# ROBUST LOCAL PROBLEM ERROR ESTIMATION FOR A SINGULARLY PERTURBED PROBLEM ON ANISOTROPIC FINITE ELEMENT MESHES

## Gerd Kunert[1]

**Abstract.** Singularly perturbed problems often yield solutions with strong directional features, *e.g.* with boundary layers. Such anisotropic solutions lend themselves to adapted, anisotropic discretizations. The quality of the corresponding numerical solution is a key issue in any computational simulation.

To this end we present a new robust error estimator for a singularly perturbed reaction–diffusion problem. In contrast to conventional estimators, our proposal is suitable for *anisotropic* finite element meshes. The estimator is based on the solution of a local problem, and yields error bounds uniformly in the small perturbation parameter. The error estimation is efficient, *i.e.* a lower error bound holds. The error estimator is also reliable, *i.e.* an upper error bound holds, provided that the anisotropic mesh discretizes the problem sufficiently well.

A numerical example supports the analysis of our anisotropic error estimator.

## 1. Introduction

Adaptive algorithms form nowadays an indispensable tool for most finite element simulations. They basically consist of the ingredients *Solve – Estimate error – Refine mesh* which are repeated until the desired accuracy is achieved, see also [2, 25]. The present work is part of a series of related endeavors in a particular field of finite element analysis. While standard finite element meshes employ *isotropic* (or shape regular) elements, we investigate so–called *anisotropic* (or stretched) elements. They are characterized by a large stretching ratio (also called aspect ratio). Equivalently, the ratio of the diameters of the circumscribed and inscribed spheres can be arbitrarily large. Such anisotropic meshes are particularly useful when the differential equation gives rise to a solution with strong directional features, such as boundary layers or interior layers. Application of anisotropic meshes as well as theoretical investigations can be found for example in [4, 5, 20, 21, 27, 28] and [11, 13, 14, 19, 23].

It is a natural desire to incorporate anisotropic meshes into adaptive algorithms. Clearly, additional ingredients are required then, namely *anisotropic information extraction* (*e.g.* find the (quasi) optimal stretching direction and stretching ratio of the anisotropic elements), and *anisotropic mesh refinement*. Less obvious but equally important is the *error estimation* part. Unfortunately most of the conventional *a posteriori* error estimators for isotropic meshes fail when applied on anisotropic meshes. Therefore the derivation and analysis of estimators which are suitable for anisotropic elements is of vital importance for any adaptive anisotropic algorithm.

[1] TU Chemnitz, Fakultät für Mathematik, 09107 Chemnitz, Germany. e-mail: `kunert@mathematik.tu-chemnitz.de`

Fortunately this challenging venture has seen some success recently [11,13,14,16,19,23]. For the Poisson model problem it has been shown that anisotropic error estimation is possible, and the methodology and analytical tools have been developed, proposed and refined. Now anisotropic error estimation has to prove its potential for more realistic settings. Singularly perturbed problems offer ideal test fields since they often induce boundary layers where anisotropic elements can be employed favourably.

From now on let us consider a singularly perturbed reaction–diffusion model problem, see (1) below, which usually gives rise to boundary layers whenever a non–vanishing right hand side meets homogeneous Dirichlet boundary data.

Although (1) forms a comparatively simple model problem, the knowledge of *robust* error estimators has been unsatisfactory for a long time. The first estimators with error bounds that are uniform in the small perturbation parameter $\varepsilon$ were due to Angermann [3], Verfürth [26] and Ainsworth and Babuška [1]; all of them considered isotropic meshes. Angermann measures the error in a somewhat strange norm (which seems to be mainly of theoretical interest) whereas Verfürth and Ainsworth and Babuška concentrate on the energy norm (which is the most natural norm). For *anisotropic* meshes Kunert [18] recently succeeded in deriving a robust residual error estimator, also for the energy norm. As a corollary of that work a further estimator has been derived and included in [19].

In our present work we propose a new error estimator for the singularly perturbed reaction–diffusion problem (1) which is suitable for anisotropic meshes, and that is based on the solution of a local problem. The roots of this local problem error estimator are twofold. Firstly it relies on the anisotropic residual error estimator [18] whose results are partly the foundation for the present analysis. Secondly we utilize the methodology of local problem error estimation. For the Poisson problem (on isotropic meshes) this is fairly well understood, see *e.g.* the exposition in [25]; hence the general *framework* of the proofs can be derived relatively easily. The precise definition and analysis of our estimator, however, are much more difficult and technical. This concerns for example the choice of the local problem, the careful calibration of all ingredients, or certain equivalence lemmas. Although we could exploit some experience from anisotropic local problem error estimation for the Poisson problem [16], the "extension" to the singularly perturbed problem requires several new ingredients and is by no means straight–forward. Note that in [1] also a local problem error estimator is derived (for isotropic meshes). However the local problem there is infinite dimensional whereas our proposal here involves an (at most) five dimensional local space.

When comparing with the anisotropic residual error estimator, our newly proposed local problem error estimator is certainly more expensive since a local problem has to be computed and solved. Nonetheless the disadvantage of any residual based estimator is that the proof of the error bound is based on several intermediate steps, such as interpolation estimates and the Cauchy Schwarz inequality. In contrast to this the local problem error estimator requires less auxiliary steps, and thus contains less constants (which are unknown in general). This can also be observed numerically where the qualitative behaviour of both error estimators is comparable but the local problem error estimation is much closer to the true error.

Finally note that all known anisotropic error estimators require that the anisotropy of the mesh and the anisotropy of the solution correspond sufficiently well. As in previous work, this correspondence is measured by a so–called matching function which is explained in our exposition.

The remainder of the paper is organized as follows. After presenting the model problem in Section 2, we repeat in Section 3 some notation, basic tools and lemmas that have been applied successfully in previous anisotropic investigations. The transformation technique is of particular importance, and several specific bubble functions play a major role. Furthermore the residual error estimator of [18] is recalled for self–containment. Next, Section 4 is devoted to the local problem error estimator and its analysis. Reliable upper and lower error bounds are proven and a stable basis of the local problem is presented. Additionally a further, face oriented local problem error estimator is given. Computational aspects are discussed in Section 5, and the numerical experiments of Section 6 confirm the analysis. The summary in Section 7 and a technical proof in Appendix A conclude this work.

## 2. The model problem and its discretization

Our focus is on a singularly perturbed reaction–diffusion model problem with Dirichlet–Neumann boundary conditions

$$
\begin{aligned}
-\varepsilon\Delta u + u &= f \quad \text{in } \Omega \\
u &= 0 \quad \text{on } \Gamma_{\mathrm{D}} \\
\varepsilon \cdot \partial u/\partial n &= g \quad \text{on } \Gamma_{\mathrm{N}}
\end{aligned} \tag{1}
$$

in a bounded, polyhedral domain $\Omega \subset \mathbb{R}^d$, $d = 2, 3$, with boundary $\partial\Omega = \Gamma_{\mathrm{D}} \cup \Gamma_{\mathrm{N}}$, $\Gamma_{\mathrm{D}} \cap \Gamma_{\mathrm{N}} = \emptyset$.

Assume $f \in \mathrm{L}^2(\Omega)$, $g \in \mathrm{L}^2(\Gamma_{\mathrm{N}})$ and $\mathrm{meas}_{d-1}(\Gamma_{\mathrm{D}}) > 0$. Let $\mathrm{H}^1(\Omega)$ be the usual Sobolev space. Its subspace of functions with zero trace on $\Gamma_{\mathrm{D}}$ is denoted by $\mathrm{H}_{\mathrm{o}}^1(\Omega)$. The corresponding variational formulation for (1) becomes:

$$
\left.
\begin{aligned}
&\text{Find } u \in \mathrm{H}_{\mathrm{o}}^1(\Omega): \quad a(u,v) = \langle \mathfrak{f}, v\rangle \qquad \forall\, v \in \mathrm{H}_{\mathrm{o}}^1(\Omega) \\
&\text{with} \quad a(u,v) := \int_\Omega \varepsilon \cdot (\nabla u)^\top \nabla v \, + \, u\,v \qquad \langle \mathfrak{f}, v\rangle := \int_\Omega fv + \int_{\Gamma_{\mathrm{N}}} gv.
\end{aligned}
\right\} \tag{2}
$$

We utilize a family $\mathcal{F} = \{\mathcal{T}_h\}$ of triangulations $\mathcal{T}_h$ of $\Omega$. Let $V_{o,h} \subset \mathrm{H}_{\mathrm{o}}^1(\Omega)$ be the space of continuous, piecewise linear functions over $\mathcal{T}_h$ that vanish on $\Gamma_{\mathrm{D}}$. Then the finite element solution $u_h \in V_{o,h}$ is uniquely defined by

$$
a(u_h, v_h) = \langle \mathfrak{f}, v_h\rangle \qquad \forall\, v_h \in V_{o,h}. \tag{3}
$$

Due to the Lax–Milgram Lemma both problems (2) and (3) admit unique solutions.

The main purpose of our analysis is to bound the error $u - u_h$ uniformly in the small perturbation parameter $\varepsilon$. Here we concentrate on the most natural norm related to (2), namely the *energy norm*

$$
\vert\!\vert\!\vert v \vert\!\vert\!\vert^2 := a(v,v) = \varepsilon\|\nabla v\|^2 \, + \, \|v\|^2
$$

which has been used also by other authors [1, 22, 26]. This energy norm is well–suited to produce appropriately refined meshes. This can be easily verified on some 1D model problem, *e.g.* for $-\varepsilon u'' + u = 0$ in $\Omega = (0, 1)$ with $u(0) = 1, u(1) = 0$. Even the optimal order of convergence can be achieved, *cf.* the exposition in [17].

## 3. Notation, basic tools and lemmas

In order to analyse error estimators on anisotropic meshes we will now introduce certain notation as well as important tools, all of which have proven to be advantageous in previous works [13, 16, 18]. All expositions are given for the more technical three dimensional case. The application to the simpler 2D case is readily possible.

From now on, let $\mathbb{P}^k(\omega)$ be the space of polynomials of order $k$ at most over some domain $\omega \subset \mathbb{R}^3$ or $\omega \subset \mathbb{R}^2$. Instead of $x \leq c \cdot y$ or $c_1 x \leq y \leq c_2 x$ (with positive constants independent of $x, y$ and $\varepsilon, \mathcal{T}_h$) we use the abbreviation $x \lesssim y$ and $x \sim y$, respectively. By $\|\cdot\|_\omega$ we denote the $\mathrm{L}^2$ norm of a function over some domain $\omega$. For $\omega = \Omega$ the subscript is omitted. Let $|\omega| := \mathrm{meas}\,(\omega)$ be the measure of a domain $\omega$. Finally for some vector $\mathbf{p}$ let $|\mathbf{p}| := \sqrt{\mathbf{p}^\top \mathbf{p}}$ be its Euclidean norm (*i.e.* length).

### 3.1. Tetrahedron – Subdomains – Mesh requirements

**Tetrahedron.** Let a triangulation $\mathcal{T}_h$ be given which satisfies the usual admissibility conditions (see Ciarlet [10], Chap. 2). The four vertices of an arbitrary tetrahedron $T \in \mathcal{T}_h$ are denoted by $P_0, \ldots, P_3$ such that $P_0 P_1$ is the longest edge of $T$, $\mathrm{meas}_2(\triangle P_0 P_1 P_2) \geq \mathrm{meas}_2(\triangle P_0 P_1 P_3)$, and $\mathrm{meas}_1(P_1 P_2) \geq \mathrm{meas}_1(P_0 P_2)$.

Additionally define three pairwise orthogonal vectors $\mathbf{p}_i$ with lengths $h_{i,T} := |\mathbf{p}_i|$, see Figure 1. Observe $h_{1,T} > h_{2,T} \geq h_{3,T}$ and set $h_{\min,T} := h_{3,T}$. The circumscribed hexahedron may facilitate the visualization.
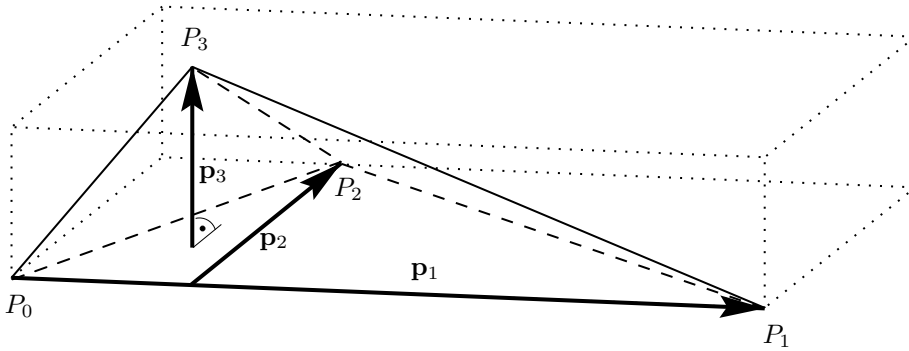
FIGURE 1.   Notation of tetrahedron $T$.

Tetrahedra are denoted by $T, T'$ or $T_i$. Faces of a tetrahedron are denoted by $E$ or $E_i$. Let $h_{E,T} := 3|T|/|E|$ be the length of the *height* over a face $E$. Because of the geometrical properties of the tetrahedron one infers $h_{\min,T} < 2 \cdot h_{E,T}$ for all faces $E$ of $T$.

**Auxiliary subdomains.** Let $T \in \mathcal{T}_h$ be an arbitrary tetrahedron. Let $\omega_T$ be that domain that is formed by $T$ and all tetrahedra that have a common face with $T$. Note that $\omega_T$ consists of less than five tetrahedra if $T$ has a boundary face.

Let $E$ be an inner face (triangle) of $\mathcal{T}_h$, *i.e.* there are two tetrahedra $T_1$ and $T_2$ having the common face $E$. Set the domain $\omega_E := T_1 \cup T_2$. If $E$ is a boundary face set $\omega_E := T$ with $T \supset E$.

**Mesh requirements.** In addition to the usual conformity conditions of the mesh (see Ciarlet [10], Chap. 2) we demand the following two assumptions.
  1. The number of tetrahedra containing a node $x_j$ is bounded uniformly.
  2. The dimensions of adjacent tetrahedra must not change rapidly, *i.e.*

$$h_{i,T'} \sim h_{i,T} \qquad \forall\, T, T' \text{ with } T \cap T' \neq \emptyset\,,\, i = 1 \ldots d.$$

**Remark 3.1.** In certain situations we do not want to use *element based* quantities (such as $h_{\min,T}$) but utilize *face related* terms instead. For example consider an interior face $E = T_1 \cap T_2$, and define the terms

$$h_E := (h_{E,T_1} + h_{E,T_2})/2, \qquad h_{\min,E} := (h_{\min,T_1} + h_{\min,T_2})/2.$$

Their advantage is that they are no longer related to $T_1$ or $T_2$ but to $E$. They clearly satisfy $h_E \sim h_{E,T_i}$ and $h_{\min,E} \sim h_{\min,T_i}$. For a boundary face $E \subset \partial T \cap \Gamma$ define similarly $h_E := h_{E,T}$ and $h_{\min,E} := h_{\min,T}$. Similar to above one can infer $h_{\min,E} < 2 \cdot h_E$ for all faces $E$.

**Transformations.** The usual transformation technique between a tetrahedron $T$ and a standard tetrahedron plays a vital role in many proofs (*cf.* [10]). However, a refined analysis has shown that *two different transformations* facilitate matters considerably, see *e.g.* [13,14]. Hence define the matrices $H_T, A_T, C_T \in \mathbb{R}^{3\times3}$ by

$$\left.\begin{aligned}
H_T &:= \operatorname{diag}(h_{1,T}, h_{2,T}, h_{3,T}), \\
A_T &:= \left( \overrightarrow{P_0P_1}, \overrightarrow{P_0P_2}, \overrightarrow{P_0P_3} \right), \\
C_T &:= \left(\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3\right),
\end{aligned}\right\} \tag{4}$$

and introduce affine linear mappings

$$F_A(\mu) := A_T \cdot \mu + \overrightarrow{P_0} \qquad \text{and} \qquad F_C(\mu) := C_T \cdot \mu + \overrightarrow{P_0} \quad, \mu \in \mathbb{R}^3.$$

These mappings implicitly define the *standard tetrahedron* $\bar{T} := F_A^{-1}(T)$ and the *reference tetrahedron* $\hat{T} := F_C^{-1}(T)$. Then $\bar{T}$ the has vertices $\bar{P}_0 = (0,0,0)^\top$ and $\bar{P}_i = \mathbf{e}_i^\top, i = 1\dots 3$, whereas $\hat{T}$ has vertices at $\hat{P}_0 = (0,0,0)^\top$, $\hat{P}_1 = (1,0,0)^\top$, $\hat{P}_2 = (\hat{x}_2, 1, 0)^\top$ and $\hat{P}_3 = (\hat{x}_3, \hat{y}_3, 1)^\top$. The conditions on the $P_i$ yield immediately $0 < \hat{x}_2 \le 1/2$, $0 < \hat{x}_3 < 1$ and $-1 < \hat{y}_3 < 1$. Figures 1 and 2 illustrate this definition.
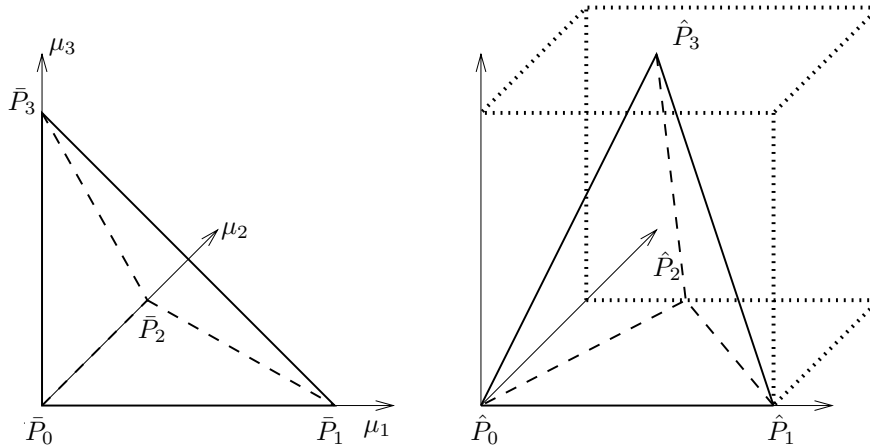


FIGURE 2. Standard tetrahedron $\bar{T}$ and reference tetrahedron $\hat{T}$.

Variables and operators that are related to the standard tetrahedron $\bar{T}$ and the reference tetrahedron $\hat{T}$ are referred to with a *bar* and a *hat*, respectively (*e.g.* $\bar{\nabla}$, $\hat{v}$). The determinants of both mappings are $|\det(A_T)| = |\det(C_T)| = 6|T|$, and the transformed derivatives satisfy $\bar{\nabla}\bar{v} = A_T^\top \nabla v$ and $\hat{\nabla}\hat{v} = C_T^\top \nabla v$.

Although $C_T$ is naturally associated with our analysis, it transforms $\hat{T}$ into $T$. Inequality constants would thus depend on $\hat{T}$. To overcome this drawback, the transformation *via* $A_T$ is used in conjunction with $C_T$ (*cf.* the compactness arguments in the proof of Lem. 4.2).

Finally, $H_T^{-1} C_T^\top$ is orthogonal since $C_T^\top \cdot C_T = H_T^2$. Hence

$$\|H_T^{-1} C_T^\top \nabla v\|_T = \|\nabla v\|_T. \tag{5}$$

**Squeezed tetrahedron $T_{E,\delta}$.** The concept of the squeezed tetrahedron has been introduced in [18] and originates from [26] (in a simpler, modified form there). Here we repeat the definition and only state the required results.

Because of the singular perturbation character of the differential equation we can favourably employ a sub–tetrahedron $T_{E,\delta} \subset T$ which depends on a face $E$ of $T$ and a real number $\delta \in (0,1]$. In an attempt to use a vivid name we will refer to $T_{E,\delta}$ as a *squeezed tetrahedron*. For its precise definition, let $T$ be an arbitrary but fixed tetrahedron, and enumerate temporarily its vertices such that $E = Q_1 Q_2 Q_3$ and $T = OQ_1 Q_2 Q_3$, *cf.* Figure 3. Introduce barycentric coordinates such that $\lambda_0$ is related to $O$, and $\lambda_1$, $\lambda_2$, $\lambda_3$ correspond to $Q_1, Q_2, Q_3$, respectively.

Let $P$ be that point with barycentric coordinates

$$\lambda_0(P) = \delta \qquad \text{and} \qquad \lambda_1(P) = \lambda_2(P) = \lambda_3(P) = \frac{1-\delta}{3}.$$

Then let $T_{E,\delta}$ be the tetrahedron with vertices $P$ and $Q_1, Q_2, Q_3$, *i.e.* $T_{E,\delta}$ has the same face $E$ as $T$ but the fourth vertex is moved towards $E$ with the rate $\delta$.

An alternative description is as follows. With $S_E$ being the midpoint (*i.e.* center of gravity) of face $E$, point $P$ lies on the line $S_E O$ such that $|\vec{S_E P}| = \delta \cdot |\vec{S_E O}|$. Note that for $\delta = 1$ one gets $T_{E,\delta} \equiv T$ whereas in the limiting case $\delta \to 0$ the tetrahedron $T_{E,\delta}$ collapses to the face $E$.
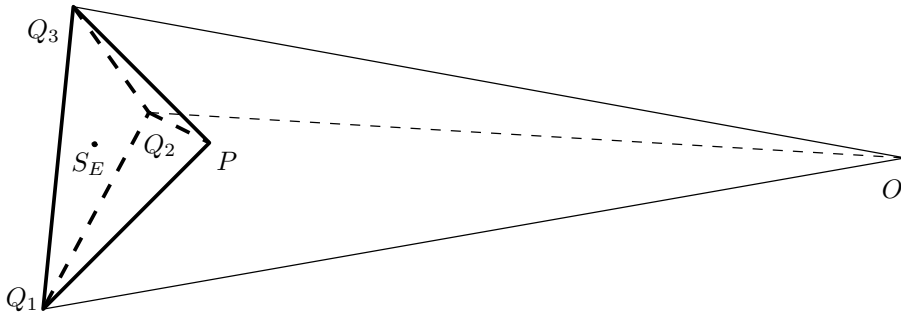
FIGURE 3.    Tetrahedra $T = OQ_1Q_2Q_3$ and $T_{E,\delta} = PQ_1Q_2Q_3$.

In order to utilize $T_{E,\delta}$ efficiently, we also require an affine linear transformation $F_{T,E,\delta}$ that maps the standard tetrahedron $\bar{T}$ onto $T_{E,\delta}$. This affine linear mapping is unique (up to permutations of the enumeration of the vertices of $\bar{T}$ and $T_{E,\delta}$). In [18] the following relations has been proven.

**Lemma 3.2.** *The radius $\varrho(T_{E,\delta})$ of the largest inscribed sphere of $T_{E,\delta}$ is equivalent to*

$$\varrho(T_{E,\delta}) \ \sim \ \min\{\delta \cdot h_{E,T} \ , \ h_{\min,T}\} \ \sim \ h_{\min,T_{E,\delta}}. \tag{6}$$

*The norm of the transformation matrix $F_{T,E,\delta}^{-1}$ is bounded by*

$$\|F_{T,E,\delta}^{-1}\|_{\mathbb{R}^{3\times3}} \ \lesssim \ \min\{\delta \cdot h_{E,T} \ , \ h_{\min,T}\}^{-1}.$$

### 3.2. Bubble functions

Another useful and important tool are so-called bubble functions which are applied, for example, for defining the local problem and its ansatz space but also for the analysis. The bubble functions were already partially introduced in [25] and [18].

Denote by $\lambda_{T,1}, \cdots, \lambda_{T,4}$ the barycentric coordinates of an arbitrary tetrahedron $T$. The *element bubble function* $b_T$ is defined by

$$b_T := 4^4 \cdot \lambda_{T,1} \cdot \lambda_{T,2} \cdot \lambda_{T,3} \cdot \lambda_{T,4} \in \mathbb{P}^4(T) \qquad \text{on } T. \tag{7}$$

For simplicity assume that $b_T$ is extended by zero outside its original domain of definition.

Further we require face bubble functions. To this end let $E = T_1 \cap T_2$ be an inner face (triangle) of $\mathcal{T}_h$. Enumerate the vertices of $T_1$ and $T_2$ such that the vertices of $E$ are numbered first, and introduce the functions

$$b_{E,T_i} := 3^3 \cdot \lambda_{T_i,1} \cdot \lambda_{T_i,2} \cdot \lambda_{T_i,3} \qquad \text{on } T_i, \ i = 1, 2.$$

The *standard face bubble function* $b_E \in \mathrm{C}^0(\omega_E)$ is now defined in a piecewise fashion (with support

$\omega_E = T_1 \cup T_2$) by

$$
b_E := \begin{cases} b_{E,T_1} & \text{on } T_1 \\ b_{E,T_2} & \text{on } T_2 \\ 0 & \text{otherwise,} \end{cases}
$$

see also the middle of Figure 4. Note that $0 \le b_T(\mathbf{x}), b_E(\mathbf{x}) \le 1$ and $\|b_T\|_\infty = \|b_E\|_\infty = 1$.

For clarity of notation we also introduce a trivial extension operator $F_{\text{ext}} : \mathbb{P}^0(E) \to \mathbb{P}^0(\omega_E)$ that maps a constant function over some face $E$ to the same constant function acting on $\omega_E$. If $E$ is a boundary face then $b_E$ and $F_{\text{ext}}$ are obviously defined only on the single tetrahedron $T \supset E$.

The following anisotropic equivalences/inverse inequalities can be derived easily.

**Lemma 3.3** (Inverse inequalities I). *Assume that $\varphi_T \in \mathbb{P}^0(T)$ and $\varphi_E \in \mathbb{P}^0(E)$. Then*

$$\|b_T\|_T \sim |T|^{1/2} \tag{8}$$

$$\|b_T^{1/2} \cdot \varphi_T\|_T \sim \|\varphi_T\|_T \tag{9}$$

$$\|\nabla(b_T \cdot \varphi_T)\|_T \lesssim h_{\min,T}^{-1} \cdot \|\varphi_T\|_T \tag{10}$$

$$\|b_E^{1/2} \cdot \varphi_E\|_E \sim \|\varphi_E\|_E. \tag{11}$$

*Proof.* The proofs employ standard scaling arguments; they are also given in [13].  ☐

The bubble functions from above suffice to analyse the error estimator for the *Poisson* equation, *cf.* [16,25]. However, for the *singularly perturbed problem* considered here we have to introduce modified face bubble functions, *cf.* also [13,26].

Start with some face $E$ and let $T_1, T_2$ be its two neighbouring tetrahedra, *i.e.* $\omega_E = T_1 \cup T_2$. For an arbitrary real number $\delta \in (0,1]$ consider both squeezed tetrahedra $T_{1,E,\delta} \subset T_1$ and $T_{2,E,\delta} \subset T_2$, *cf.* Figures 3 and 4. Now we are ready to present the so–called *squeezed face bubble function* $b_{E,\delta}$ which acts only on $T_{1,E,\delta} \cup T_{2,E,\delta} \subset \omega_E$. Its piecewise definition is

$$
b_{E,\delta} := \begin{cases} b_{\bar{E}} \circ F_{T_1,E,\delta}^{-1} & \text{on } T_{1,E,\delta} \\ b_{\bar{E}} \circ F_{T_2,E,\delta}^{-1} & \text{on } T_{2,E,\delta} \\ 0 & \text{on } \omega_E \setminus (T_{1,E,\delta} \cup T_{2,E,\delta}), \end{cases} \tag{12}
$$

where $b_{\bar{E}}$ is the standard face bubble function for the face $\bar{E} = F_{T_i,E,\delta}^{-1}(E)$ of the tetrahedron $\bar{T} = F_{T_i,E,\delta}^{-1}(T_{i,E,\delta})$. Note that the squeezed face bubble function on $T_i$ can equivalently be viewed as the standard face bubble function on the squeezed tetrahedron $T_{i,E,\delta}$, *i.e.*

$$b_{E,\delta}\big|_{T_i} \equiv b_{E,T_{i,E,\delta}}.$$

Figure 4 may facilitate the understanding of the standard/squeezed face bubble function for the two–dimensional case. For boundary faces one restricts $b_{E,\delta}$ to the unique tetrahedron with $\partial T \supset E$.

Standard scaling arguments for the transformation $F_{T_i,E,\delta} : \bar{T} \to T_{i,E,\delta}$, together with the essential Lemma 3.2 yield now the inverse inequalities for the squeezed face bubble function.

**Lemma 3.4** (Inverse equivalences II). *Let $E$ be an arbitrary face of $T$, assume $\varphi_E \in \mathbb{P}^0(E)$, and let $\delta \in (0,1]$ be arbitrary. Then one has*

$$\|b_{E,\delta} \cdot F_{\text{ext}}(\varphi_E)\|_T \sim \delta^{1/2} \cdot h_{E,T}^{1/2} \cdot \|\varphi_E\|_E \tag{13}$$

$$\|\nabla(b_{E,\delta} \cdot F_{\text{ext}}(\varphi_E))\|_T \sim \delta^{1/2} \cdot h_{E,T}^{1/2} \cdot \min\{\delta \cdot h_{E,T} , h_{\min,T}\}^{-1} \cdot \|\varphi_E\|_E. \tag{14}$$
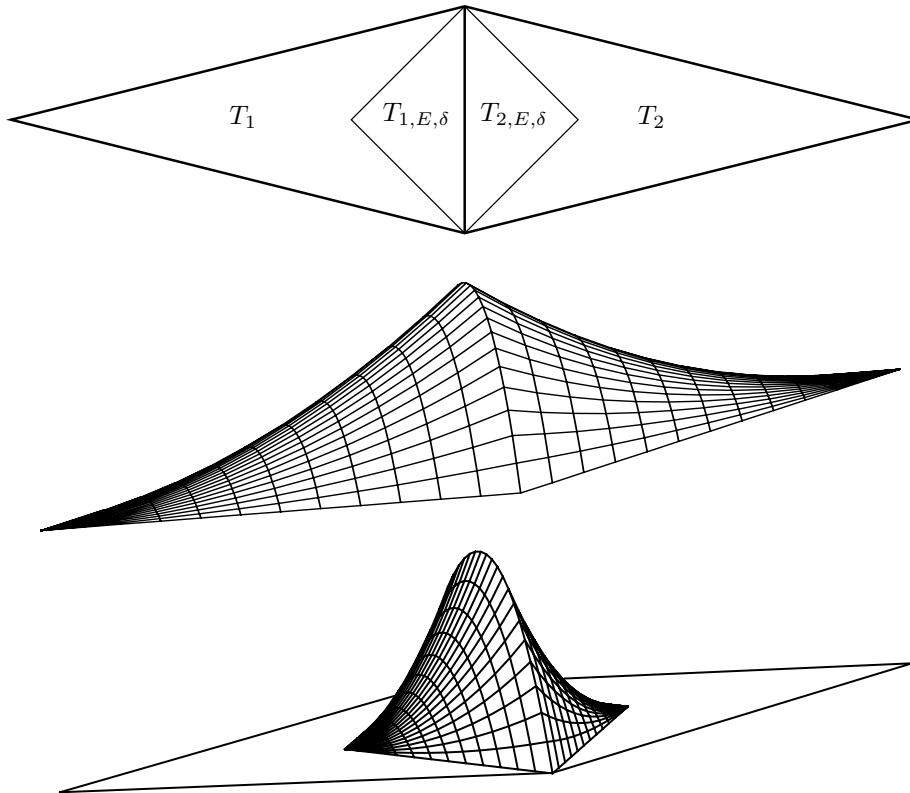
FIGURE 4. Top: $\omega_E$ and squeezed triangles $T_{i,E,\delta}$ (2D case). Middle: standard face bubble function $b_E$. Bottom: squeezed face bubble function $b_{E,\delta}$.

*Proof.* Standard scaling arguments for the transformation $F_{T,E,\delta} : \bar{T} \to T_{E,\delta}$ readily imply (13).

For (14) we start with the equivalence

$$\|\nabla b_E\|_T \sim h_{\min,T}^{-1} \cdot |T|^{1/2}$$

which has been proven (in a slightly different form) in [13], (Lem. 3.5) and [16], (Lem. 5). Above we have realized that the squeezed face bubble function can also be viewed as the standard face bubble function on the squeezed tetrahedron. Thus one can utilize the previous equivalence with the function $b_{E,\delta}$ on the tetrahedron $T_{E,\delta}$. Together with Lemma 3.2 and $|T_{E,\delta}| = \delta \cdot |T|$ this results in

$$\|\nabla b_{E,\delta}\|_T \sim h_{\min,T_{E,\delta}}^{-1} \cdot |T_{E,\delta}|^{1/2}$$
$$\overset{(6)}{\sim} \min\{\delta \cdot h_{E,T} , h_{\min,T}\}^{-1} \cdot \delta^{1/2} \cdot (h_{E,T}|E|)^{1/2}$$

which completes the proof. $\qquad \square$

### 3.3. Matching function and residual error estimator

Parts of the analysis of the local problem error estimator rely on results for the anisotropic residual error estimator of [18] which are thus repeated now for self–containment. Note that both estimators have been developed in close collaboration to enable certain equivalence properties (*cf.* Th. 4.3 below). Related aspects are discussed here as well.

Let us start with an important difference between error estimation on *isotropic* and *anisotropic* meshes. For isotropic meshes the error estimation is valid no matter what actual mesh is used. In contrast, this feature is lost on anisotropic meshes where all known error estimators require the anisotropy of the mesh to be aligned with the anisotropy of the solution. Heuristically this means that anisotropic elements (*e.g.* tetrahedra) are stretched in that direction where the solution shows little variation. If this requirement is violated then the upper and lower error bound may differ by an arbitrarily large factor (depending on the degree of the misalignment of mesh and solution).

In order to investigate this matter mathematically, let us recall the proposals from known (analytically based) anisotropic error estimators. Siebert [23] restricts the set of treatable anisotropic functions. Kunert [13, 14, 18] introduces a so–called matching function $m_1(v, \mathcal{T}_h)$ that measures the alignment of an anisotropic function $v$ and an anisotropic mesh $\mathcal{T}_h$. Lastly, in Dobrowolski *et al.* [11] a saturation assumption is necessary that implies a similar correspondence of the anisotropic mesh and the anisotropic solution.

For a rigorous analysis it is advantageous to *measure* the alignment of mesh and function. To this end the *matching function* has been proposed by Kunert [13, 14]:

**Definition 3.5** (Matching function)**.** Let $v \in \mathrm{H}^1(\Omega)$, and $\mathcal{T}_h \in \mathcal{F}$ be a triangulation of $\Omega$. Define the *matching function* $m_1 : \mathrm{H}^1(\Omega) \times \mathcal{F} \mapsto \mathbb{R}$ by

$$m_1(v, \mathcal{T}_h) \; := \; \left( \sum_{T \in \mathcal{T}_h} h_{\min,T}^{-2} \cdot \|C_T^\top \nabla v\|_T^2 \right)^{1/2} \Big/ \; \|\nabla v\|. \tag{15}$$

Note that the entries of the vector $C_T^\top \nabla v \equiv (\mathbf{p}_1^\top \nabla v, \mathbf{p}_2^\top \nabla v, \mathbf{p}_3^\top \nabla v)^\top$ can also be viewed as scaled directional derivatives along the orthogonal directions $\mathbf{p}_i$ (recall $|\mathbf{p}_i| = h_{i,T}$).
To deepen the understanding of the matching function let us briefly discuss its behaviour and influence. More details and a comprehensive discussion can be found in [13, 14].
By defining temporarily $h_{\max,T} := h_{1,T}$, one obtains

$$1 \; \leq \; m_1(v, \mathcal{T}_h) \lesssim \max_{T \in \mathcal{T}_h} \frac{h_{\max,T}}{h_{\min,T}}.$$

Although this crude upper bound is useless for practical purposes it implies $m_1 \sim 1$ on isotropic meshes. Then $m_1$ merges with other constants and becomes invisible; in this sense (15) is an extension of the theory for isotropic meshes. If an anisotropic mesh $\mathcal{T}_h$ is *well aligned* with an anisotropic function $v$ then one also obtains $m_1(v, \mathcal{T}_h) \sim 1$. If, however, the anisotropic meshes are *not aligned* with the function then the matching function can be arbitrarily large, $m_1(v, \mathcal{T}_h) \gg 1$.
The influence of the matching function $m_1$ can be seen in the error bound (19) of Lemma 3.6 and in the discussion afterwards.

Next the residual error estimator will be presented. The methodology to obtain a lower error bound requires residual terms from a finite dimensional space [13, 25]. Hence we replace the exact element residual by an approximate element residual which is constant over an element $T$ (*e.g.* by means of an $\mathrm{L}^2$ projection into $\mathbb{P}^0(T)$). Proceed analogously for the face residuals where $g$ is replaced by $g_h$ which is piecewise constant over the Neumann faces. The precise definitions are as follows.

**Element and face residual.** The exact element residual over an element $T$ is given by

$$R_T := f \; - \; (-\varepsilon \Delta u_h + u_h) \qquad \text{on } T.$$

The *(approximate) element residual* $r_T$ is any approximation to $R_T$ that is constant on $T$, *i.e.*

$$r_T \in \mathbb{P}^0(T).$$

For $x \in E$ define the *(approximate) face residual* $r_E \in \mathbb{P}^0(E)$ by

$$
r_E(x) := \begin{cases}
\varepsilon \cdot \lim\limits_{t \to +0} \left[ \dfrac{\partial u_h}{\partial n_E}(x + t n_E) - \dfrac{\partial u_h}{\partial n_E}(x - t n_E) \right] & \text{if } E \subset \Omega \setminus \Gamma \\[2ex]
g_h - \varepsilon \cdot \partial u_h / \partial n & \text{if } E \subset \Gamma_{\mathrm{N}} \\[1ex]
0 & \text{if } E \subset \Gamma_{\mathrm{D}}.
\end{cases}
$$

Here $n_E \perp E$ is any of the two unitary normal vectors whereas $n \perp E \subset \Gamma_{\mathrm{N}}$ denotes the outer unitary normal vector.

**Residual scaling factor.** The residuals are often accompanied by the factor

$$
\alpha_T := \min\{1, \varepsilon^{-1/2} \cdot h_{\min,T}\} \cdot \tag{16}
$$

This factor plays a similar role as the local Peclet number does for diffusion convection problems.

For some interior face $E = T_1 \cap T_2$ we define the corresponding *face related* term by

$$
\alpha_E := (\alpha_{T_1} + \alpha_{T_2})/2 = \min\{1, \varepsilon^{-1/2} \cdot h_{\min,E}\} \cdot \tag{17}
$$

For boundary faces $E$ set similarly $\alpha_E := \alpha_T$ for $E \subset \partial T$. Note that the mesh requirements imply $\alpha_E \sim \alpha_{T_1} \sim \alpha_{T_2}$, *cf.* also Remark 3.1.

**Local residual error estimator.** For a tetrahedron $T$, define it by

$$
\eta_{\varepsilon,\mathrm{R},T} := \left( \alpha_T^2 \cdot \|r_T\|_T^2 + \varepsilon^{-1/2} \cdot \alpha_T \cdot \sum_{E \subset \partial T \setminus \Gamma_{\mathrm{D}}} \|r_E\|_E^2 \right)^{1/2}. \tag{18}
$$

**Local data approximation term.** To shorten the notation, introduce the term

$$
\zeta_{\varepsilon,T} := \left( \alpha_T^2 \cdot \sum_{T' \subset \omega_T} \|R_{T'} - r_{T'}\|_{T'}^2 + \varepsilon^{-1/2} \cdot \alpha_T \cdot \sum_{E \subset \partial T \cap \Gamma_{\mathrm{N}}} \|g - g_h\|_E^2 \right)^{1/2}
$$

that can also be viewed as a consistency error expression. Finally, define the *global* terms

$$
\eta_{\varepsilon,\mathrm{R}}^2 := \sum_{T \in \mathcal{T}_h} \eta_{\varepsilon,\mathrm{R},T}^2 \qquad \text{and} \qquad \zeta_\varepsilon^2 := \sum_{T \in \mathcal{T}_h} \zeta_{\varepsilon,T}^2.
$$

The following residual error estimation essentially has been proven by Kunert [18]. Here we have included the treatment of Neumann boundary conditions. Additionally the approximate element residual $r_T$ is constant here (instead of linear as in [18]).

**Lemma 3.6.** *The error is bounded locally from below for all $T \in \mathcal{T}_h$ by*

$$
\eta_{\varepsilon,\mathrm{R},T} \lesssim \|u - u_h\|_{\omega_T} + \zeta_{\varepsilon,T}.
$$

*The error is bounded globally from above by*

$$
\|u - u_h\| \lesssim m_1(u - u_h, \mathcal{T}_h) \cdot \left[ \eta_{\varepsilon,\mathrm{R}}^2 + \zeta_\varepsilon^2 \right]^{1/2}. \tag{19}
$$

*Both error bounds are uniform in $\varepsilon$.*

We remark that only the upper error bound contains the matching function $m_1$. Hence only (19) is influenced by the degree of the alignment of mesh and function, *i.e.* the global error estimator $\eta_{\varepsilon,\mathrm{R}}$ is emphasized by the factor $m_1(u - u_h, \mathcal{T}_h)$. When $m_1 \sim 1$ the lower and upper error bound possess the same quality. Obviously, the smaller $m_1$ the better the upper error bound. In the case of $m_1 \gg 1$ however both error bounds differ by a large factor which results in a useless error estimation.

**Remark 3.7.** Note that the upper error bound of (19) can not be computed exactly as it contains $m_1(u-u_h, \mathcal{T}_h)$ and thus the (unknown) error $u - u_h$. As a remedy $m_1$ can be approximated, *e.g.* by means of a recovered gradient $\nabla^{\mathrm{R}} u_h \approx \nabla u$:

$$
\begin{aligned}
m_1(u - u_h, \mathcal{T}_h) &\equiv \left( \sum_{T \in \mathcal{T}_h} h_{\min,T}^{-2} \cdot \| C_T^\top \nabla(u - u_h) \|_T^2 \right)^{1/2} \Big/ \| \nabla(u - u_h) \| \\
&\approx \left( \sum_{T \in \mathcal{T}_h} h_{\min,T}^{-2} \cdot \| C_T^\top (\nabla^{\mathrm{R}} u_h - \nabla u_h) \|_T^2 \right)^{1/2} \Big/ \| \nabla^{\mathrm{R}} u_h - \nabla u_h \| \\
&=: m_1^{\mathrm{R}}(u_h, \mathcal{T}_h),
\end{aligned}
\tag{20}
$$

*cf.* [14] for a more comprehensive discussion. Nevertheless we tried to bound $m_1(u - u_h, \mathcal{T}_h)$ by computable expressions, or at least evaluate how close $m_1^{\mathrm{R}}(u_h, \mathcal{T}_h)$ is to $m_1(u - u_h, \mathcal{T}_h)$. Unfortunately neither aim has been achieved yet. We argue, however, that this failure is not as problematic as it may seem at first glance.

Let us start with a practical point of view, and recall that our singularly perturbed problem is solved by a standard finite element method. Useful results in this case can only be obtained if the anisotropic mesh reflects the anisotropic behaviour of the solution. Implicitly this means that the different sizes of the components of $\nabla(u - u_h)$ are balanced by the different length scales of the elements (*via* $C_T$). Recalling the definition of the matching function, $m_1(u - u_h, \mathcal{T}_h)$ is then likely to be small, and $m_1(u - u_h, \mathcal{T}_h) \approx m_1^{\mathrm{R}}(u_h, \mathcal{T}_h)$. Numerical experiments described below clearly indicate that $m_1^{\mathrm{R}}$ is a robust approximation to $m_1$ in practice.

Secondly, if one seeks a theoretical approach, one would have to exploit the balance between $\nabla(u-u_h)$ and $C_T$. This is a rather delicate issue. Indeed, the saturation assumption of [11] is just a different formulation of the same difficulty. Alternatively one might want to employ superconvergence properties. However, superconvergence occurs only in an asymptotic sense, and it requires highly structured meshes that have to be aligned with the anisotropic solution. Hence this approach seems to be unsuitable for practical anisotropic problems.

## 4. LOCAL PROBLEM ERROR ESTIMATORS

### 4.1. Definition of the error estimator $\eta_{\mathrm{D}}$

The main ideas behind local problem error estimation have been known for a long time [2,6,8,24,25]. Basically the problem is solved locally but with higher accuracy, and the difference between the new solution and the original finite element solution serves as error estimator.

In each one of the aforementioned sources an *isotropic* mesh is assumed. In contrast to this the author has shown in [16] that reliable local problem error estimation is possible on *anisotropic* meshes as well. There a Poisson model problem has been investigated, the methodology of the analysis has been presented, and some important new tools and results have been developed.

In our work here we demonstrate that anisotropic local problem error estimation is not restricted to the Poisson problem but that it can be extended to the singularly perturbed reaction–diffusion problem (1). We propose a new error estimator for the latter problem. Note that the only other local problem error estimator for a singularly perturbed reaction–diffusion problem is due to [1] where an isotropic mesh is assumed, and where the local problem is infinite dimensional.

While the general structure of the proofs here is similar to the ones for the local problem error estimator for the Poisson problem [16], the actual ingredients differ. This mainly concerns the squeezed tetrahedron and

its properties as well as the squeezed face bubble functions which play a vital role in almost all analysis. The definitions of the error estimators $\eta_{\varepsilon,\mathrm{D},T}$ and $\eta_{\varepsilon,\mathrm{R},T}$ require a very careful balancing of all scaling factors (*e.g.* $\alpha_T$ from (16)) and of the "squeezing" parameter $\delta_E$ from (21). Consequently the proof of the vital Lemma 4.2 is even more technical than in [16], see also Appendix A. The derivation of a stable basis for the local problem is different from [16]. Furthermore special care has to be taken to obtain a feasible implementation of the estimator. Hence computational aspects and difficulties are addressed.

The remainder of this section is devoted to the definition of the local problem and the error estimator. Then Lemma 4.2 gives two central inequalities for the local space. Next, Theorem 4.3 states the equivalence of the local problem error estimator $\eta_{\varepsilon,\mathrm{D},T}$ and the residual error estimator $\eta_{\varepsilon,\mathrm{R},T}$. The main results, namely lower and upper bounds of the error, are given in Theorem 4.4. In Theorem 4.7 it is shown that a certain basis of the local space $V_T$ is stable (*i.e.* the local Dirichlet problem is well–conditioned). Finally other choices of local problem error estimators are feasible as well. This is demonstrated exemplarily in Section 4.4 for a face based estimator.

When deriving the error estimator, the corresponding local problem should be cheap to solve but simultaneously be rich enough to extract information on the error $e := u - u_h$. Here the subdomain of the local problem is chosen to be $\omega_T$. Let

$$\mathrm{H}_{\mathrm{o}}^1(\omega_T) := \left\{ v \in \mathrm{H}^1(\Omega) : \operatorname{supp} v \subseteq \omega_T, \quad v = 0 \ \text{ on } \partial \omega_T \setminus \Gamma_{\mathrm{N}} \right\}.$$

For an arbitrary function $v \in \mathrm{H}_{\mathrm{o}}^1(\omega_T)$ the error then satisfies

$$a(u - u_h, v)\Big|_{\omega_T} = \int\limits_{\omega_T} f \cdot v + \int\limits_{\partial \omega_T \cap \Gamma_{\mathrm{N}}} g \cdot v - \int\limits_{\omega_T} \varepsilon (\nabla u_h)^\top \nabla v - \int\limits_{\omega_T} u_h v.$$

The local problem is obtained by approximating the space $\mathrm{H}_{\mathrm{o}}^1(\omega_T)$ by some local, finite dimensional space $V_T \subset \mathrm{H}_{\mathrm{o}}^1(\omega_T)$ which is spanned by an element bubble function and some squeezed face bubble functions. Their "squeezing" parameters $\delta_E$ (*cf.* (12)) are now specified to be

$$\delta_E := \min\left\{ 1, \frac{h_{\min,E}}{h_E}, \frac{\sqrt{\varepsilon}}{h_E} \right\}. \tag{21}$$

Recalling $h_{\min,E} < 2 \cdot h_E$ from Remark 3.1 we conclude $\delta_E \sim \min\{h_{\min,E}/h_E, \sqrt{\varepsilon}/h_E\} = \sqrt{\varepsilon} h_E^{-1} \alpha_E$.

The local space $V_T$ is defined by

$$V_T := \operatorname{span}\{b_T, b_{E,\delta_E} : E \subset \partial T \setminus \Gamma_{\mathrm{D}}\}. \tag{22}$$

The local problem can be formulated most conveniently by means of the (approximate) residuals.

**Definition 4.1** (Local Dirichlet problem error estimator).
Find a solution $e_T \in V_T$ of the local variational problem:

$$a(e_T, v_T) \equiv \int\limits_{\omega_T} \varepsilon (\nabla e_T)^\top \nabla v_T + e_T \, v_T$$

$$= \sum_{T' \subset \omega_T} \int_{T'} r_{T'} \cdot v_T + \sum_{E \subset \partial T \setminus \Gamma_{\mathrm{D}}} \int_E r_E \cdot v_T \tag{23}$$

for all $v_T \in V_T$. The *local* and *global error estimators* then become

$$\eta_{\varepsilon,\mathrm{D},T} := \|e_T\|_{\omega_T} \qquad \text{and} \qquad \eta_{\varepsilon,\mathrm{D}}^2 := \sum_{T \in \mathcal{T}_h} \eta_{\varepsilon,\mathrm{D},T}^2. \tag{24}$$

Note that the particular choice of the local ansatz space $V_T$ (namely $v_T = 0$ on $\partial\omega_T \setminus \partial T$) reduces certain boundary integrals and norms. An equivalent formulation of the local problem is derived by partial integration.

**Alternative:** Find $e_T \in V_T$ such that

$$a(e_T, v_T) \;=\; a(u - u_h, v_T) - \sum_{T' \subset \omega_T} \int_{T'} (R_{T'} - r_{T'})\, v_T - \int_{\partial T \cap \Gamma_{\mathrm{N}}} (g - g_h)\, v_T \qquad \forall\, v_T \in V_T. \tag{25}$$

## 4.2. Equivalence and bounds of the local problem estimator

The methodology of the error estimator partly utilizes ideas that have already been introduced for the anisotropic local problem estimator for the Poisson problem [16], and for the anisotropic residual estimator for a singularly perturbed reaction–diffusion equation [18]. All the details however are original and new. The first lemma plays a central role in the analysis of the estimator.

**Lemma 4.2.** *The following relations hold for all $v_T \in V_T$.*

$$\|v_T\|_{\omega_T} \lesssim h_{\min,T} \cdot \|\nabla v_T\|_{\omega_T} \tag{26}$$

$$\|v_T\|_E \lesssim h_E^{-1/2}\, \delta_E^{-1/2} \cdot \min\{h_{\min,T}, \delta_E\, h_E\} \cdot \|\nabla v_T\|_{\omega_T} \qquad \forall\, E \subset \partial T. \tag{27}$$

*The inequalities are uniform in the squeezing parameters $\delta_E \in (0,1]$ which define the space $V_T$.*

*If $T$ has at least two Neumann boundary faces then the constants in (26), (27) can depend on the shape of the Neumann boundary (but do not depend on the triangulation $\mathcal{T}_h$ nor on $T$). More precisely, this Neumann boundary forms an edge at $T$, and the angle between the Neumann faces at this edge determine the constants. The smaller this angle, the worse the constants may be.*

*Proof.* The technical proof is postponed to the appendix. $\qquad\square$

**Theorem 4.3** (Equivalence with residual error estimator). *The local problem error estimator $\eta_{\varepsilon,\mathrm{D},T}$ is equivalent to the residual error estimator $\eta_{\varepsilon,\mathrm{R},T}$ in the following sense:*

$$\eta_{\varepsilon,\mathrm{D},T}^2 \lesssim \sum_{T' \subset \omega_T} \eta_{\varepsilon,\mathrm{R},T'}^2 \tag{28}$$

$$\eta_{\varepsilon,\mathrm{R},T}^2 \lesssim \sum_{T' \subset \omega_T} \eta_{\varepsilon,\mathrm{D},T'}^2. \tag{29}$$

*Both inequalities are uniform in $\varepsilon$.*

*If $T$ has at least two Neumann boundary faces then the constant in (28) can depend on the shape of the Neumann boundary (but does not depend on the triangulation $\mathcal{T}_h$ nor on $T$).*

*Proof.* Recall the definition (24) of $\eta_{\varepsilon,\mathrm{D},T}$, observe that $e_T = 0$ on $\partial\omega_T \setminus \partial T$, and take into account the modifications for boundary faces. By integration by parts one obtains

$$\eta_{\varepsilon,\mathrm{D},T}^2 = \|\!|e_T|\!\|_{\omega_T}^2 \;=\; a(e_T, e_T) \;\overset{(23)}{=}\; \sum_{T' \subset \omega_T} \int_{T'} r_{T'} \cdot e_T + \sum_{E \subset \partial T \setminus \Gamma_{\mathrm{D}}} \int_E r_E \cdot e_T$$

$$\leq \left( \sum_{T' \subset \omega_T} \|r_{T'}\|_{T'}^2 \right)^{1/2} \cdot \|e_T\|_{\omega_T} \;+\; \sum_{E \subset \partial T \setminus \Gamma_{\mathrm{D}}} \|r_E\|_E \cdot \|e_T\|_E.$$

Now $\|e_T\|_{\omega_T}$ and $\|e_T\|_E$, $E \subset \partial T$, are to be bounded. Recall the definition of $\alpha_T$ and $\delta_E$ and apply Lemma 4.2 to obtain

$$\|e_T\|_{\omega_T} \leq \||e_T\||_{\omega_T}$$

$$\|e_T\|_{\omega_T} \overset{(26)}{\lesssim} h_{\min,T} \cdot \|\nabla e_T\|_{\omega_T} \leq h_{\min,T} \cdot \varepsilon^{-1/2} \||e_T\||_{\omega_T}$$

$$\Rightarrow \quad \|e_T\|_{\omega_T} \lesssim \min\{1, \varepsilon^{-1/2} \cdot h_{\min,T}\} \cdot \||e_T\||_{\omega_T} \equiv \alpha_T \cdot \||e_T\||_{\omega_T} \tag{30}$$

$$\text{and} \quad \|e_T\|_E \overset{(27)}{\lesssim} h_E^{-1/2} \delta_E^{-1/2} \cdot \min\{h_{\min,T}, \delta_E h_E\} \cdot \|\nabla e_T\|_{\omega_T}$$

$$\sim \varepsilon^{1/4} \alpha_E^{1/2} \|\nabla e_T\|_{\omega_T}$$

$$\lesssim \varepsilon^{-1/4} \alpha_T^{1/2} \||e_T\||_{\omega_T}. \tag{31}$$

Inserting these inequalities and utilizing $\alpha_T \sim \alpha_{T'}$ for neighboring tetrahedra results in

$$\eta_{\varepsilon,\mathrm{D},T}^2 \lesssim \left( \sum_{T' \subset \omega_T} \alpha_{T'}^2 \cdot \|r_{T'}\|_{T'}^2 + \varepsilon^{-1/2} \alpha_T \cdot \sum_{E \subset \partial T \setminus \Gamma_\mathrm{D}} \|r_E\|_E \right)^{1/2} \cdot \||e_T\||_{\omega_T}$$

which, together with $\||e_T\||_{\omega_T} = \eta_{\varepsilon,\mathrm{D},T}$, proves (28).

For the proof of (29) we require bounds of $\eta_{\varepsilon,\mathrm{R},T}$, and thus of $\|r_T\|_T$ and $\|r_E\|_E$. The structure of the proof is similar to our analysis for the Poisson equation [16].

We first bound the term $\|r_{T'}\|_{T'}$, with $T' \subset \omega_T$ being an arbitrary tetrahedron. Recall definition (7) of the bubble function $b_{T'}$ and set $v_{T'} := b_{T'} \cdot r_{T'}$. Then $b_{T'}$ and $v_{T'}$ belong to the finite element space $V_{T'}$. Hence the local problem related to $T'$ has to be invoked. The local problem (23) and equivalence (9) imply

$$\|r_{T'}\|_{T'}^2 \overset{(9)}{\sim} \|b_{T'}^{1/2} \cdot r_{T'}\|_{T'}^2 = \int_{T'} r_{T'} \cdot v_{T'} \quad \text{since } v_{T'} \in \mathrm{H}_\mathrm{o}^1(T')$$

$$\overset{(23)}{=} a(e_{T'}, v_{T'}) \leq \||e_{T'}\||_{T'} \cdot \||v_{T'}\||_{T'},$$

where $e_{T'} \in V_{T'}$ denotes the solution of the local problem over $\omega_{T'}$. Inequality (10) results in

$$\||v_{T'}\||_{T'}^2 = \varepsilon \|\nabla(b_{T'} \cdot r_{T'})\|_{T'}^2 + \|b_{T'} \cdot r_{T'}\|_{T'}^2$$

$$\overset{(10)}{\sim} \varepsilon h_{\min,T'}^{-2} \cdot \|r_{T'}\|_{T'}^2 + \|r_{T'}\|_{T'}^2 \sim \alpha_{T'}^{-2} \|r_{T'}\|_{T'}^2.$$

Combining both inequalities yields

$$\|r_{T'}\|_{T'} \lesssim \alpha_T^{-1} \cdot \||e_{T'}\||_{T'} \leq \alpha_T^{-1} \cdot \eta_{\varepsilon,\mathrm{D},T'} \qquad \forall T' \subset \omega_T \tag{32}$$

since $\alpha_{T'}$ does not change rapidly across adjacent tetrahedra $T'$.

The norm of $r_E \in \mathbb{P}^0(E)$ for an interior face $E \subset \partial T \setminus \Gamma$ is bounded similarly. Let us recall the definition (12) of the squeezed face bubble function $b_{E,\delta}$, and set $v_E := b_{E,\delta} \cdot F_{\mathrm{ext}}(r_E) \in V_T \cap \mathrm{H}_\mathrm{o}^1(\omega_E)$. Integration by parts

and $v_E = 0$ on $\partial T \cap \Gamma$ imply

$$
\begin{aligned}
\|r_E\|_E^2 &\overset{(11)}{\sim} \|b_E^{1/2} \cdot r_E\|_E^2 = \int_E r_E \cdot v_E \\
&\overset{(23)}{=} a(e_T, v_E) - \sum_{T' \subset \omega_E} \int_{T'} r_{T'} \, v_E \\
&\leq \|e_T\|_{\omega_E} \cdot \|v_E\|_{\omega_E} + \sum_{T' \subset \omega_E} \|r_{T'}\|_{T'} \cdot \|v_E\|_{T'}.
\end{aligned}
$$

Now the norms of $v_E$ are bounded by means of inverse inequalities, and by using the specific value of $\delta_E$ from (21). This leads to

$$
\begin{aligned}
\|v_E\|_{T'} &= \|b_{E,\delta} \cdot F_{\mathrm{ext}}(r_E)\|_{T'} \overset{(13)}{\lesssim} \delta_E^{1/2} \cdot h_{E,T'}^{1/2} \cdot \|r_E\|_E \overset{(21)}{\sim} \varepsilon^{1/4} \alpha_T^{1/2} \cdot \|r_E\|_E \\
\|\nabla v_E\|_{T'} &= \|\nabla(b_{E,\delta} \cdot F_{\mathrm{ext}}(r_E))\|_{T'} \\
&\overset{(14)}{\lesssim} \delta_E^{1/2} \cdot h_{E,T'}^{1/2} \cdot \min\{\delta_E \cdot h_{E,T'}, h_{\min,T'}\}^{-1} \cdot \|r_E\|_E \\
&\overset{(21)}{\sim} \varepsilon^{-1/4} \alpha_T^{-1/2} \cdot \|r_E\|_E \\
\Rightarrow \quad \|v_E\|_{\omega_E} &= (\varepsilon\|\nabla(v_E)\|_{\omega_E}^2 + \|v_E\|_{\omega_E}^2)^{1/2} \lesssim \varepsilon^{1/4} \alpha_T^{-1/2} \cdot \|r_E\|_E.
\end{aligned}
$$

Next one utilizes the previous bound (32) of $\|r_{T'}\|_{T'}$ for both tetrahedra $T' \subset \omega_E$. Combining all estimates yields

$$
\|r_E\|_E \lesssim \varepsilon^{1/4} \alpha_T^{-1/2} \cdot \sum_{T' \subset \omega_E} \eta_{\varepsilon,\mathrm{D},T'} \qquad \forall E \subset \partial T \setminus \Gamma. \tag{33}
$$

The norm of $r_E \in \mathbb{P}^0(E)$ for a Neumann boundary face $E \subset \partial T \cap \Gamma_\mathrm{N}$ is bounded similarly (*cf.* [16]) and gives analogously

$$
\|r_E\|_E \lesssim \varepsilon^{1/4} \alpha_T^{-1/2} \cdot \eta_{\varepsilon,\mathrm{D},T} \qquad \forall E \subset \partial T \cap \Gamma_\mathrm{N}.
$$

Collecting all the results for $\|r_T\|_T$ and $\|r_E\|_E$ and inserting them into the definition of $\eta_{\varepsilon,\mathrm{R},T}$ gives (29). $\qquad \square$

With the help of Theorem 4.3 we easily derive the main result, namely upper and lower error bounds by means of the local problem error estimator.

**Theorem 4.4** (Local problem error estimation).
*The error is bounded locally from below by*

$$
\eta_{\varepsilon,\mathrm{D},T} \leq \|u - u_h\|_{\omega_T} + c \cdot \zeta_{\varepsilon,T} \qquad \forall T \in \mathcal{T}_h. \tag{34}
$$

*The error is bounded globally from above by*

$$
\|u - u_h\| \lesssim m_1(u - u_h, \mathcal{T}_h) \cdot \left[\eta_{\varepsilon,\mathrm{D}}^2 + \zeta_\varepsilon^2\right]^{1/2}. \tag{35}
$$

*Both inequalities are uniform in $\varepsilon$.*

*The lower error bound (34) is a strict inequality where the only constant $c$ is at the data approximation term $\zeta_{\varepsilon,T}$. As always, this constant $c$ is independent of $\varepsilon$, $T$, $u$ and $u_h$. However, if $T$ has at least two Neumann boundary faces then $c$ can depend on the shape of the Neumann boundary (but does not depend on the triangulation $\mathcal{T}_h$ nor on $T$).*

Note that by analogy with Lemma 3.6 only the upper error bound is influenced by the matching function $m_1(u - u_h, \mathcal{T}_h)$.

*Proof.* For (34), apply formulation (25) of the local problem, recall how $\omega_T$ and $V_T$ are modified if $T$ has a boundary face, and observe in particular that $e_T = 0$ on $\partial \omega_T \setminus \partial T$. Then one obtains

$$
\begin{aligned}
\eta_{\varepsilon,\mathrm{D},T}^2 \;&=\; \interleave e_T \interleave_{\omega_T}^2 \;=\; a(e_T, e_T) \\
&\overset{(25)}{=}\; a(u - u_h, e_T) - \sum_{T' \subset \omega_T} \int_{T'} (R_{T'} - r_{T'}) \cdot e_T \;-\; \int_{\Gamma_{\mathrm{N}} \cap \partial T} (g - g_h) \cdot e_T \\
&\leq\; \interleave u - u_h \interleave_{\omega_T} \cdot \interleave e_T \interleave_{\omega_T} \;+\; \left( \sum_{T' \subset \omega_T} \| R_{T'} - r_{T'} \|_{T'} \right) \cdot \| e_T \|_{\omega_T} \;+\; \| g - g_h \|_{\Gamma_{\mathrm{N}} \cap \partial T} \cdot \| e_T \|_{\Gamma_{\mathrm{N}} \cap \partial T}.
\end{aligned}
$$

With the previous bounds (30) and (31) one readily obtains the desired estimate (34). Finally inequality (35) follows immediately from the error bound (19) of the residual error estimator, and from relation (29) between $\eta_{\varepsilon,\mathrm{R},T}$ and $\eta_{\varepsilon,\mathrm{D},T}$. $\qquad\square$

**Remark 4.5.** An anisotropic adaptive solution algorithm as described in Section 1 consists of the main ingredients: *solve; estimate error; extract anisotropic information; refine anisotropic mesh.* There exist isolated investigations for each of these ingredients, see for example, [15] and the citations therein. The implementation and analysis of the *complete* adaptive algorithm remains a challenging endeavour, particularly for complicated domains in three dimensions.

The anisotropic error estimation presented in this work answers the question of the optimal element size. The other pieces of information such as the anisotropic stretching direction and stretching ratio have to be obtained (currently) by means of other (heuristic) procedures, *e.g. via* the popular Hessian strategy, *cf.* [20] or [15].

The next ingredient, *i.e.* the *anisotropic mesh refinement*, has been considered by several authors as well. A particularly demanding task is to combine anisotropic meshes with hierarchies that are required for fast multilevel solvers. Research is increasingly focusing on such topics.

Summarizing, our anisotropic error estimation is a major step towards a general anisotropic solution algorithm.

**Remark 4.6.** The approximate element residual $r_T$ is defined here by a *constant* approximation to the exact element residual $R_T$, *cf.* Section 3.3. One can also use higher order approximations, for example a *linear* approximated element residual $r_T$. Then the residual error estimate of Lemma 3.6 remains valid (actually this case has been treated in [18]).

For the definition of the local problem error estimator, however, one would have to employ a larger local space $V_T$. This requirement stems from the analysis where the function $r_T \cdot b_T$ has to be contained in the space $V_T$ (*cf.* proof of Th. 4.3).

Hence one can achieve a more accurate error estimation in conjunction with a smaller data approximation term $\zeta_\varepsilon$. The expense is a larger local problem that is considerably more technical and expensive to implement.

### 4.3. A stable basis for the local problem

Here we will present a stable basis for the local problem under consideration. An equivalent description of this aim is that the variational problem is well-conditioned, *i.e.* the condition number of the corresponding finite element matrix is bounded independently of the perturbation parameter $\varepsilon$ and of the aspect ratio of the elements.

Recall that the local ansatz space is $V_T = \mathrm{span}\{b_T, b_{E,\delta_E} : E \subset \partial T \setminus \Gamma_{\mathrm{D}}\}$. As a basis of $V_T$ we choose

$$
\Phi \;:=\; \left( b_T \,,\; \delta_E^{-1/2} \cdot b_{E,\delta} \;:\; E \subset \partial T \setminus \Gamma_{\mathrm{D}} \right). \tag{36}
$$

For simplicity of notation enumerate the faces of $T$ such that interior and Neumann faces come first, and denote them by $E_i$, $i = 1, \ldots, m$, $m \leq 4$. Denote the parameter of the squeezed face bubble functions temporarily by $\delta_i := \delta_{E_i}$ . Hence any function $v_T \in V_T$ can be expressed as

$$v_T \; = \; \beta_0 b_T \; + \; \sum_{i=1}^{m} \beta_i \cdot \delta_i^{-1/2} \cdot b_{E_i, \delta_i} \; = \; \Phi \cdot \mathbf{v}$$

$$\text{with} \qquad \mathbf{v} := (\beta_0, \beta_1 \ldots \beta_m)^\top.$$

The stiffness matrix $K_T \in \mathbb{R}^{(1+m) \times (1+m)}$ of the local problem is given by means of the finite element isomorphism $a(v_T, w_T) \; = \; (K_T \mathbf{v}, \mathbf{w})$ for all $w_T = \Phi \cdot \mathbf{w} \in V_T$.

**Theorem 4.7** (Stable basis). *The basis (36) of $V_T$ is stable, i.e. the condition number $\kappa(K_T)$ of the local problem stiffness matrix $K_T$ is bounded uniformly in $\varepsilon$ and $T$:*

$$\kappa(K_T) \sim 1 \qquad \forall\, T \in \mathcal{T}_h.$$

*Proof.* The condition number is given by

$$\kappa(K_T) \; = \; \Big[ \max_{\mathbf{v} \neq \mathbf{0}} (K_T \mathbf{v}, \mathbf{v})/(\mathbf{v}, \mathbf{v}) \Big] \Big/ \Big[ \min_{\mathbf{w} \neq \mathbf{0}} (K_T \mathbf{w}, \mathbf{w})/(\mathbf{w}, \mathbf{w}) \Big].$$

Thus investigate the scalar product $(K_T \mathbf{v}, \mathbf{v})$ which equals

$$(K_T \mathbf{v}, \mathbf{v}) \; = \; a(v_T, v_T) \; = \; \|\|v_T\|\|_{\omega_T}^2 \; = \; \varepsilon \|\nabla v_T\|_{\omega_T}^2 \; + \; \|v_T\|_{\omega_T}^2.$$

We start by bounding $\|\|v_T\|\|_{\omega_T}$ from above. The triangle inequality readily implies

$$\|\|v_T\|\|_{\omega_T} \; \leq \; |\beta_0| \cdot \|\|b_T\|\|_{\omega_T} \; + \; \sum_{i=1}^{m} |\beta_i| \cdot \delta_i^{-1/2} \, \|\|b_{E_i, \delta_i}\|\|_{\omega_T} .$$

Using inverse inequalities (8) and (10) one derives

$$\|\|b_T\|\|_{\omega_T}^2 \; = \; \varepsilon \|\nabla b_T\|_T^2 + \|b_T\|_T^2 \; \lesssim \; \varepsilon h_{\min,T}^{-2} \, |T| + |T| \; \overset{(16)}{\sim} \; \alpha_T^{-2} \cdot |T|.$$

The second inverse equivalences (13) and (14) and the particular choice of $\delta_i \equiv \delta_{E_i}$ from (21) yield

$$\begin{aligned} \|\|b_{E_i, \delta_i}\|\|_{\omega_T}^2 \; &= \; \varepsilon \|\nabla b_{E_i, \delta_i}\|_{\omega_E}^2 + \|b_{E_i, \delta_i}\|_{\omega_E}^2 \\ &\sim \; \delta_i \cdot |T| \cdot \Big( \varepsilon \min\{\delta_i \cdot h_{E,T}, h_{\min,T}\}^{-2} + 1 \Big) \\ &\overset{(21)}{\sim} \; \delta_i \cdot |T| \cdot (1 + \alpha_T^{-2}) \; \sim \; \delta_i \cdot |T| \cdot \alpha_T^{-2}. \end{aligned}$$

Altogether one obtains

$$\|\|v_T\|\|_{\omega_T} \; \lesssim \; |\beta_0| \cdot \alpha_T^{-1} \, |T|^{1/2} + \sum_{i=1}^{m} |\beta_i| \cdot \alpha_T^{-1} \, |T|^{1/2} \; \sim \; \alpha_T^{-1} \, |T|^{1/2} \cdot \|\mathbf{v}\|_{\mathbb{R}^{1+m}}.$$

To bound $\|\|v_T\|\|_{\omega_T}$ from below, apply Lemma 4.2 giving

$$\|\nabla v_T\|_{\omega_T} \; \overset{(26)}{\gtrsim} \; h_{\min,T}^{-1} \cdot \|v_T\|_{\omega_T}$$

$$\text{and} \qquad \|\|v_T\|\|_{\omega_T}^2 \; \gtrsim \; (\varepsilon \, h_{\min,T}^{-2} + 1) \cdot \|v_T\|_{\omega_T}^2 \; \sim \; \alpha_T^{-2} \cdot \|v_T\|_{\omega_T}^2.$$

In the proof of Lemma 4.2 in Appendix A it is shown that

$$\|v_T\|_{\omega_T}^2 \overset{(41)}{\sim} |T| \cdot \sum_{i=0}^m \beta_i^2 \sim |T| \cdot \|\mathbf{v}\|_{\mathbb{R}^{1+m}}^2$$

which completes the lower bound of $\|\|v_T\|\|_{\omega_T}$. Summarizing all results, one ends up with

$$(K_T\mathbf{v}, \mathbf{v}) = \|\|v_T\|\|_{\omega_T}^2 \sim \alpha_T^{-2} |T| \cdot \|\mathbf{v}\|_{\mathbb{R}^{1+m}}^2$$

which immediately yields $\lambda_{\min}(K_T) \sim \lambda_{\max}(K_T) \sim \alpha_T^{-2} |T|$ and the desired assertion $\kappa(K_T) \sim 1$.  $\square$

### 4.4.  **A further, face based local problem error estimator**

With the methodology presented so far one can derive further local problem error estimators. This will be demonstrated here for a *face based* local problem error estimator. Such an estimator can be advantageous when other ingredients of an adaptive algorithm are face based too (*e.g.* the refinement procedure).

We start again with a corresponding residual error estimator. For an arbitrary but fixed face $E$ define the face based residual error estimator and the approximation term by

$$\eta_{\varepsilon,\mathrm{R},E} := \left(\alpha_E^2 \cdot \sum_{T \subset \omega_E} \|r_T\|_T^2 + \varepsilon^{-1/2}\,\alpha_E \cdot \|r_E\|_E^2\right)^{1/2} \tag{37}$$

$$\zeta_{\varepsilon,E} := \alpha_E \cdot \left(\sum_{T \subset \omega_E} \|R_T - r_T\|_T^2\right)^{1/2} + \varepsilon^{-1/4}\alpha_E^{1/2} \cdot \|g - g_h\|_{E \cap \Gamma_\mathrm{N}}, \tag{38}$$

respectively (the norm $\|\cdot\|_{E \cap \Gamma_\mathrm{N}}$ here is to be evaluated only when $E \subset \Gamma_\mathrm{N}$).

Utilizing the techniques and most of the results of [18] one can comparatively easily prove the following residual error estimation.

**Lemma 4.8.** *The error is bounded locally from below for all faces $E$ of $\mathcal{T}_h$ by*

$$\eta_{\varepsilon,\mathrm{R},E} \lesssim \|\|u - u_h\|\|_{\omega_E} + \zeta_{\varepsilon,E}.$$

*The error is bounded globally from above by*

$$\|\|u - u_h\|\| \lesssim m_1(u - u_h, \mathcal{T}_h) \cdot \left(\sum_{E \in \mathcal{T}_h} \eta_{\varepsilon,\mathrm{R},E}^2 + \zeta_{\varepsilon,E}^2\right)^{1/2},$$

*where the sum over $E \in \mathcal{T}_h$ includes interior and boundary faces of the triangulation. Both error bounds are uniform in $\varepsilon$.*

Note that the residual error estimator can be modified such that it contains only the face residual but not the element residuals. Then a very similar result is achieved, *cf.* [19]. Since this modification is not suitable for our subsequent analysis, we omit a detailed description.

The local space associated to a face $E$ is now set to

$$V_E := \mathrm{span}\{b_{E,\delta_E} \text{ if } E \not\subset \Gamma_\mathrm{D} \,,\ b_T \,\forall T \subset \omega_E\},$$

*i.e.* $V_E$ is three dimensional for interior faces $E$. The local problem is: Find $e_E \in V_E$ such that

$$a(e_E, v_E) = \sum_{T \subset \omega_E} \int_T r_T \cdot v_E + \int_E r_E \cdot v_E \qquad \forall\, v_E \in V_E.$$

The *local face error estimator* then becomes

$$\eta_{\varepsilon,\mathrm{D},E} := \|\!|e_E|\!\|_{\omega_E}.$$

Again an alternative, equivalent description of the local problem is possible and advantageous.
Alternative: Find $e_E \in V_E$ such that

$$a(e_E, v_E) = a(u - u_h, v_E) - \sum_{T \subset \omega_E} \int_T (R_T - r_T)\, v_E \ - \int_{E \cap \Gamma_{\mathrm{N}}} (g - g_h)\, v_E$$

holds for all $v_E \in V_E$.

Using the techniques and even some results of the previous analysis of the element based local problem error estimator the following theorem can be shown. Because of the similarities of the proofs we only state the result.

**Theorem 4.9** (Face based local problem error estimator)**.**

*The face based residual error estimator and local problem error estimator are equivalent:*

$$\eta_{\varepsilon,\mathrm{D},E} \ \sim \ \eta_{\varepsilon,\mathrm{R},E} \qquad \forall\, E \in \mathcal{T}_h.$$

*The error is bounded locally from below for all faces $E$ of $\mathcal{T}_h$ by*

$$\eta_{\varepsilon,\mathrm{D},E} \ \leq \ \|\!|u - u_h|\!\|_{\omega_E} \ + \ c \cdot \zeta_{\varepsilon,E}.$$

*The error is bounded globally from above by*

$$\|\!|u - u_h|\!\| \lesssim m_1(u - u_h, \mathcal{T}_h) \cdot \left( \sum_{E \in \mathcal{T}_h} \eta_{\varepsilon,\mathrm{D},E}^2 + \zeta_{\varepsilon,E}^2 \right)^{1/2}.$$

*All relations are uniform in $\varepsilon$.*

## 5. Computational implementation

### 5.1. Difficulties and their solution

It is a major demand that the local problem can be constructed and solved as fast as possible since usually the error estimation is as expensive as the assembly of the global finite element stiffness matrix and the solution process for $u_h$. Therefore one encounters two main problems when applying our error estimator. Both difficulties are related to the computation of the local stiffness matrix $K_T$ which arises from the bilinear form $a(\cdot, \cdot)$, see also Section 4.3.

**Problem 1.** The support of the squeezed face bubble function $b_{E,\delta}$ is not $\omega_E$ but some ($\delta$–dependent) part of it. For example the computation of $a(b_{E,\delta}, \cdot)$ implies a comparatively complicated domain of integration. In particular when evaluating $a(b_{E_1,\delta_{E_1}}, b_{E_2,\delta_{E_2}})$ for two different squeezed face bubble functions, the domain of integration becomes

$$\mathrm{supp}(b_{E_1,\delta_{E_1}}) \cap \mathrm{supp}(b_{E_2,\delta_{E_2}})$$

which might be empty, or a single tetrahedron, or the union of two tetrahedra, depending on $\delta_{E_1}$ and $\delta_{E_2}$ (*cf.* also Figs. 3 and 4). Even to determine and describe the domain of integration is not trivial, save the actual integration.

**Remedy.** We modify the parameter for the squeezed face bubble function to be

$$\tilde{\delta}_E := \min\left\{\frac{1}{4}, \frac{h_{\min,E}}{h_E}, \frac{\sqrt{\varepsilon}}{h_E}\right\} \equiv \min\left\{\frac{1}{4}, \delta_E\right\} \sim \delta_E.$$

Then all results remain valid, only the inequality constants may be slightly worse (but they are still uniform in $\varepsilon$). The main advantage now is that

$$\text{supp}(b_{E_1,\tilde{\delta}_{E_1}}) \cap \text{supp}(b_{E_2,\tilde{\delta}_{E_2}}) = \emptyset.$$

Hence the computation of the modified local matrix $\tilde{K}_T$ is less expensive, as the matrix now contains several zero entries. Even more, the sparsity pattern

$$\tilde{K}_T = \tilde{K}_T^\top = \begin{bmatrix} * & * & * & * & * \\ * & * & 0 & 0 & 0 \\ * & 0 & * & 0 & 0 \\ * & 0 & 0 & * & 0 \\ * & 0 & 0 & 0 & * \end{bmatrix}$$

allows a particularly fast and simple solution of the local problem.

**Problem 2.** The basis functions of $V_T$ are polynomials of a relatively high degree. Hence numerical integration rules to compute $a(\cdot,\cdot)$ are far too expensive and thus unsuitable.

**Remedy.** Instead we propose a direct computation of the integrals involved. The procedure is explained exemplarily for $a(b_T, b_T)$. Using the transformation technique *via* $F_A : \bar{T} \to T$, one obtains

$$a(b_T, b_T) = \varepsilon \int_T (\nabla b_T)^\top \cdot \nabla b_T + \int_T b_T^2$$
$$= 6|T|\,\varepsilon \int_{\bar{T}} (\bar{\nabla} b_{\bar{T}})^\top \cdot A_T^{-1} A_T^{-\top} \cdot \bar{\nabla} b_{\bar{T}} + 6|T| \int_{\bar{T}} b_{\bar{T}}^2.$$

with $b_{\bar{T}}$ being the element bubble function for the standard tetrahedron $\bar{T}$. A straight–forward computation yields

$$\int_{\bar{T}} b_{\bar{T}}^2 = \frac{4096}{155\,925}.$$

In order to obtain the remaining integral, define the matrices

$$M := (m_{ij})_{i,j=1}^3 = A_T^{-1} A_T^{-\top}$$
$$\text{and} \qquad N := (n_{ij})_{i,j=1}^3 = \int_{\bar{T}} \bar{\nabla} b_{\bar{T}} \cdot (\bar{\nabla} b_{\bar{T}})^\top, \qquad i.e. \quad n_{ij} = \int_{\bar{T}} \bar{\partial}_{\bar{x}_1} b_{\bar{T}} \cdot \bar{\partial}_{\bar{x}_i} b_{\bar{T}}$$

and observe that

$$\int_{\bar{T}} (\bar{\nabla} b_{\bar{T}})^\top \cdot A_T^{-1} A_T^{-\top} \cdot \bar{\nabla} b_{\bar{T}} = N : M = \sum_{i,j=1}^3 n_{ij} \cdot m_{ij}.$$

The matrices $A_T$ and $M$ are determined by the geometry of $T$ whereas $N$ can be computed directly giving

$$N = \frac{2048}{2835} \cdot \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}.$$

Hence one only has to determine $A_T$ and $|T|$ and perform the aforementioned operations to obtain $a(b_T, b_T)$.

The remaining values $a(\cdot, \cdot)$ are computed similarly, thus a brief description of the result should suffice.

$a(b_T, b_{E,\delta})$: Use an affine linear transformation $F_{\tilde{A}} : \bar{T} \to T$ such that $\bar{E} = F_{\tilde{A}}^{-1}(E)$ lies in the $\bar{x}_1 \bar{x}_2$ plane. The parameter $\delta \in (0, 1]$ of the squeezed face bubble function can be arbitrary here. Then

$$\int_{\bar{T}} b_{\bar{T}} \cdot b_{\bar{E},\delta} = \frac{4}{4725} \delta^2 (-2\delta^3 + 15\delta^2 - 42\delta + 47)$$

$$\int_{\bar{T}} \bar{\nabla} b_{\bar{T}} \cdot (\bar{\nabla} b_{\bar{E},\delta})^\top = \begin{bmatrix} 2c & c & c \\ c & 2c & c \\ c & c & d \end{bmatrix}$$

$$\text{with} \qquad c := -\frac{4}{35}\delta^2(\delta - 4) \qquad \text{and} \qquad d := \frac{8}{105}(\delta - 1)(\delta^2 - 7\delta + 18).$$

$a(b_{E,\delta}, b_{E,\delta})$: Utilize the same transformation $F_{\tilde{A}}$ as before which implies for arbitrary $\delta \in (0, 1]$

$$\int_{\bar{T}} b_{\bar{E},\delta}^2 = \frac{9}{560}\delta$$

$$\int_{\bar{T}} \bar{\nabla} b_{\bar{E},\delta} \cdot (\bar{\nabla} b_{\bar{E},\delta})^\top = \begin{bmatrix} 2c & c & c \\ c & 2c & c \\ c & c & d \end{bmatrix}$$

$$\text{with} \qquad c := \frac{81}{280}\delta \qquad \text{and} \qquad d := \frac{27}{140}(\delta + 2/\delta).$$

$a(b_{E_1,\delta_1}, b_{E_2,\delta_2})$: Since we propose to use $\delta_i := \tilde{\delta}_{E_i} \le 1/4$, the supports of both squeezed face bubble functions are distinct, thus $a(b_{E_1,\delta_1}, b_{E_2,\delta_2}) = 0$.

Collecting all the previous results, the local stiffness matrix $K_T$ can now be assembled. The right–hand side is computed similarly (actually, the procedure is even simpler since the integrals do not involve derivatives). In the next paragraphs we show that this direct computation of the local problem is indeed much cheaper than the numerical integration rules.

## 5.2. Computational effort

The comparison will not investigate every detail and every possible optimization, as the difference between both approaches will turn out to be overwhelming. Even more, a precise operation count would be computer dependent. For example present processors may be able to combine one multiplication and one addition to a single operation.

### 5.2.1. *Computational effort for direct computation*

The suggested approach utilizes four different transformations $F_{\tilde{A}}$ on the element $T$ and four transformations on each neighbouring tetrahedron (*i.e.* to compute $b_{E,\delta}$ on this neighbour). Hence $|\det \tilde{A}|$ has to be computed five times, and $\tilde{A}^{-1}$ and $\tilde{A}^{-1} \cdot \tilde{A}^{-\top}$ are to be computed eight times. The operation count is roughly

| Operation | Operation count | Total |
|---|---|---|
| $\det \tilde{A}$ | $5 \times (\ 9\times\ \ 5+) =$ | $(\ \ 45\times\ \ 25+)$ |
| $\tilde{A}^{-1}$ | $8 \times (19\times\ \ 9+) =$ | $(152\times\ \ 72+)$ |
| $\tilde{A}^{-1} \cdot \tilde{A}^{-\top}$ (symmetric) | $8 \times (18\times\ 12+) =$ | $(144\times\ \ 96+)$ |
| | | $\Sigma : (341\times\ 193+)$ |

where $5 \times (9\times\ 5+)$ stands for 9 multiplications and 5 additions which are performed five times.

Next, the values of $\int_{\bar{T}} b_{\bar{T}}^2$, $\int_{\bar{T}} \bar{\nabla} b_{\bar{T}} \cdot (\bar{\nabla} b_{\bar{T}})^\top$ etc are determined. Some of these values (which contain $b_{E,\delta}$) depend on $\delta$. The computational effort is roughly

| | | |
|---|---|---|
| $\int_{\bar{T}} b_{\bar{T}}^2$ | | $(\ 0\times\ \ 0+)$ |
| $\int_{\bar{T}} b_{\bar{T}} \cdot b_{\bar{E},\delta}$ | $4 \times (8\times\ 3+) =$ | $(32\times\ 12+)$ |
| $\int_{\bar{T}} b_{\bar{E},\delta}^2$ | $4 \times (1\times\ 0+) =$ | $(\ 4\times\ \ 0+)$ |
| $\int_{\bar{T}} \bar{\nabla} b_{\bar{T}} \cdot (\bar{\nabla} b_{\bar{T}})^\top$ | | $(\ 0\times\ \ 0+)$ |
| $\int_{\bar{T}} \bar{\nabla} b_{\bar{T}} \cdot (\bar{\nabla} b_{\bar{E},\delta})^\top$ | $4 \times (6\times\ 4+) =$ | $(24\times\ 16+)$ |
| $\int_{\bar{T}} \bar{\nabla} b_{\bar{E},\delta} \cdot (\bar{\nabla} b_{\bar{E},\delta})^\top$ | $4 \times (1\times\ 1+) =$ | $(16\times\ \ 4+)$ |
| | | $\Sigma : (76\times\ \ 32+)$ |

Subsequently $N : M$ is to be determined, with $N$, $M$ being symmetric matrices. The computational effort is approximately $(7\times\ 5+)$ which has to be repeated 9 times (*i.e.* once for each matrix entry of $K_T$). The final value of $a(\cdot,\cdot)$ is obtained by adding both sub–integrals and multiplying it by $|\det\tilde{A}|$. This adds $9 \times (1\times\ 1+)$.

Summarizing all results, the total effort required to assemble the local stiffness matrix is approximately

$$(341\times\ 193+)\ +\ (76\times\ 32+)\ +\ (63\times\ 45+)\ +\ (9\times\ 9+)\ =\ \mathbf{(489\times\ \ 279+)}\ .$$

### 5.2.2. *Computational effort for numerical integration*

Here we will exemplarily investigate $a(b_T, b_T) = \int_T b_T^2 + \varepsilon \int_T (\nabla b_T)^2$. Computation by means of numerical integration is based on

$$\int_T b_T^2 = 6|T| \int_{\bar{T}} b_{\bar{T}}^2 \ \approx\ 6|T| \cdot \sum_i \omega_i \cdot b_{\bar{T}}^2(\bar{x}_i)$$

$$\int_T (\nabla b_T)^2 = 6|T| \int_{\bar{T}} (\bar{\nabla} b_{\bar{T}})^\top \cdot A_T^{-1} A_T^{-\top} \cdot \bar{\nabla} b_{\bar{T}}$$

$$\approx 6|T| \cdot \sum_i \omega_i \cdot \left( (\bar{\nabla} b_{\bar{T}})^\top \cdot A_T^{-1} A_T^{-\top} \cdot \bar{\nabla} b_{\bar{T}} \right)(\bar{x}_i)$$

where $(\omega_i, \bar{x}_i)_i$ denotes some numerical integration rule for the standard tetrahedron $\bar{T}$ with weights $\omega_i$ and evaluation points $\bar{x}_i$. Exactly as for the direct computation above, one requires the matrices $A_T$ and $A_T^{-1}$ (computational effort is $(19\times\ \ 9+)$) as well as $6|T| = |\det A_T|$ (which leads to $(9\times\ 5+)$).

Consider $\int \bar{T} b_{\bar{T}}^2$ next. Since $b_{\bar{T}} \in \mathbb{P}^4(\bar{T})$ one requires an integration rule which is exact for $\mathbb{P}^8(\bar{T})$. The simplest rule that we know of involves 43 evaluation points [9]. The evaluation of

$$\omega_i \cdot b_{\bar{T}}^2(\bar{x}_i) \ =\ \omega_i \cdot (256 \cdot \lambda_1 \lambda_2 \lambda_3 \lambda_4)^2(\bar{x}_i)$$

at a single point $\bar{x}_i$ requires $(6\times\ 3+)$. Hence the total amount for $\int_T b_T^2$ is about $43 \times (6\times\ 4+) = (258\times\ 172+)$. Similarly $\int_T (\nabla b_T)^2$ is investigated where an integration rule with 24 evaluation points suffices [12]. After some consideration one ends up with an effort of roughly $(480\times\ 384+)$. Thus the total effort to compute $a(b_T, b_T)$ amounts approximately to $(750\times\ 550+)$. We expect the other eight scalar products $a(\cdot,\cdot)$ to be cheaper because of the smaller polynomial degree. Then, however, the computational domain involves $T_{E,\delta}$ which requires more considerations. Altogether the effort to compute the whole matrix of the local problem is likely to be of order

$$(\mathcal{O}(5000)\times\ \ \mathcal{O}(4000)+)\ .$$

Even without a precise calculation of the computational effort for the numerical integration it is absolutely clear that this approach is far too expensive. Equivalently, if the computational effort should be of the same size $(\mathcal{O}(500)\times \mathcal{O}(300)+)$ as for the direct computation then only one tenth of the required evaluation points can

be used. This would render the integrals and, subsequently, the matrix to be very inaccurate and thus useless. Hence direct computation of the local problem is a must.

Furthermore we note that our procedure for the direct computation is very similar to the computation of the local problem for the *Poisson equation, cf.* [16]. The computational effort is roughly the same, *i.e.* the singularly perturbed character of our differential equation here is no disadvantage.

## 6. Numerical experiments

Here we investigate the performance of the local problem error estimator $\eta_{\varepsilon,\mathrm{D},T}$ of (24) by means of numerical experiments. We utilize a model problem which has already been applied in [18] to analyse the element based residual error estimator $\eta_{\varepsilon,\mathrm{R},T}$, and which has been employed in [19] to investigate a face based residual error estimator. Thus even the interesting comparison between different types of error estimators is possible.

Let us consider the 3D model problem

$$-\varepsilon\Delta u + u = 0 \quad \text{in } \Omega := (0,1)^3, \qquad u = u_0 \quad \text{on } \Gamma_\mathrm{D} := \partial\Omega$$

where the perturbation parameter is set to $\varepsilon = 10^{-4}$. The results for $\varepsilon = 10^{-8}$ are omitted because they are very similar; thus they confirm the $\varepsilon$–robustness of our results. Next, prescribe the exact solution

$$u = e^{-x/\sqrt{\varepsilon}} + e^{-y/\sqrt{\varepsilon}} + e^{-z/\sqrt{\varepsilon}}.$$

which displays typical boundary layers along the planes $x = 0$, $y = 0$, and $z = 0$. The Dirichlet boundary data $u_0$ are chosen accordingly.

We apply the finite element method with a sequence of meshes $\mathcal{T}_k$, each of which is the tensor product of three one–dimensional Bakhvalov–like meshes [7] with $2^k$ intervals in [0,1], $k = 1 \ldots 6$. To describe the 1D nodal distribution properly, denote the transition point of the boundary layer by $\tau := \sqrt{\varepsilon}|\ln\sqrt{\varepsilon}|$. Then $2^{k-1}$ nodes are *exponentially* distributed in the boundary layer interval $[0, \tau]$ whereas the remaining interval $[\tau, 1]$ is divided into $2^{k-1}$ *equidistant* intervals, *cf.* Figure 5. More precisely, the (1D) nodal coordinate of the $m$-th node is

$$x_m := \begin{cases} -\beta\sqrt{\varepsilon}\ln\left[1 - \dfrac{m}{2^{k-1}}(1 - e^{-\tau/\beta/\sqrt{\varepsilon}})\right] & \text{for } m = 0 \ldots 2^{k-1}, \beta = 3/2 \\ \tau + (1-\tau)\cdot\left(\dfrac{m}{2^{k-1}} - 1\right) & \text{for } m = 2^{k-1} + 1 \ldots 2^k. \end{cases}$$

Note that the original (1D) Bakhvalov mesh utilizes a slightly different transition point $\tau$. Furthermore we do not know whether these tensor product type meshes are optimal (which, of course, also depends on the optimality criterion).
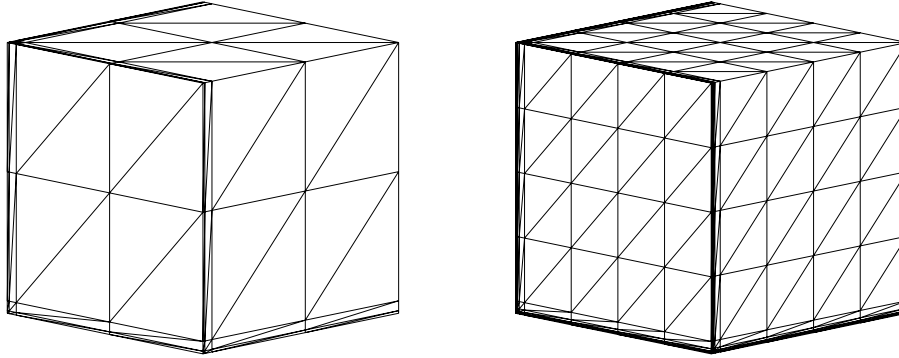
FIGURE 5. Mesh $\mathcal{T}_2$ – Mesh $\mathcal{T}_3$.

The first table below presents some information about the meshes and their maximum aspect ratio. The last two columns give the exact value of the matching function $m_1(u - u_h, \mathcal{T}_k)$ as well as its approximation $m_1^R(u_h, \mathcal{T}_k)$ from (20).

| Mesh $\mathcal{T}_k$ | # Elements | Aspect ratio | $m_1(u - u_h, \mathcal{T}_k)$ | $m_1^R(u_h, \mathcal{T}_k)$ |
|---|---|---|---|---|
| 1 | 48 | 29.4 | 1.55 | 1.68 |
| 2 | 384 | 69.5 | 1.62 | 1.52 |
| 3 | 3 072 | 82.6 | 1.69 | 1.69 |
| 4 | 24 576 | 88.6 | 1.88 | 1.86 |
| 5 | 196 608 | 91.5 | 2.37 | 2.03 |
| 6 | 1 572 864 | 92.9 | 3.04 | 2.29 |

Since the size of $m_1$ is comparatively small and grows only mildly, the chosen meshes discretize the problem sufficiently well. Additionally the approximation $m_1^R$ is satisfactorily close to the exact value. Hence the matching function and its approximation are useful tools for the theoretical analysis as well as for assessing the mesh quality in numerical computations. This topic has already been discussed for the Poisson equation in [14].

Next our main analytical results are to be confirmed numerically, namely the error bounds of Theorem 4.4. Therefore we present the ratios of left–hand side and right–hand side of (35) and (34), respectively, in the table below. These ratios have to be bounded from above (*cf.* Th. 4.4).

| Mesh $\mathcal{T}_k$ | $\|\|u - u_h\|\|$ | $\dfrac{\|\|u - u_h\|\|}{m_1 \cdot \left(\eta_{\varepsilon,\mathrm{D}}^2 + \zeta_\varepsilon^2\right)^{1/2}}$ | $\displaystyle\max_{T \in \mathcal{T}_k} \dfrac{\eta_{\varepsilon,\mathrm{D},T}}{\|\|u - u_h\|\|_{\omega_T} + \zeta_{\varepsilon,T}}$ |
|---|---|---|---|
| 1 | 0.154E + 0 | 0.517 | 0.308 |
| 2 | 0.536E − 1 | 0.399 | 0.321 |
| 3 | 0.229E − 1 | 0.413 | 0.415 |
| 4 | 0.110E − 1 | 0.437 | 0.480 |
| 5 | 0.553E − 2 | 0.396 | 0.507 |
| 6 | 0.282E − 2 | 0.330 | 0.511 |

Start with the second column which yields a convergence rate of the error $\|\|u - u_h\|\|$ of approximately $N^{-0.324}$, with $N$ being the number of elements. This is almost the optimal rate of $N^{-1/3}$ which indicates that the meshes under consideration discretize the singular problem well. Next, the ratios of the third and fourth column are related to the upper and lower error bound, respectively. These ratios are bounded from above and thus confirm

the predictions of Theorem 4.4. Note that from a practical point of view the moderately decreasing values of the upper error bound (third column) imply that the error is increasingly overestimated.

In the last table we examine the equivalence of the local problem error estimator and the residual error estimator, as described in Theorem 4.3. Again we compute the ratios related to (28) and (29). Since all values are bounded from above, this impressively underpins our analytical results.

| Mesh $\mathcal{T}_k$ | $\displaystyle\max_{T\in\mathcal{T}_k} \frac{\eta_{\varepsilon,\mathrm{D},T}}{\left(\sum_{T'\subset\omega_T}\eta_{\varepsilon,\mathrm{R},T'}^2\right)^{1/2}}$ | $\displaystyle\max_{T\in\mathcal{T}_k} \frac{\eta_{\varepsilon,\mathrm{R},T}}{\left(\sum_{T'\subset\omega_T}\eta_{\varepsilon,\mathrm{D},T'}^2\right)^{1/2}}$ |
|---|---|---|
| 1 | 0.302 | 1.556 |
| 2 | 0.443 | 4.278 |
| 3 | 0.507 | 4.956 |
| 4 | 0.386 | 4.851 |
| 5 | 0.275 | 4.708 |
| 6 | 0.202 | 4.770 |

Since the same numerical example has been considered for the residual error estimator of [18] we can easily compare both estimator. Qualitatively both estimators behave similarly whereas from a quantitative viewpoint one observes roughly $\eta_{\varepsilon,\mathrm{R}} \approx 4 \cdot \eta_{\varepsilon,\mathrm{D}}$. Furthermore the residual error estimator $\eta_{\varepsilon,\mathrm{R}}$ overestimates the true error more than the local problem error estimator $\eta_{\varepsilon,\mathrm{D}}$ does. This indeed can be expected since the derivation of $\eta_{\varepsilon,\mathrm{R}}$ requires more intermediate steps (such as interpolation estimates and Cauchy Schwarz inequalities).

## 7. SUMMARY

We have considered a singularly perturbed reaction–diffusion problem and proposed a new error estimator that can be applied to *anisotropic* finite element meshes. The rigorous analysis confirms that the error estimation is uniform in the small perturbation parameter. Furthermore tight error bounds are obtained provided the anisotropic mesh is chosen according to the anisotropy of the solution. Thus reliable and efficient error estimation is possible on anisotropic meshes.

Then a stable basis of the local problem has been derived and an additional, face oriented local problem error estimator has been proposed. Finally implementational aspects have been discussed and analysed. A numerical experiment complements the theory.

## Appendix A. PROOF OF LEMMA 4.2

First we state Lemma 4.2 again.

**Lemma 4.2.** *The following relations below hold for all $v \in V_T$.*

$$\|v\|_{\omega_T} \lesssim h_{\min,T} \cdot \|\nabla v\|_{\omega_T} \tag{26}$$

$$\|v\|_E \lesssim h_E^{-1/2}\, \delta_E^{-1/2} \cdot \min\{h_{\min,T}, \delta_E\, h_E\} \cdot \|\nabla v\|_{\omega_T} \qquad \forall E \subset \partial T. \tag{27}$$

*The inequalities are uniform in the squeezing parameters $\delta_E \in (0,1]$ which define the space $V_T$.*

*If $T$ has at least two Neumann boundary faces then the constants in (26), (27) can depend on the shape of the Neumann boundary (but do not depend on the triangulation $\mathcal{T}_h$ nor on $T$).*

*Proof.* The proof here utilizes some key ideas that were already applied in [13, Lem. 3.5] and [16]. Our exposition here requires several non–trivial extensions which are due to the singularly perturbed problem, and the use of *squeezed* face bubble functions in particular. In order to facilitate the understanding of the proof, each major step will be given a distinctive name.

Set $T_0 := T$ and enumerate the remaining tetrahedra of $\omega_T \setminus T$ by $T_1 \ldots T_k$. If $T$ has boundary faces then $k < 4$. The faces of $T$ are denoted accordingly by $E_i := T_i \cap T$.

**Transformation.** In order to prove (26) consider the tetrahedron $T_i$ and rewrite $\|v\|_{T_i}$ by means of the transformation $A_{T_i}$, and $\|\nabla v\|_{T_i}$ *via* the transformations $C_{T_i}, H_{T_i}$. Utilizing $|T| \sim |T_i|$, and with certain abbreviations given below this yields

$$\|v\|_{T_i}^2 = 6|T_i| \cdot \|\bar{v}\|_{\hat{T}_i}^2 \ \sim \ |T| \cdot r_i$$

$$\text{and} \qquad \|v\|_{\omega_T}^2 \sim |T| \cdot \sum_{i=0}^{k} r_i \ = \ |T| \cdot r \tag{39}$$

where we have introduced

$$r_i := \|\bar{v}\|_{\hat{T}_i}^2 \ \geq 0 \qquad \text{and} \qquad r := \sum_{i=0}^{k} r_i \ \geq 0.$$

Using the matrix $H_{T_i}$ from (4) in conjunction with (5), the term $\|\nabla v\|_{\omega_T}$ is transformed similarly to give

$$\|\nabla v\|_{T_i}^2 \overset{(5)}{=} \|H_{T_i}^{-1} C_{T_i}^\top \nabla v\|_{T_i}^2 \ = \ 6|T_i| \cdot \|H_{T_i}^{-1} \cdot \hat{\nabla} \hat{v}\|_{\hat{T}_i}^2$$

$$\text{which implies} \qquad \|\nabla v\|_{\omega_T}^2 \ = \ \sum_{i=0}^{k} \|\nabla v\|_{T_i}^2 \ = \ \sum_{i=0}^{k} 6|T_i| \cdot \|H_{T_i}^{-1} \cdot \hat{\nabla} \hat{v}\|_{\hat{T}_i}^2$$

$$= \ 6 \cdot \sum_{i=0}^{k} |T_i| \cdot h_{\min,T_i}^{-2} \cdot \big\|\mathrm{diag}(\gamma_{1,i}, \gamma_{2,i}, 1) \cdot \hat{\nabla} \hat{v}\big\|_{\hat{T}_i}^2$$

$$\sim \ h_{\min,T}^{-2} \cdot |T| \cdot \sum_{i=0}^{k} s_i \ = \ h_{\min,T}^{-2} \cdot |T| \cdot s \tag{40}$$

$$\text{with} \quad \gamma_{1,i} := h_{\min,T_i}/h_{1,T_i} \qquad \text{and} \qquad \gamma_{2,i} := h_{\min,T_i}/h_{2,T_i}$$

$$s_i := \big\|\mathrm{diag}(\gamma_{1,i}, \gamma_{2,i}, 1) \cdot \hat{\nabla} \hat{v}\big\|_{\hat{T}_i}^2 \ \geq 0 \quad \text{and} \quad s := \sum_{i=0}^{k} s_i \ \geq 0.$$

A rough outline of the proof is as follows. Realize first that $r$ and $s$ depend on various variables (*e.g.* on the geometry of $T$, the parameters $\delta_i := \delta_{E_i}$ etc.) Then consider $r, s$ over some *compact* set of variables. Since both terms turn out to be continuous, one can investigate their maxima and minima which eventually provide the assertion.

**Extend domain of definition to a compact set.** Let us start with the case where $T$ has no Neumann boundary face. Assume further that $T$ has $m$ interior faces and $4 - m$ Dirichlet faces. The local space $V_T$ is spanned by $b_T$ and $b_{E_i,\delta_i}$, $i = 1 \ldots m$. For our purpose we utilize an expansion of $v \in V_T$ where the squeezed face bubble functions are additionally scaled by $\delta_i^{-1/2}$, namely

$$v \ = \ \beta_0 \cdot b_T + \sum_{i=1}^{m} \beta_i \cdot \delta_i^{-1/2} b_{E_i,\delta_i} \qquad \beta_i \in \mathbb{R}.$$

Without loss of generality assume $v \not\equiv 0$ and $\sum_{i=0}^{m} \beta_i^2 = 1$. After transformation *via* $A_{T_i}$ the representation of $\bar{v}$ becomes

$$\bar{v}\big|_{\bar{T}_i} = v\big|_{T_i} \circ F_{A_{T_i}} = \begin{cases} \beta_0 \cdot b_{\bar{T}} + \sum_{i=1}^{m} \beta_i \cdot \delta_i^{-1/2} b_{\bar{E}_i, \delta_i} & \text{for } i = 0 \ (i.e. \text{ on } \bar{T}) \\ \beta_i \cdot \delta_i^{-1/2} b_{\bar{E}_i, \delta_i} & \text{for } i = 1 \ldots m \ (i.e. \text{ on } \bar{T}_i). \end{cases}$$

Hence $\bar{v}$ depends on $\beta_0 \ldots \beta_m$ and $\delta_1 \ldots \delta_m$, $\delta_i \in (0, 1]$. Note further that $\delta_i$ influences $r_0 = \|\bar{v}\|_{\bar{T}}^2$ and $r_i = \|\bar{v}\|_{\bar{T}_i}^2$ but not the other values $r_j$.

Next $\bar{v}$ is to be considered over a *compact* set. Thus introduce

$$B := \left\{ (\beta_0, \ldots, \beta_m) \ : \ \sum_{i=0}^{m} \beta_i^2 = 1 \right\} \quad \text{and} \quad D := \left\{ (\delta_1, \ldots, \delta_m) \ : \ \delta_i \in [0, 1] \ \forall \, i \right\}.$$

The case $\delta_i = 0$ requires additional consideration. While the function $\delta_i^{-1/2} b_{E_i, \delta_i}$ has a well–defined meaning for $\delta_i \in (0, 1]$, this is no longer true for $\delta_i \to 0$. Then $\operatorname{supp}(b_{E_i, \delta_i}|_{T_i}) = T_{i, E_i, \delta_i}$ degenerates, and $\delta_i^{-1/2} \to \infty$. Therefore the value of $r_j = \|\bar{v}\|_{\bar{T}_j}^2$ for $\delta_i = 0$ is defined as the limit for $\delta_i \to 0$:

$$r_j(\delta_i = 0) := \lim_{\delta_i \to 0} r_j(\delta_i), \quad j = 0, i.$$

This limit is well–defined since the vanishing support $T_{i, E_i, \delta_i}$ of the squeezed face bubble function and its scaling factor $\delta_i^{-1/2}$ are exactly balanced. For the outer tetrahedra this can be easily seen by utilizing two transformations, namely *via* $F_{A_{T_i}} : \bar{T} \to T_i$ and *via* $F_{T_i, E_i, \delta_i}^{-1} : T_{i, E_i, \delta_i} \to \bar{T}_i$. By using $|\det(A_{T_i})| = 6|T|$ and $|\det(F_{T_i, E_i, \delta_i})| = 6\delta_i|T|$ one obtains

$$\begin{aligned} r_i &= \beta_i^2 \, \|\delta_i^{-1/2} b_{\bar{E}_i, \delta_i}\|_{\bar{T}_i}^2 \\ &\overset{A_{T_i}}{=} \beta_i^2 \, \delta_i^{-1} \cdot (6|T|)^{-1} \|b_{E_i, \delta_i}\|_{T_i}^2 = \beta_i^2 \, \delta_i^{-1} \cdot (6|T|)^{-1} \|b_{E_i, \delta_i}\|_{T_{i, E_i, \delta_i}}^2 \\ &\overset{F_{T_i, E_i, \delta_i}^{-1}}{=} \beta_i^2 \, \delta_i^{-1} \cdot (6|T|)^{-1} \cdot 6\delta_i|T| \cdot \|b_{\bar{E}_i}\|_{\bar{T}_i}^2 = \frac{9}{560} \beta_i^2 \ \sim \ \beta_i^2 \end{aligned}$$

since $b_{\bar{E}_i}$ is a standard face bubble function on the standard tetrahedron $\bar{T}_i$. Therefore $\lim_{\delta_i \to 0} r_i$ exists,

$$\lim_{\delta_i \to 0} r_i = \frac{9}{560} \beta_i^2.$$

For $r_0$ proceed similarly.

**Consider Maximum and Minimum.** As a consequence we can consider $r_i$ and $r$ on $B \times D$, and $r_i, r$ vary continuously over that compact set. Therefore $r$ attains its maximum and minimum. To show that this minimum is positive, assume the contrary which implies $r_i = 0$ for all $i = 0 \ldots m$. On the outer tetrahedra $T_i$, $i = 1 \ldots m$, proceed exactly as in the last paragraph to obtain

$$0 = r_i = \frac{9}{560} \beta_i^2$$

which implies $\beta_i = 0$, $i = 1 \ldots m$. On the main tetrahedron $T$ then $\bar{v}$ is reduced to $\bar{v} = \beta_0 \, b_{\bar{T}}$ giving

$$0 = r_0 = \|\bar{v}\|_{\bar{T}}^2 = \frac{4096}{155\,925} \beta_0^2 \ \sim \ \beta_0^2$$

and $\beta_0 = 0$ too. This contradicts $\sum_{i=0}^m \beta_i^2 = 1$, hence

$$\min_{B \times D} r > 0.$$

Together with $\max_{B \times D} r \sim 1$ we obtain

$$r \sim 1$$

or, equivalently,

$$\|v\|^2_{\omega_T} \sim |T| \cdot \sum_{i=0}^m \beta_i^2. \tag{41}$$

**Investigation of $s$.** The investigation of $s$ and $s_i$ relies on the same basic ideas as before. The details however are much more technical because derivatives are involved (*i.e.* $\nabla v$) and the transformation $C_{T_i}$ is applied.

Consider $s_i = \|\mathrm{diag}(\gamma_{1,i}, \gamma_{2,i}, 1) \cdot \hat{\nabla}\hat{v}\|^2_{\hat{T}_i} \geq 0$ which depends on $\hat{T}_i$, $\gamma_{j,i}$ and $\hat{v}$. The restrictions on $\hat{T}_i$ and $T_i$ imply $0 < \gamma_{1,i}, \gamma_{2,i} \leq 1$ and, for the nodal coordinates of $\hat{T}_i$, $0 < \hat{x}_{2,T_i} \leq 1/2$, $0 < \hat{x}_{3,T_i} < 1$, $-1 < \hat{y}_{3,T_i} < 1$. Similar as before we omit the actual meaning that stands behind $s_i$, and view it instead as a purely analytical term that depends on $\hat{x}_{j,T_i}(j = 2,3)$, $\hat{y}_{3,T_i}$, $\gamma_{j,i}(j = 1,2)$, $\delta_j$ and $\beta_j$. Next consider $s_i$ over the compact set $X_i \times G_i \times D \times B$, with

$$X_i := \left\{ (\hat{x}_{2,T_i}, \hat{x}_{3,T_i}, \hat{y}_{3,T_i}) : 0 \leq \hat{x}_{2,T_i} \leq \frac{1}{2}, \, 0 \leq \hat{x}_{3,T_i} \leq 1, \, -1 \leq \hat{y}_{3,T_i} \leq 1 \right\},$$

$$G_i := \left\{ (\gamma_{1,i}, \gamma_{2,i}) : \quad 0 \leq \gamma_{1,i}, \gamma_{2,i} \leq 1 \right\}.$$

It is obvious that $s_i$ is continuous on $X_i, G_i, B$ and for $\delta_i \in (0,1]$. Note again that $\delta_i$ influences only $s_0$ and $s_i$. The only cause for discontinuity of $s_i$ is $\delta_i \to 0$ which may lead to $s_i \to \infty$ (because $\delta_i^{-1/2} \to \infty$ and $|\hat{\nabla}\hat{b}_{E_i,\delta_i}| \to \infty$, see below). Nevertheless such a discontinuity does not disturb our analysis since we want to bound $s_i$ from below. For a precise investigation we define again

$$s_j(\delta_i = 0) := \lim_{\delta_i \to 0} s_j(\delta_i)$$

and consider then the term $\min\{1, s_i\}$ which is continuous for $\delta_i \in [0,1]$.

Since $s = \sum_{i=0}^m s_i$, this term is continuous as well, and it attains its minimum over the compact set

$$K := \bigtimes_{i=0}^m X_i \quad \times \quad \bigtimes_{i=0}^m G_i \quad \times \quad D \quad \times \quad B.$$

In order to show that this minimum is positive assume the contrary, namely $s = s_i = 0$ for all $i = 0 \ldots m$. Start with any of the outer tetrahedra $T_i, i = 1 \ldots m$. The representation of $\hat{v}$ there is $\hat{v}|_{\hat{T}_i} = \beta_i \delta_i^{-1/2} \cdot \hat{b}_{E_i,\delta_i}|_{\hat{T}_i}$. Then

$$0 = s_i = \|\mathrm{diag}(\gamma_{1,i}, \gamma_{2,i}, 1) \cdot \hat{\nabla}\hat{v}\|^2_{\hat{T}_i}$$
$$\geq \|\mathbf{e}_3^\top \cdot \hat{\nabla}\hat{v}\|^2_{\hat{T}_i} = \beta_i^2 \, \delta_i^{-1} \, \|\mathbf{e}_3^\top \cdot \hat{\nabla}\hat{b}_{E_i,\delta_i}\|^2_{\hat{T}_i},$$

with $\mathbf{e}_3 := (0,0,1)^\top$. The latter norm is analysed similarly as for $r_i$ by using two transformations *via* $F_{C_{T_i}}$ : $\hat{T} \to T_i$ and *via* $F_{T_i,E_i,\delta_i}^{-1} : T_{i,E_i,\delta_i} \to \bar{T}_i$. In contrast to $r_i$, however, we cannot evaluate $s_i$ exactly but bound it

instead. From $|\det(C_{T_i})| = 6|T|$ and $|\det(F_{T_i,E_i,\delta_i})| = 6\delta_i|T|$ one derives

$$
\begin{aligned}
\delta_i^{-1}\,\|\mathbf{e}_3^\top \cdot \hat{\nabla}\hat{b}_{E_i,\delta_i}\|_{\hat{T}_i}^2 \;&\overset{C_{T_i}}{=}\; \delta_i^{-1}\cdot(6|T|)^{-1}\|\mathbf{e}_3^\top\cdot C_{T_i}^\top\nabla b_{E_i,\delta_i}\|_{T_i}^2 \\
&=\; \delta_i^{-1}\cdot(6|T|)^{-1}\|\mathbf{e}_3^\top\cdot C_{T_i}^\top\nabla b_{E_i,\delta_i}\|_{T_{i,E_i,\delta_i}}^2 \\
&\overset{F_{T_i,E_i,\delta_i}^{-1}}{=}\; \delta_i^{-1}\cdot(6|T|)^{-1}\cdot 6\delta_i|T|\cdot\|\mathbf{e}_3^\top\cdot C_{T_i}^\top F_{T_i,E_i,\delta_i}^{-\top}\bar{\nabla}b_{\bar{E}_i}\|_{\bar{T}_i}^2 \\
&=\; \|(F_{T_i,E_i,\delta_i}^{-1}C_{T_i}\cdot\mathbf{e}_3)^\top\cdot\bar{\nabla}b_{\bar{E}_i}\|_{\bar{T}_i}^2.
\end{aligned}
$$

Recalling the definition of $C_{T_i}$ from (4) yields $C_{T_i}\mathbf{e}_3 = \mathbf{p}_{3,T_i}$ which is a vector from a vertice to the opposite face in the tetrahedron $T_i$, see Figure 1. Hence $F_{T_i,E_i,\delta_i}^{-1}C_{T_i}\cdot\mathbf{e}_3$ is a vector from a vertice to the opposite face in the tetrahedron $F_{T_i,E_i,\delta_i}^{-1}T_i$. Using $T_i \supset T_{i,E_i,\delta_i}$ one obtains $F_{T_i,E_i,\delta_i}^{-1}T_i \supset F_{T_i,E_i,\delta_i}^{-1}T_{i,E_i,\delta_i} \equiv \bar{T}_i$ and thus

$$
|F_{T_i,E_i,\delta_i}^{-1}C_{T_i}\cdot\mathbf{e}_3|_{\mathbb{R}^3} \;\geq\; \varrho(\bar{T}_i) \;=\; \frac{1}{3+\sqrt{3}} \;\sim\; 1
$$

where $\varrho(\bar{T}_i)$ denotes the diameter of the inscribed ball of $\bar{T}_i$. Then

$$
\|(F_{T_i,E_i,\delta_i}^{-1}C_{T_i}\cdot\mathbf{e}_3)^\top\cdot\bar{\nabla}b_{\bar{E}_i}\|_{\bar{T}_i}^2 \;\gtrsim\; \min_{|\mathbf{q}|_{\mathbb{R}^3}=1}\|\mathbf{q}^\top\cdot\bar{\nabla}b_{\bar{E}_i}\|_{\bar{T}_i}^2 \;=\; 81/280 \;\sim\; 1.
$$

Summarizing the previous results, we end up with

$$
0 = s_i \;\geq\; \beta_i^2\,\delta_i^{-1}\,\|\mathbf{e}_3^\top\cdot\hat{\nabla}\hat{b}_{E_i,\delta_i}\|_{\hat{T}_i}^2 \;\gtrsim\; \beta_i^2.
$$

This holds for $\delta_i \in (0,1]$ and therefore also for the limit $\delta_i = 0$. Hence one concludes

$$
\beta_i \;=\; 0 \qquad \forall i = 1\ldots m.
$$

Next consider the main tetrahedron $T$ where $v$ is now reduced to $v = \beta_0\cdot b_T$. Then

$$
0 = s_0 \;\geq\; \beta_0^2\,\|\mathbf{e}_3^\top\cdot\hat{\nabla}b_{\hat{T}}\|_{\hat{T}}^2
$$

immediately implies $\beta_0 = 0$ which contradicts the assumption $\sum_{i=1}^m \beta_i^2 = 1$. Therefore the minimum of $s$ is positive giving

$$
s \;\gtrsim\; 1 \;\sim\; r.
$$

Together with equivalences (39), (40) from the beginning this finishes off the technical proof of assertion (26).

**$T$ has Neumann faces.** In this case $\omega_T$ consists of less than five tetrahedra, and $\dim V_T < 5$. Although the representation of $v$ changes as well, the main ideas from above can still be applied to show the assertion. Thus we omit the proof.

It is noteworthy that the case of two or more Neumann boundary faces of $T$ gives rise to a particular phenomenon. If one can guarantee $\delta_i \leq \delta_* < 1 \,\forall i$ (with some parameter $\delta_*$ which is the same for all elements) then the resulting inequality is as before. Otherwise the inequality constant in (26) may depend on the shape of the Neumann boundary but does not depend on the triangulation $\mathcal{T}_h$ nor on $T$, *cf.* also [13, Lem. 3.5].

**Proof of (27).** Assume first that $E_i$ is an interior face, and consider the corresponding outer tetrahedron $T_i$. Apply (14) with $\varphi_E \equiv 1$ to obtain

$$
\|\nabla b_{E_i,\delta_i}\|_{T_i} \;\sim\; \delta_i^{1/2}\cdot h_{E_i,T_i}^{1/2}\cdot\min\{\delta_i\,h_{E_i,T_i}\,,\,h_{\min,T_i}\}^{-1}\cdot|E_i|^{1/2}.
$$

Together with $\|b_{E_i,\delta_i}\|_{E_i} = \|b_{E_i}\|_{E_i} \sim |E_i|^{1/2}$ this yields immediately

$$\|b_{E_i,\delta_i}\|_{E_i} \;\sim\; \delta_i^{-1/2} \cdot h_{E_i,T_i}^{-1/2} \cdot \min\{\delta_i\, h_{E_i,T_i}\,,\, h_{\min,T_i}\} \cdot \|\nabla b_{E_i,\delta_i}\|_{T_i}.$$

From $v\big|_{T_i} = \beta_i \cdot b_{E_i,\delta_i}\big|_{T_i}$ and $h_{E_i,T_i} \sim h_{E_i}$, $h_{\min,T_i} \sim h_{\min,T}$ one concludes

$$\|v\|_{E_i} \sim \delta_i^{-1/2} \cdot h_{E_i,T_i}^{-1/2} \cdot \min\{\delta_i\, h_{E_i,T_i}\,,\, h_{\min,T_i}\} \cdot \|\nabla v\|_{T_i}$$
$$\lesssim \delta_i^{-1/2} \cdot h_{E_i}^{-1/2} \cdot \min\{\delta_i\, h_{E_i}\,,\, h_{\min,T}\} \cdot \|\nabla v\|_{\omega_T}$$

which proves the assertion.

If $E_i$ is a Dirichlet face then $v|_{E_i} \equiv 0$, and (27) holds trivially. Finally, if $E_i$ is a Neumann face then the proof becomes more technical since no outer tetrahedron $T_i$ exists. Then one has to utilize similar ideas as for proving (26). The details are omitted here.     □

## References

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

[2] M. Ainsworth and J. Oden, *A Posteriori Error Estimation in Finite Element Analysis*. John Wiley & Sons, New York (2000).

[3] L. Angermann, Balanced *a-posteriori* error estimates for finite volume type discretizations of convection-dominated elliptic problems. *Computing* **55** (1995) 305–323.

[4] T. Apel and G. Lube, Anisotropic mesh refinement in stabilized Galerkin methods. *Numer. Math.* **74** (1996) 261–282.

[5] T. Apel and S. Nicaise, The finite element method with anisotropic mesh grading for elliptic problems in domains with corners and edges. *Math. Methods Appl. Sci.* **21** (1998) 519–549.

[6] I. Babuška and W.C. Rheinboldt, Error estimates for adaptive finite element computations. *SIAM J. Numer. Anal.* **15** (1978) 736–754.

[7] N.S. Bakhvalov, Optimization of methods for the solution of boundary value problems in the presence of a boundary layer. *Zh. Vychisl. Mat. i Mat. Fiz.* **9** (1969) 841–859. In Russian.

[8] R.E. Bank and A. Weiser, Some *a posteriori* error estimators for elliptic partial differential equations. *Math. Comput.* **44** (1985) 283–301.

[9] M. Beckers, *Numerical Integration in High Dimensions*. Ph.D. Thesis, Katholieke Universiteit Leuven / Louvain, Belgium (1992).

[10] P.G. Ciarlet, *The Finite Element Method for Elliptic Problems*. North-Holland Publishing Company, Amsterdam, New York, Oxford (1978).

[11] M. Dobrowolski, S. Gräf and C. Pflaum, On *a posteriori* error estimators in the finite element method on anisotropic meshes. *ETNA, Electron. Trans. Numer. Anal.* **8** (1999) 36–45.

[12] P. Keast, Moderate-degree tetrahedral quadrature formulas. *Comput. Methods Appl. Mech. Engrg.* **55** (1986) 339–348.

[13] G. Kunert, *A Posteriori Error Estimation for Anisotropic Tetrahedral and Triangular Finite Element Meshes*. Logos Verlag, Berlin (1999). Also Ph.D. Thesis, TU Chemnitz, `http://archiv.tu-chemnitz.de/pub/1999/0012/index.html`

[14] G. Kunert, An *a posteriori* residual error estimator for the finite element method on anisotropic tetrahedral meshes. *Numer. Math.* **86** (2000) 471–490. DOI 10.1007/s002110000170.

[15] G. Kunert, Towards anisotropic mesh construction and error estimation in the finite element method. To appear in *Numer. Meth. Partial Differential Equations*. Preprint SFB393/00_01, TU Chemnitz (2000). Also `http://archiv.tu-chemnitz.de/pub/2000/0066/index.html`

[16] G. Kunert, A local problem error estimator for anisotropic tetrahedral finite element meshes. *SIAM J. Numer. Anal.* **39** (2001) 668–689.

[17] G. Kunert, *A note on the energy norm for a singularly perturbed model problem*. Preprint SFB393/01-02, TU Chemnitz (2001). Also `http://archiv.tu-chemnitz.de/pub/2001/0006/index.html`

[18] G. Kunert, Robust *a posteriori* error estimation for a singularly perturbed reaction-diffusion equation on anisotropic tetrahedral meshes. To appear in *Adv. Comp. Math.*

[19] G. Kunert and R. Verfürth, Edge residuals dominate *a posteriori* error estimates for linear finite element methods on anisotropic triangular and tetrahedral meshes. *Numer. Math.* **86** (2000) 283–303. DOI 10.1007/s002110000152.

[20] J. Peraire, M. Vahdati, K. Morgan and O.C. Zienkiewicz, Adaptive remeshing for compressible flow computation. *J. Comput. Phys.* **72** (1987) 449–466.

[21] W. Rick, H. Greza and W. Koschel, FCT-solution on adapted unstructured meshes for compressible high speed flow computations. in *Flow Simulation with High-Performance Computers* **I**, in Notes Numer. Fluid Mech. **38**, E.H. Hirschel, Ed., Vieweg (1993) 334–438 .

[22] H.-G. Roos, M. Stynes and L. Tobiska, *Numerical Methods for Singularly Perturbed Differential Equations. Convection-Diffusion and Flow Problems.* Springer, Berlin (1996).

[23] K.G. Siebert, An *a posteriori* error estimator for anisotropic refinement. *Numer. Math.* **73** (1996) 373–398.

[24] R. Verfürth, *A posteriori* error estimation and adaptive mesh-refinement techniques. *J. Comput. Appl. Math.* **50** (1994) 67–83.

[25] R. Verfürth, *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques.* Wiley-Teubner, Chichester, Stuttgart (1996).

[26] R. Verfürth, Robust *a posteriori* error estimators for singularly perturbed reaction–diffusion equations. *Numer. Math.* **78** (1998) 479–493.

[27] R. Vilsmeier and D. Hänel, Computational aspects of flow simulation in three dimensional, unstructured, adaptive grids, in *Flow Simulation with High-Performance Computers* **II**, in Notes Numer. Fluid Mech. **52**, E.H. Hirschel, Ed., Vieweg (1996) 431–44.

[28] O.C. Zienkiewicz and J. Wu, Automatic directional refinement in adaptive analysis of compressible flows. *Internat. J. Numer. Methods Engrg.* **37** (1994) 2189–2210 .