up with equality

$$\int_{\partial K} g_K \, \mathrm{d}\boldsymbol{s} = - \int_K (f - \kappa_K^2 u_h) \, \mathrm{d}\boldsymbol{x} = - \int_K (\Pi_K f - \kappa_K^2 u_h) \, \mathrm{d}\boldsymbol{x}.$$

Observing that $\operatorname{div}(\lambda_m \lambda_n \boldsymbol{t}_{mn}) = \lambda_m - \lambda_n$ and that

$$\sum_{n=1}^{d+1} \sum_{\substack{m=2 \\ m>n}}^{d+1} (\lambda_m(\boldsymbol{x}) - \lambda_n(\boldsymbol{x})) (\boldsymbol{x}_n - \boldsymbol{x}_m) = -(d+1)(\boldsymbol{x} - \overline{\boldsymbol{x}}_K),$$

where $\boldsymbol{x} = \sum_{m=1}^{d+1} \boldsymbol{x}_m \lambda_m(\boldsymbol{x})$ and $\overline{\boldsymbol{x}}_K = \sum_{m=1}^{d+1} \boldsymbol{x}_m/(d+1)$, we can compute the divergence of $\boldsymbol{\tau}_K^{\mathrm{Q}}$ as

$$\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{Q}}(\boldsymbol{x}) = \frac{1}{d+1} \sum_{n=1}^{d+1} \sum_{\substack{m=2 \\ m>n}}^{d+1} (\lambda_m(\boldsymbol{x}) - \lambda_n(\boldsymbol{x}))(\boldsymbol{x}_n - \boldsymbol{x}_m) \cdot \boldsymbol{\nabla} r(\overline{\boldsymbol{x}}_K)$$
$$= (\overline{\boldsymbol{x}}_K - \boldsymbol{x}) \cdot \boldsymbol{\nabla} r(\overline{\boldsymbol{x}}_K).$$

Using the fact that $r = \Pi_K f - \kappa_K^2 u_h$ is affine and the centroid quadrature rule for simplices that is exact for all linear functions, we obtain

$$\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{Q}}(\boldsymbol{x}) = -r(\boldsymbol{x}) + r(\overline{\boldsymbol{x}}_K) = -r(\boldsymbol{x}) + \frac{1}{|K|} \int_K r \, \mathrm{d}\boldsymbol{x}.$$

The statement of the lemma follows by summing the above equations. $\qquad \square$

The next result shows that $\boldsymbol{\tau}_K^{(1)}$ gives and efficient estimate of the local error in element $K$:

**Lemma 4.** *If $K \in \mathcal{T}_h$ is such that $\kappa_K \rho_K \leq 1$ then*

$$\eta_K\left(\boldsymbol{\tau}_K^{(1)}\right) \leq C \left( \|u - u_h\|_{\widetilde{K}} + h_K \|f - \Pi f\|_{\widetilde{K}} + h_K^{1/2} \left\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\right\|_{\Gamma_{\mathrm{N}} \cap \partial K} \right).$$

*Proof.* Let

$$\boldsymbol{c}_n = \sum_{\substack{m=1 \\ m \neq n}}^{d+1} R_{|\gamma_m}(\boldsymbol{x}_n) \left|\boldsymbol{\nabla} \lambda_m\right| \boldsymbol{t}_{nm},$$

then $\boldsymbol{\tau}_K^{\mathrm{L}} = -\sum_{n=1}^{d+1} \lambda_n \boldsymbol{c}_n$ and we have

$$|\boldsymbol{c}_n| \leq \sum_{\substack{m=1 \\ m \neq n}}^{d+1} |R_{|\gamma_m}(\boldsymbol{x}_n)|\, |\boldsymbol{\nabla}\lambda_m|\, |\boldsymbol{t}_{nm}| \leq C \sum_{\substack{m=1 \\ m \neq n}}^{d+1} |R_{|\gamma_m}(\boldsymbol{x}_n)|,$$

because $d|K||\boldsymbol{\nabla}\lambda_m| = |\gamma_m|$, $|\boldsymbol{t}_{nm}| \leq h_K$, and due to the shape regularity assumption (3). Further, thanks to the linearity of $R$,

$$\sum_{\substack{n=1 \\ n \neq m}}^{d+1} |R_{|\gamma_m}(\boldsymbol{x}_n)|^2 \leq C \frac{1}{|\gamma_m|} \|R\|_{\gamma_m}^2 \,.$$

We utilize these results to bound

$$\left\|\boldsymbol{\tau}_K^{\mathrm{L}}\right\|_K^2 \leq C|K| \sum_{n=1}^{d+1} |\boldsymbol{c}_n|^2 \leq C|K| \sum_{m=1}^{d+1} \sum_{\substack{n=1 \\ n \neq m}}^{d+1} |R_{|\gamma_m}(\boldsymbol{x}_n)|^2 \leq C h_K \|R\|_{\partial K}^2 \,. \tag{34}$$

Similarly, we have

$$\left\|\boldsymbol{\tau}_K^{\mathrm{Q}}\right\|_K^2 \leq C \sum_{n=1}^{d+1} \sum_{\substack{m=2 \\ m>n}}^{d+1} \int_K \lambda_m^2 \lambda_n^2 \,\mathrm{d}\boldsymbol{x} |\boldsymbol{t}_{mn}|^4 |\boldsymbol{\nabla}r|^2 \leq C h_K^4 |K|\, |\boldsymbol{\nabla}r|^2 \leq C h_K^4 \|\boldsymbol{\nabla}r\|_K^2 \,,$$

where we used the fact that $\boldsymbol{\nabla}r$ is constant over $K$. Since $r \in \mathbb{P}^1(K)$, we have the inverse inequality $\|\boldsymbol{\nabla}r\|_K \leq C h_K^{-1} \|r\|_K$ and we obtain

$$\left\|\boldsymbol{\tau}_K^{\mathrm{Q}}\right\|_K \leq C h_K \|r\|_K = C h_K \|\Pi_K r_h\|_K \,, \tag{35}$$

because $r_h = f - \kappa_K^2 u_h + \Delta u_h$ and $\Pi_K r_h = r$ on $K$.

Finally, using estimate (30) in (34) and estimate (28) in (35), we derive

$$\left\|\boldsymbol{\tau}_K^{\mathrm{L}}\right\|_K \leq C \left[\|u - u_h\|_{\widetilde{K}} + h_K \|f - \Pi f\|_{\widetilde{K}} + h_K^{1/2} \left\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\right\|_{\Gamma_{\mathrm{N}} \cap \partial K}\right], \tag{36}$$

$$\left\|\boldsymbol{\tau}_K^{\mathrm{Q}}\right\|_K \leq C \left[\|u - u_h\|_K + h_K \|f - \Pi f\|_K\right]. \tag{37}$$

Notice that the assumption $\kappa_K \rho_K \leq 1$ and the shape regularity (3) imply the existence of a constant $C > 0$ such that $C h_K \leq \min\{h_K, \kappa_K^{-1}\} \leq h_K$. Due to Lemma 3 and definition (15) we have

$$\eta_K(\boldsymbol{\tau}_K^{(1)}) = \left\|\boldsymbol{\tau}_K^{(1)} - \boldsymbol{\nabla}u_h\right\|_K = \left\|\boldsymbol{\tau}_K^{\mathrm{L}} + \boldsymbol{\tau}_K^{\mathrm{Q}}\right\|_K \leq \left\|\boldsymbol{\tau}_K^{\mathrm{L}}\right\|_K + \left\|\boldsymbol{\tau}_K^{\mathrm{Q}}\right\|_K$$

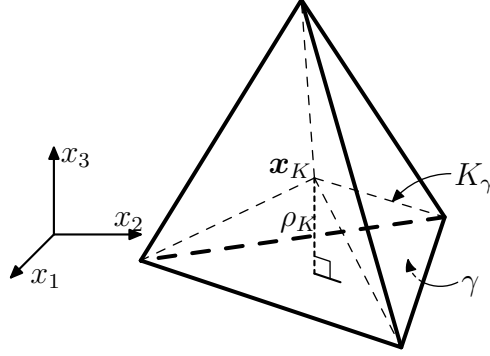and estimates (36)–(37) finish the proof. $\qquad\square$

FIGURE 2. Division of $K$ into subsimplices $K_\gamma$ and the local Cartesian coordinates.

3.6. **Flux reconstruction #2.** On elements $K$ for which $\kappa_K \rho_K > 1$, we use a flux reconstruction given by

$$\boldsymbol{\tau}_K^{(2)} = \boldsymbol{\nabla} u_h|_K + \boldsymbol{\tau}_K^{\mathrm{O}}, \tag{38}$$

where the vector field $\boldsymbol{\tau}_K^{\mathrm{O}}$ is defined piecewise on each element $K$ in the following way. We consider $d+1$ subsimplices $K_\gamma$ of $K$ that are defined as convex hulls of the incentre $\boldsymbol{x}_K$ and facets $\gamma$ of $K$. In each subsimplex $K_\gamma$ we define $\boldsymbol{\tau}_K^{\mathrm{O}}$ to be

$$\boldsymbol{\tau}_K^{\mathrm{O}}(\boldsymbol{x}) = \rho_K^{-1}(1 - \kappa_K x_d)^+(\boldsymbol{x} - \boldsymbol{x}_K)R(x_1, \ldots, x_{d-1}) \quad \text{in } K_\gamma,$$

where $z^+ = (|z| + z)/2$ stands for the positive part of $z$, $R = g_K - \boldsymbol{\nabla} u_h|_K \cdot \boldsymbol{n}_K$ as before, $\rho_K$ is the inradius of $K$, and $\boldsymbol{x} = (x_1, x_2, \ldots, x_d)$ are local Cartesian coordinates defined in such a way that points $(x_1, \ldots, x_{d-1})$ lie in the plane of $\gamma$ and $x_d$ corresponds to the direction perpendicular to $\gamma$ aiming inwards $K$, see Figure 2 for a three-dimensional illustration.

Clearly, $\boldsymbol{\tau}_K^{\mathrm{O}}$ vanishes for $x_d \geq \kappa_K^{-1}$. We also observe that the normal component of $\boldsymbol{\tau}_K^{\mathrm{O}}$ vanishes on $\partial K_\gamma \setminus \gamma$ whilst on $\gamma$

$$\boldsymbol{\tau}_K^{\mathrm{O}} \cdot \boldsymbol{n}_K|_\gamma = R(x_1, \ldots, x_{d-1}),$$

because $x_d = 0$ on $\gamma$ and $(\boldsymbol{x} - \boldsymbol{x}_K) \cdot \boldsymbol{n}_K = \rho_K$ on $\gamma$. These conditions guarantee that $\boldsymbol{\tau}_K^{(2)} \in \boldsymbol{H}(\mathrm{div}, K)$ and $\boldsymbol{\tau}_K^{(2)} \cdot \boldsymbol{n}_K = g_K$ on $\partial K$. Moreover, the flux leads to a locally efficient estimator for the error in the case $\kappa_K \rho_K > 1$:

**Lemma 5.** Let $K \in \mathcal{T}_h$ then $\boldsymbol{\tau}_K^{(2)} \cdot \boldsymbol{n}_K = g_K$ on $\partial K$ and, if $\kappa_K \rho_K > 1$, then

$$\eta_K\big(\boldsymbol{\tau}_K^{(2)}\big) \leq C \left( \|u - u_h\|_{\widetilde{K}} + \kappa_K^{-1} \|f - \Pi f\|_{\widetilde{K}} + \kappa_K^{-1/2} \big\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\big\|_{\Gamma_{\mathrm{N}} \cap \partial K} \right).$$

*Proof.* The first assertion has been shown above. Suppose $\kappa_K \rho_K > 1$, then since $|\boldsymbol{x} - \boldsymbol{x}_K| \leq h_K$ and $h_K/\rho_K$ is bounded uniformly thanks to (3), we obtain

$$\big\|\boldsymbol{\tau}_K^{\mathrm{O}}\big\|_{K_\gamma}^2 \leq \frac{h_K^2}{\rho_K^2} \int_0^{\kappa_K^{-1}} (1 - \kappa_K x_d)^2 \, \mathrm{d}x_d \|R\|_\gamma^2 = \frac{h_K^2}{\rho_K^2} \frac{1}{3\kappa_K} \|R\|_\gamma^2 \leq C\kappa_K^{-1} \|R\|_\gamma^2. \tag{39}$$

For $x_d \leq \kappa_K^{-1}$, a simple computation yields inequality

$$|\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{O}}| \leq \rho_K^{-1}(1 - \kappa_K x_d)^+ (d|R| + h_K|\boldsymbol{\nabla}_\gamma R|) + \rho_K^{-1}\kappa_K h_K|R|,$$

where $\boldsymbol{\nabla}_\gamma$ denotes the gradient with respect to $x_1, \ldots, x_{d-1}$ only. Consequently,

$$\left\|\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{O}}\right\|_{K_\gamma}^2 \leq \frac{C}{\rho_K^2} \left( \int_0^{\kappa_K^{-1}} (1 - \kappa_K x_d)^2 \, \mathrm{d}x_d \left[ \|R\|_\gamma^2 + h_K^2 \|\boldsymbol{\nabla}_\gamma R\|_\gamma^2 \right] + \kappa_K h_K^2 \|R\|_\gamma^2 \right).$$

Since $R \in \mathbb{P}^1(\gamma)$, we use the shape regularity and the inverse estimate $\|\boldsymbol{\nabla}_\gamma R\|_\gamma \leq C h_\gamma^{-1} \|R\|_\gamma$ to derive

$$\left\|\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{O}}\right\|_{K_\gamma}^2 \leq \frac{C}{\rho_K^2} \frac{1}{\kappa_K} \max\{1, \kappa_K h_K\}^2 \|R\|_\gamma^2 \leq C \kappa_K \|R\|_\gamma^2, \tag{40}$$

where the last inequality follows from the shape regularity (3), from (5), and from the assumption $\kappa_K \rho_K > 1$.

Hence, thanks to (30) and (39):

$$\left\|\boldsymbol{\tau}_K^{(2)} - \boldsymbol{\nabla} u_h\right\|_K = \left\|\boldsymbol{\tau}_K^{\mathrm{O}}\right\|_K \leq C \kappa_K^{-1/2} \|R\|_{\partial K}$$
$$\leq C \left( \|u - u_h\|_{\widetilde{K}} + \kappa_K^{-1} \|f - \Pi f\|_{\widetilde{K}} + \kappa_K^{-1/2} \left\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\right\|_{\Gamma_{\mathrm{N}} \cap \partial K} \right).$$

Similarly, estimates (40), (28), and (30) yield the bound

$$\kappa_K^{-1} \left\|\Pi_K f - \kappa_K^2 u_h + \operatorname{div} \boldsymbol{\tau}_K^{(2)}\right\|_K \leq \kappa_K^{-1} \|\Pi_K r_h\|_K + \kappa_K^{-1} \left\|\operatorname{div} \boldsymbol{\tau}_K^{\mathrm{O}}\right\|_K$$
$$\leq \kappa_K^{-1} \|\Pi_K r_h\|_K + C \kappa_K^{-1/2} \|R\|_{\partial K}$$
$$\leq C \left( \|u - u_h\|_{\widetilde{K}} + \kappa_K^{-1} \|f - \Pi f\|_{\widetilde{K}} + \kappa_K^{-1/2} \left\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\right\|_{\Gamma_{\mathrm{N}} \cap \partial K} \right).$$

Combining these estimates gives the result claimed. $\qquad \square$

3.7. **Main result.** We combine the flux reconstructions $\boldsymbol{\tau}_K^{(1)}$ and $\boldsymbol{\tau}_K^{(2)}$ in a natural way and construct $\boldsymbol{\tau} \in \boldsymbol{H}(\operatorname{div}, \Omega)$ elementwise as

$$\boldsymbol{\tau}|_K = \begin{cases} \boldsymbol{\tau}_K^{(1)} & \text{if } \kappa_K \rho_K \leq 1, \\ \boldsymbol{\tau}_K^{(2)} & \text{if } \kappa_K \rho_K > 1, \end{cases} \tag{41}$$

where $\boldsymbol{\tau}_K^{(1)}$ and $\boldsymbol{\tau}_K^{(2)}$ are defined in (31) and (38). The following theorem shows that the associated error estimator provides a guaranteed upper bound on the error, which is robust with respect to $\kappa$ and $h$.

**Theorem 6.** *Let $u$ be the exact weak solution given by (2) and and $u_h \in V_h$ be its finite element approximation (10). Let the flux reconstruction $\boldsymbol{\tau} \in \boldsymbol{H}(\operatorname{div}, \Omega)$ be given by (41). Then the error in $u_h$ is bounded by*

$$\|u - u_h\|^2 \leq \eta^2(\boldsymbol{\tau}) = \sum_{K \in \mathcal{T}_h} [\eta_K(\boldsymbol{\tau}) + \operatorname{osc}_K(f) + \operatorname{osc}_{\Gamma_{\mathrm{N}} \cap \partial K}(g_{\mathrm{N}})]^2 .$$

*Moreover, there exists a positive constant $C$, independent of any mesh-size or any values $\kappa_K$ satisfying (6)–(7), such that*

$$\eta_K(\boldsymbol{\tau}) \leq C\Big( \|\!|u - u_h|\!\|_{\widetilde{K}} + \min\{h_K, \kappa_K^{-1}\} \|f - \Pi f\|_{\widetilde{K}}$$

$$+ \min\{h_K, \kappa_K^{-1}\}^{1/2} \left\|g_{\mathrm{N}} - \Pi_\gamma^K g_{\mathrm{N}}\right\|_{\Gamma_{\mathrm{N}} \cap \partial K} \Big).$$

*Proof.* It follows immediately from Lemmas 2, 4, and 5. $\qquad\square$

In view of convention (9), this result holds even if $\kappa_K = 0$ for any number of elements $K \in \mathcal{T}_h$. Theorem 6 provides a robust, computable upper bound, but it is possible to improve the bound at the expense of having to compute both $\eta_K(\boldsymbol{\tau}_K^{(1)})$ and $\eta_K(\boldsymbol{\tau}_K^{(2)})$ on every element. The associated flux is defined by

$$\boldsymbol{\tau}^*|_K = \begin{cases} \boldsymbol{\tau}_K^{(1)} & \text{if } \kappa_K = 0 \text{ or if } \eta_K(\boldsymbol{\tau}_K^{(1)}) \leq \eta_K(\boldsymbol{\tau}_K^{(2)}), \\ \boldsymbol{\tau}_K^{(2)} & \text{otherwise.} \end{cases} \tag{42}$$

and the corresponding estimator is given by $\eta(\boldsymbol{\tau}^*)$, which in turn involves the local indicator $\eta_K(\boldsymbol{\tau}^*) = \min\big\{\eta_K(\boldsymbol{\tau}_K^{(1)}), \eta_K(\boldsymbol{\tau}_K^{(2)})\big\}$. This flux reconstruction is slightly more expensive to compute, but it yields more accurate estimator than $\boldsymbol{\tau}$, because $\eta_K(\boldsymbol{\tau}^*) \leq \eta_K(\boldsymbol{\tau})$. Clearly, if we replace $\boldsymbol{\tau}$ by $\boldsymbol{\tau}^*$ in Theorem 6, both its statements remain valid.

## 4. NUMERICAL EXAMPLE

This section illustrates numerical performance of the a posteriori error estimators $\eta(\boldsymbol{\tau})$ and $\eta(\boldsymbol{\tau}^*)$ for a three dimensional example. In particular, the example confirms the robustness of both estimators with respect to the discontinuous reaction coefficient $\kappa$ and with respect to the mesh size.

We consider problem (1) in a cube $\Omega = (-1, 1)^3$, with piecewise constant coefficient $\kappa$ defined by

$$\kappa(x_1, x_2, x_3) = \begin{cases} \kappa_1 & \text{for } x_1 < 0, \\ \kappa_2 & \text{for } x_1 \geq 0, \end{cases}$$

where $0 < \kappa_1 \leq \kappa_2$ are constants. The right-hand side is $f = \kappa_1^2$. Homogeneous Dirichlet boundary conditions are assumed on $\Gamma_{\mathrm{D}} = \{(x_1, x_2, x_3) \in \partial\Omega : x_1 = \pm 1\}$ and homogeneous Neumann boundary conditions are prescribed on $\Gamma_{\mathrm{N}} = \partial\Omega \backslash \Gamma_{\mathrm{D}}$.

Its exact solution can be expressed as

$$u(x_1, x_2, x_3) = \begin{cases} A_1 \mathrm{e}^{-\kappa_1 x_1} + A_2 \mathrm{e}^{\kappa_1 x_1} + 1 & \text{for } x_1 < 0, \\ A_3 \mathrm{e}^{-\kappa_2 x_1} + A_4 \mathrm{e}^{\kappa_2 x_1} + \kappa_1^2/\kappa_2^2 & \text{for } x_1 \geq 0, \end{cases}$$

where constants $A_1, \ldots, A_4$ are uniquely determined by the Dirichlet boundary conditions and by the requirement of $C^1$ continuity of $u(x_1, x_2, x_3)$ for $x_1 = 0$. In the subsequent computations we fix $\kappa_2 = 10^6$ and hence the solution has a boundary layer at least in the vicinity of the face $x_1 = 1$. Although the true solution has a univariate nature, this plays no role in the computations.
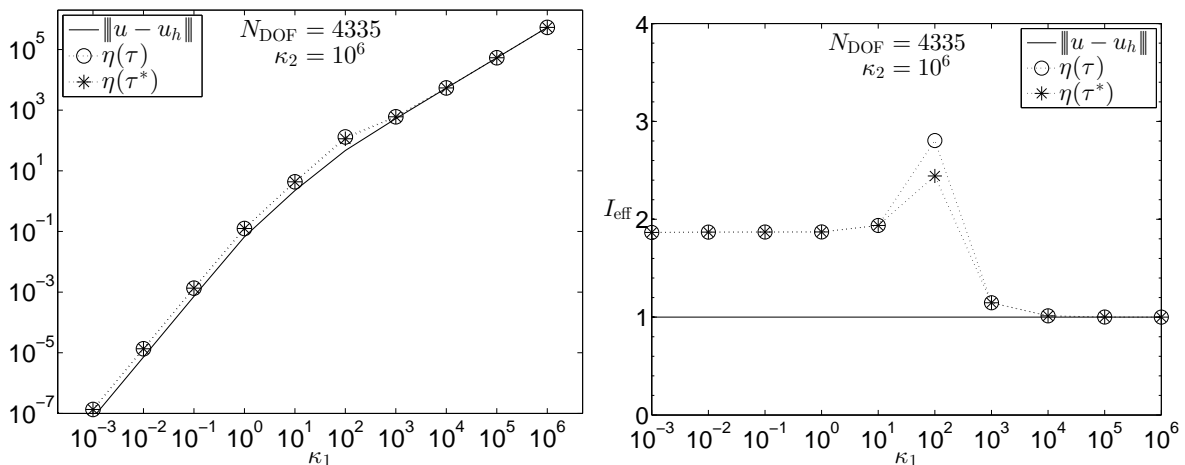
FIGURE 3. Dependence of $\|u - u_h\|$, $\eta(\boldsymbol{\tau})$, and $\eta(\boldsymbol{\tau}^*)$ on $\kappa_1$ (left) and corresponding effectivity indices (right). These results correspond to $\kappa_2 = 10^6$ and to a mesh with $N_{\mathrm{DOF}} = 4335$ ($M = 16$).

We approximate this problem using linear finite elements on uniform tetrahedral meshes that are constructed in two steps. First, the cube $\Omega$ is uniformly divided into $M^3$ subcubes and then each subcube is split into 6 tetrahedrons along its diagonal. The resulting mesh then has $N_{\mathrm{DOF}} = (M - 1)(M + 1)^2$ degrees of freedom.

Figure 3 presents the results for a fixed mesh ($M = 16$, $N_{\mathrm{DOF}} = 4335$). The left panel shows the dependence of the true error $\|u - u_h\|$ and the error estimators $\eta(\boldsymbol{\tau})$ and $\eta(\boldsymbol{\tau}^*)$ as $\kappa_1$ is varied in the range $(0, \kappa_2]$. The right panel presents the effectivity indices $I_{\mathrm{eff}} = \eta/\|u - u_h\|$. We observe that both estimators provide upper bound on the error and that they robustly capture the behaviour of the error in the whole range of values of $\kappa_1$. Thus, they are independent of the ratio $\kappa_1/\kappa_2$ in this case. As expected, the effectivity index for $\eta(\boldsymbol{\tau}^*)$ is smaller than for $\eta(\boldsymbol{\tau})$. Both indices exhibit values around 2 for small values of $\kappa_1$ and they are close to 1 for $\kappa_1 \geq 1000$.

Similarly, Figure 4 demonstrates the behaviour of these error estimators and of the true error with respect to the number of degrees of freedom. In this case we fix $\kappa_1 = 100$ and solve the problem on a series of meshes with $M = 2, 2^2, 2^3, \ldots, 2^7$. We have chosen the most unfavourable value $\kappa_1 = 100$ for which both error estimators exhibit the highest overestimation in Figure 3.

As above, the left panel of Figure 4 presents the values of the true error $\|u - u_h\|$ and of the estimators $\eta(\boldsymbol{\tau})$ and $\eta(\boldsymbol{\tau}^*)$, while the right panel shows the effectivity indices. Again, we verify the upper bound property of the estimators and observe their robust behaviour with respect to the mesh size. The effectivity indices have values around 1 and 2 with an exception of the intermediate case, where the mesh size $h$ is comparable to $1/\kappa_1$.
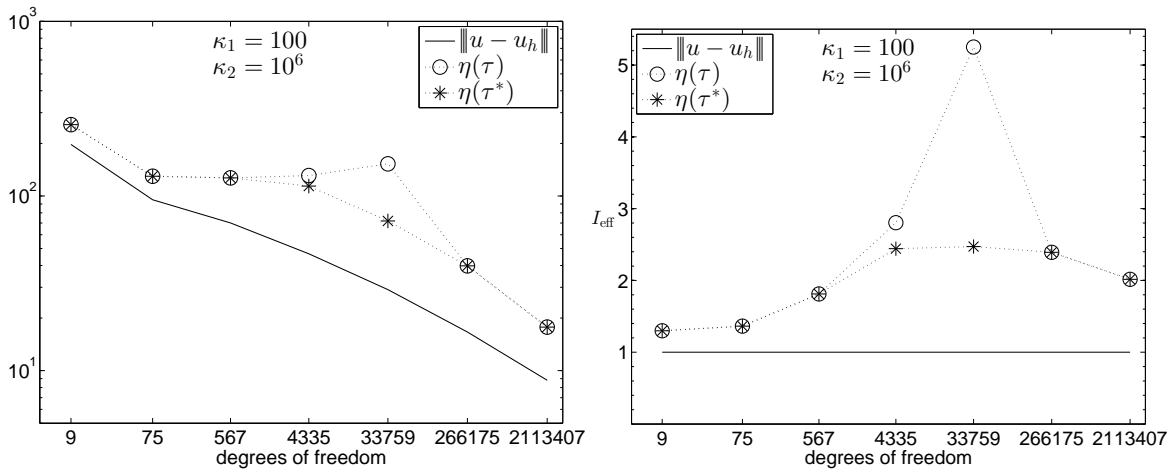
FIGURE 4. Dependence of $\|u - u_h\|$, $\eta(\boldsymbol{\tau})$, and $\eta(\boldsymbol{\tau}^*)$ on the number of degrees of freedom (left) and corresponding effectivity indices (right). These results were computed on a sequence of uniformly refined meshes with $\kappa_1 = 100$ and $\kappa_2 = 10^6$.

## 5. Conclusions

We presented a robust a posteriori error estimator on the energy norm of the approximation error for a reaction-diffusion problem in arbitrary dimension. The reaction coefficient $\kappa$ is assumed to be piecewise constant and mixed Dirichlet-Neumann boundary conditions are allowed. The estimator is robust with respect to the reaction coefficient $\kappa$, including the singularly perturbed case, and it provides a computable upper bound on the error. The upper bound is guaranteed up to round-off errors and quadrature errors in the evaluation of $\eta(\boldsymbol{\tau})$.

The approach is suitable for the piecewise linear finite element approximations. The Galerkin condition (10) is required in order to guarantee the exact equilibration condition (21) in case of small values of the reaction coefficient $\kappa$. On the other hand, the exact equilibration is not needed for large values of $\kappa$ and the presented error estimator can be used for an arbitrary (conforming) approximate solution $u_h \in V$.

Finally, we note that whilst we have assumed conformity of the approximation, this is not essential. Methodologies derived in [23, 24, 25] could be used to extend the error bound to any piecewise linear non-conforming approximation.

## References

### References

[1] M. Ainsworth, J. T. Oden, A posteriori error estimation in finite element analysis, Wiley, New York, 2000.

[2] M. Ainsworth, T. Vejchodský, Fully computable robust a posteriori error bounds for singularly perturbed reaction–diffusion problems, Numer. Math. 119 (2) (2011) 219–243.

[3] W. Dörfler, A convergent adaptive algorithm for Poisson's equation, SIAM J. Numer. Anal. 33 (3) (1996) 1106–1124.

[4] R. Verfürth, A posteriori error estimators for convection-diffusion equations, Numer. Math. 80 (4) (1998) 641–663.

[5] M. Ainsworth, I. Babuška, Reliable and robust a posteriori error estimating for singularly perturbed reaction-diffusion problems, SIAM J. Numer. Anal. 36 (2) (1999) 331–353 (electronic).

[6] S. Grosman, An equilibrated residual method with a computable error approximation for a singularly perturbed reaction-diffusion problem on anisotropic finite element meshes, M2AN Math. Model. Numer. Anal. 40 (2) (2006) 239–267.

[7] I. Cheddadi, R. Fučík, M. I. Prieto, M. Vohralík, Guaranteed and robust a posteriori error estimates for singularly perturbed reaction–diffusion problems, M2AN Math. Model. Numer. Anal. 43 (2009) 867–888.

[8] T. Linss, A posteriori error estimation for arbitrary-order FEM applied to singularly perturbed one-dimensional reaction-diffusion problems, Appl. Math. To appear, 2014.

[9] J. L. Synge, The hypercircle in mathematical physics: a method for the approximate solution of boundary value problems, Cambridge University Press, New York, 1957.

[10] J. P. Aubin, H. G. Burchard, Some aspects of the method of the hypercircle applied to elliptic variational problems, in: Numerical Solution of Partial Differential Equations, II (SYNSPADE 1970) (Proc. Sympos., Univ. of Maryland, College Park, Md., 1970), Academic Press, New York, 1971, pp. 1–67.

[11] J. Haslinger, I. Hlaváček, Convergence of a finite element method based on the dual variational formulation, Apl. Mat. 21 (1) (1976) 43–65.

[12] B. F. de Veubeke, Displacement and equilibrium models in the finite element method, in: O. Zienkiewicz, G. Hollister (Eds.), Stress Analysis, Wiley, London, 1965, pp. 145–197.

[13] S. Repin, A posteriori estimates for partial differential equations, Vol. 4 of Radon Series on Computational and Applied Mathematics, Walter de Gruyter GmbH & Co. KG, Berlin, 2008.

[14] D. Braess, J. Schöberl, Equilibrated residual error estimator for edge elements, Math. Comp. 77 (262) (2008) 651–672.

[15] Z. Cai, S. Zhang, Flux recovery and a posteriori error estimators: conforming elements for scalar elliptic equations, SIAM J. Numer. Anal. 48 (2) (2010) 578–602.

[16] P. Jiránek, Z. Strakoš, M. Vohralík, A posteriori error estimates including algebraic error and stopping criteria for iterative solvers, SIAM J. Sci. Comput. 32 (3) (2010) 1567–1590.

[17] D. W. Kelly, The self-equilibration of residuals and complementary a posteriori error estimates in the finite element method, Internat. J. Numer. Methods Engrg. 20 (8) (1984) 1491–1506.

[18] P. Ladevèze, D. Leguillon, Error estimate procedure in the finite element method and applications, SIAM J. Numer. Anal. 20 (3) (1983) 485–509.

[19] N. Parés, H. Santos, P. Díez, Guaranteed energy error bounds for the Poisson equation using a flux-free approach: Solving the local problems in subdomains, Internat. J. Numer. Methods Engrg. 79 (10) (2009) 1203–1244.

[20] M. Vohralík, Guaranteed and fully robust a posteriori error estimates for conforming discretizations of diffusion problems with discontinuous coefficients, J. Sci. Comput. 46 (3) (2011) 397–438.

[21] L. E. Payne, H. F. Weinberger, An optimal Poincaré inequality for convex domains, Arch. Rational Mech. Anal. 5 (1960) 286–292 (1960).

[22] I. Šebestová, T. Vejchodský, Two-sided bounds for eigenvalues of differential operators with applications to Friedrichs, Poincaré, trace, and similar constants, SIAM J. Numer. Anal. To appear.

[23] M. Ainsworth, Robust a posteriori error estimation for nonconforming finite element approximation, SIAM J. Numer. Anal. 42 (6) (2005) 2320–2341 (electronic).

[24] M. Ainsworth, A posteriori error estimation for discontinuous Galerkin finite element approximation, SIAM J. Numer. Anal. 45 (4) (2007) 1777–1798 (electronic).

[25] A. Ern, A. F. Stephansen, M. Vohralík, Guaranteed and robust discontinuous Galerkin a posteriori error estimates for convection-diffusion-reaction problems, J. Comput. Appl. Math. 234 (1) (2010) 114–130.

MARK AINSWORTH, DIVISION OF APPLIED MATHEMATICS, BROWN UNIVERSITY, 182 GEORGE STREET PROVIDENCE, RI 02912, USA

*E-mail address*: mark_ainsworth@brown.edu

TOMÁŠ VEJCHODSKÝ, MATHEMATICAL INSTITUTE, UNIVERSITY OF OXFORD, ANDREW WILES BUILDING, RADCLIFFE OBSERVATORY QUARTER, WOODSTOCK ROAD, OXFORD, OX2 6GG, UK, AND INSTITUTE OF MATHEMATICS, ACADEMY OF SCIENCES, ŽITNÁ 25, CZ-115 67 PRAGUE 1, CZECH REPUBLIC.

*E-mail address*: vejchod@math.cas.cz