



**HAL**  
open science

# A framework for robust a posteriori error control in unsteady nonlinear advection-diffusion problems

Vít Dolejší, Alexandre Ern, Martin Vohralík

► **To cite this version:**

Vít Dolejší, Alexandre Ern, Martin Vohralík. A framework for robust a posteriori error control in unsteady nonlinear advection-diffusion problems. 2011. hal-00652979v1

**HAL Id: hal-00652979**

**<https://hal.science/hal-00652979v1>**

Submitted on 16 Dec 2011 (v1), last revised 5 Sep 2012 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A framework for robust a posteriori error control in unsteady nonlinear advection-diffusion problems\*

Vít Dolejší<sup>†</sup>

Alexandre Ern<sup>‡</sup>

Martin Vohralík<sup>§</sup>

December 16, 2011

## Abstract

We derive a framework for a posteriori error estimates in unsteady, nonlinear, possibly degenerate, advection-diffusion problems. Our estimators are based on space-time equilibrated flux reconstruction and are locally computable. They are derived for the error measured in a space-time mesh-dependent dual norm stemming from the problem and meshes at hand augmented by a jump-based contribution measuring possible nonconformities in space. Owing to this choice, a guaranteed upper bound, as well as global efficiency and robustness with respect to all model (e.g., nonlinearities, advection dominance, final time) and discretization parameters are achieved. Local-in-time and in-space efficiency is also shown for a localized upper bound of the error measure. In order to apply the framework to a given numerical method, two simple conditions, local space-time mass conservation and an approximation property of the reconstructed fluxes, need to be verified. We show how to do this for the discontinuous Galerkin method in space and the Crank–Nicolson scheme in time. Numerical experiments illustrate the theory.

**Key words:** unsteady nonlinear advection-diffusion problem, a posteriori estimate, dual norm, flux reconstruction, flux equilibration, unified framework, robustness, discontinuous Galerkin method

## 1 Introduction

We consider the unsteady nonlinear problem

$$\partial_t u - \nabla \cdot \boldsymbol{\sigma}(u, \nabla u) = f \quad \text{in } Q := \Omega \times (0, t_F), \quad (1.1a)$$

$$u = 0 \quad \text{on } \partial\Omega \times (0, t_F), \quad (1.1b)$$

$$u(\cdot, 0) = u_0 \quad \text{in } \Omega, \quad (1.1c)$$

with  $\Omega \subset \mathbb{R}^d$ ,  $d \geq 2$ , a polygonal (polyhedral) domain,  $t_F > 0$  the final time,  $f$  the source term, and  $u_0$  the initial datum. The function  $\boldsymbol{\sigma}(u, \nabla u)$  takes the form

$$\boldsymbol{\sigma}(u, \nabla u) := \underline{\mathbf{K}}(u) \nabla u - \boldsymbol{\phi}(u), \quad (1.2)$$

where  $\underline{\mathbf{K}}(\cdot)$  is a nonlinear, possibly degenerate, tensor-valued function associated with diffusive transport and  $\boldsymbol{\phi}(\cdot)$  a nonlinear vector-valued function associated with advective transport. Precise assumptions on  $\underline{\mathbf{K}}(\cdot)$  and  $\boldsymbol{\phi}(\cdot)$  are specified in Section 2.1. In what follows,  $u$  is termed *potential* and  $-\boldsymbol{\sigma}(u, \nabla u)$  advective-diffusive *flux*. We assume that the problem (1.1a)–(1.1c) admits a unique weak solution  $u$  and that a fully

---

\*This work was partly supported by the Groupement MoMaS (PACEN/CNRS, ANDRA, BRGM, CEA, EdF, IRSN) and by the ERT project “Enhanced oil recovery and geological sequestration of CO<sub>2</sub>: mesh adaptivity, a posteriori error control, and other advanced techniques” (LJLL/IPPEN). The research of V. Dolejší is supported by the Grant No. 201/08/0012 of the Czech Science Foundation.

<sup>†</sup>Department of Numerical Mathematics, Charles University in Prague, Sokolovská 83, 186 75 Praha 8, Czech Republic (dolejsi@karlin.mff.cuni.cz).

<sup>‡</sup>Université Paris-Est, CERMICS, Ecole des Ponts, 77455 Marne-la-Vallée, France (ern@cermics.enpc.fr).

<sup>§</sup>UPMC Univ. Paris 06, UMR 7598, Laboratoire Jacques-Louis Lions, 75005, Paris, France & CNRS, UMR 7598, Laboratoire Jacques-Louis Lions, 75005, Paris, France (vohralik@ann.jussieu.fr).

discrete approximate potential, say  $u_{h\tau}$ , is available. The space discretization scheme can be nonconforming, so that  $u_{h\tau}$  can exhibit jumps across the interior spatial mesh faces.

Quite a broad literature has been devoted to the a posteriori error analysis for problems of the form (1.1a)–(1.1c). One of the main issues is to simultaneously prove global upper and (local) lower bounds for the error (reliability and (local) efficiency). For the linear heat equation, reliability is achieved in the energy norm by Picasso [32] and Repin [35], and in higher-order norms by Makridakis and Nochetto [26]. In Verfürth [39] and Bergam, Bernardi, and Mghazli [6], both reliability and efficiency are proven by augmenting the energy norm by a dual norm of the time derivative. The lower bound is then local in time but global in space. A further analysis unifying various space discretization schemes has been recently given in [20]. For nonlinear parabolic problems, reliability and local and global efficiency have been derived by Verfürth [37, 38] under a restriction on the relative size of space and time steps. This restriction has been lifted by Verfürth in [40] for convex two-dimensional spatial domains thereby achieving global efficiency; the price to pay is a solution of a linear diffusion problem by the finite element method on each time step. Robustness with respect to advection dominance or nonlinearities is not addressed in [40]. For unsteady linear advection-diffusion problems, Verfürth [41] has proved robustness with respect to advection dominance while augmenting the energy norm by a dual norm of the material derivative. A solution of a finite element reaction-diffusion problem on each time step is, again, necessary. All the above works concern the nondegenerate case; still less results are available on a posteriori error estimation for degenerate parabolic problems. We cite in particular Nochetto *et al.* [28] and Ohlberger [29]. In these works, only the error upper bound is derived.

The purpose of the present work is to derive *guaranteed*, *(locally) efficient*, and *robust* a posteriori error estimates for the problem (1.1a)–(1.1c). Here, *guaranteed* means that the estimates represent an explicitly computable upper bound on the error and *robustness* means that the efficiency is independent of all parameters (that is, spatial domain, final time, space-time discretization parameters, size of nonlinearity, and size of advection). Our key idea is the introduction of a *space-time mesh-dependent dual norm* to measure the error, see (2.6). This norm includes the full space-time nonlinear advection-diffusion operator, which stands in contrast to previous work where only parts of the differential operator at hand are included (i.e., the time derivative or the material derivative as in [39, 41]). Such approaches have been used recently by Chaillou and Suri [9, 10] and in [15] in the context of steady nonlinear diffusion problems. Moreover, evaluating the dual norm with respect to a specific mesh-dependent norm for test functions with bounded time and space derivatives in  $L^2(Q)$ , see (2.5a)–(2.5b), also allows to *simplify* substantially the *proof* of the *error lower bound*, whereby space-time bubble functions can be considered instead of the more usual space bubble functions at fixed times. This point also presents the practically crucial advantage that an error lower bound can be achieved using *locally computable estimators*, in contrast to [40, 41] where the solution of a global diffusion/reaction-diffusion problem on each time step is necessary to evaluate the estimators. The error measure (2.6) may seem rather weak at a first glance. However, this measure admits an easily and locally computable upper bound which consists of weighted  $L^2$ -norms of the potential error and of the error in the nonlinear advection-diffusion flux, see (2.9)–(2.10). Our efficiency results carry over to this norm and, moreover, can then be *localized in space and in time*. The numerical experiments of Section 8 actually show that our estimators generally provide efficient estimates for this weighted, physical  $L^2$ -norm.

Our results are derived in a *unified framework* where the actual numerical scheme used to obtain  $u_{h\tau}$  need not be specified. The error upper bound hinges on an advection-diffusion *flux reconstruction* and its local *space-time equilibration*, see Assumption 3.1, while the error lower bound requires a local *approximation property* on this flux, see Assumption 4.1. Applying the present framework to a given numerical scheme simply boils down to verifying these two assumptions. The reconstruction of the flux depends on the space discretization scheme, as discussed in [20] for the linear heat equation. Our approach is thus in line with the developments relying for linear model problems on the Prager and Synge equality [33] and pursued later by, e.g., Ladevèze [24], Neittaanmäki and Repin [27], Luce and Wohlmuth [25], Braess *et al.* [7], or Ainsworth [1].

This paper is structured as follows. We present the continuous and discrete settings and define our error measure  $\mathcal{J}_u(u_{h\tau})$  in Section 2. This is a sum of the above-discussed space-time mesh-dependent dual norm  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$  of (2.6) and of a weighted jumps term  $\mathcal{J}_{u,\text{NC}}(u_{h\tau})$  which measures possible nonconformities in space, see (2.11). We state our a posteriori error estimate, Theorem 3.3, in Section 3 yielding the error upper bound  $\mathcal{J}_u(u_{h\tau}) \leq \eta_{\text{FR}} + \eta_{\text{NC}} + \eta_{\text{IC}}$ . Theorem 4.4 of Section 4 then establishes the error lower bound  $\eta_{\text{FR}} + \eta_{\text{NC}} \lesssim \mathcal{J}_u(u_{h\tau})$  under the assumption that the quadrature error caused by the nonlinearity of the flux  $\sigma$  is small enough. Here,  $\lesssim$  means up to a generic constant which is independent of all discretization

and model parameters. Moreover, Theorem 4.2 provides a space-time localized version in terms of the computable error upper bound (2.9)–(2.10) on the error measure  $\mathcal{J}_{u, \text{FR}}(u_{h\tau})$ . We devote Sections 5 and 6 to the proofs of the results of Sections 3 and 4, respectively. To illustrate the developed abstract framework, we apply it in Section 7 to the discontinuous Galerkin method in space and the Crank–Nicolson scheme in time. Numerical experiments, including nonlinear degenerate advection-diffusion problems, are presented in conclusion in Section 8.

## 2 The setting

This section briefly describes the continuous and discrete settings and discusses the error measure.

### 2.1 Continuous setting and weak solution

We consider the function spaces

$$X := L^2(0, t_{\text{F}}; H_0^1(\Omega)), \quad (2.1a)$$

$$Y := \{\varphi \in L^2(0, t_{\text{F}}; H_0^1(\Omega)); \partial_t \varphi \in L^2(Q); \varphi(\cdot, t_{\text{F}}) = 0\}, \quad (2.1b)$$

recalling that  $\varphi \in L^2(0, t_{\text{F}}; H_0^1(\Omega))$  and  $\partial_t \varphi \in L^2(0, t_{\text{F}}; H^{-1}(\Omega))$  classically implies  $\varphi \in C^0([0, t_{\text{F}}]; L^2(\Omega))$ . The weak solution  $u$  of (1.1a)–(1.1c) is sought in the space  $X$ , whereas the space  $Y$  is used as test space. Specifically, we assume that there exists a unique weak solution  $u \in X$  such that

$$\int_0^{t_{\text{F}}} \{(f, \varphi) + (u, \partial_t \varphi) - (\boldsymbol{\sigma}(u, \nabla u), \nabla \varphi)\}(t) dt + (u_0, \varphi(\cdot, 0)) = 0 \quad \forall \varphi \in Y, \quad (2.2)$$

where  $(\cdot, \cdot)$  denotes the inner product in  $L^2(\Omega)$  or  $[L^2(\Omega)]^d$ . In order to ensure that all the terms in (2.2) are well-defined, we assume  $f \in L^2(Q)$ ,  $u_0 \in L^2(\Omega)$ , and  $\boldsymbol{\sigma}(u, \nabla u) \in [L^2(Q)]^d$  which is satisfied, for  $u \in X \cap L^\infty(Q)$ , e.g., if  $\mathbf{K} \in L^\infty_{\text{loc}}(\mathbb{R}; \mathbb{R}^{d \times d})$  and  $\phi \in C^1(\mathbb{R}; \mathbb{R}^d)$ . In deriving our a posteriori error estimators in Section 3, we exploit the fact that (2.2) consists, except for the contribution of the initial condition, of space-time inner products in  $L^2(Q)$ .

The present framework also covers some particular cases of degenerate parabolic equations of the form

$$\partial_t b(v) - \nabla \cdot (\tilde{\mathbf{K}} \nabla v - \tilde{\boldsymbol{\phi}}(v)) = f, \quad (2.3)$$

where  $b(\cdot)$  is an increasing function with locally Hölder regularity such that  $b(0) = 0$ . Problem (2.3) includes slow-diffusion-type problems, e.g., the porous media equation for which  $b(v) = v^{1/m}$ ,  $1 < m < +\infty$  (so that  $b'(0) = +\infty$ ), and fast-diffusion-type problems of elliptic-parabolic form, e.g., the Richards equation for which typically  $b'(0) = 0$ . Existence, uniqueness, and regularity results for problem (2.3) can be found in the work of Alt and Luckhaus [2] and Otto [30], leading to weak solutions  $v \in X \cap L^\infty(Q)$ . Since the function  $b$  is increasing, equation (2.3) can be recast into the form (1.1a) by setting  $u := b(v)$ , yielding  $\mathbf{K}(u) := \tilde{\mathbf{K}}(b^{-1})'(u)$  and  $\boldsymbol{\phi}(u) := \tilde{\boldsymbol{\phi}}(b^{-1}(u))$ . Therefore, the present analysis can be applied to (2.3) under the assumption  $b(v) \in X$ . This assumption holds true in the fast-diffusion regime, but not necessarily in the slow-diffusion regime where it is possible that  $v \in X$  but  $b(v) \notin X$ . The Stefan problem, governed by  $\partial_t u - \Delta \beta(u) = 0$  where  $\beta(\cdot)$  is a nondecreasing Lipschitz function, is not covered by the present assumptions either, since, in this case, the weak solution  $u$  can exhibit jumps. Finally, we observe that we do not include zero-order terms in (1.1a).

**Remark 2.1** (Physical units). *In advection-diffusion problems, the physical unit of the components of the diffusion tensor  $\mathbf{K}(\cdot)$  is  $L^2 \text{T}^{-1}$  and that of the components of the transport velocity  $\boldsymbol{\phi}(\cdot)$  is  $L \text{T}^{-1}$ . Here,  $L$  stands for length and  $\text{T}$  for time.*

### 2.2 Discrete setting and approximate solution

We consider an increasing sequence of discrete times  $\{t^n\}_{0 \leq n \leq N}$  such that  $t^0 = 0$  and  $t^N = t_{\text{F}}$ . We set  $I_n := (t^{n-1}, t^n]$  and  $\tau^n := t^n - t^{n-1}$  for all  $1 \leq n \leq N$ . We consider a time-sequence of matching simplicial meshes  $\{\mathcal{T}^n\}_{0 \leq n \leq N}$  of the spatial domain  $\Omega$ : the mesh  $\mathcal{T}^0$  is used to approximate the initial datum, and,

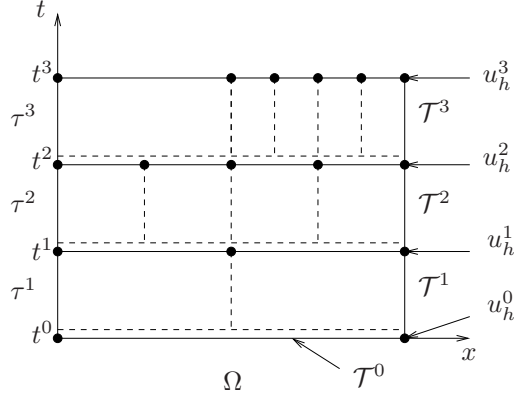


Figure 1: Time-dependent meshes and discrete solutions

for all  $1 \leq n \leq N$ , the mesh  $\mathcal{T}^n$  is used to march from  $t^{n-1}$  to  $t^n$ ; cf. Figure 1 for an illustration in one space dimension. We assume that  $\mathcal{T}^n$  is obtained from  $\mathcal{T}^{n-1}$  by refining some elements and coarsening some other ones. For all  $1 \leq n \leq N$ , we denote by  $\overline{\mathcal{T}}^{n-1,n}$  the coarsest common refinement of  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$  and by  $\underline{\mathcal{T}}^{n-1,n}$  the finest common coarsening of  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$ . Both space and time meshes can either be given a priori, or be generated by a space-time adaptive time-marching algorithm, see, e.g., [20] and the references therein.

At this stage, we do not need to specify any numerical scheme with which to obtain the approximate solution  $u_{h\tau}$ . We simply assume that  $u_{h\tau}$  is in the space

$$X_h := \{\varphi \in L^2(0, t_F; H^1(\mathcal{T})); \partial_t \varphi \in L^2(Q)\}.$$

Here,  $L^2(0, t_F; H^1(\mathcal{T}))$  is spanned by those functions  $\varphi \in L^2(Q)$  that, for all  $1 \leq n \leq N$  and all  $T \in \overline{\mathcal{T}}^{n-1,n}$ , satisfy  $\varphi|_{T \times I_n} \in L^2(I_n; H^1(T))$ . Functions in  $X_h$  can exhibit jumps across mesh interfaces, but, owing to the assumption  $\partial_t \varphi \in L^2(Q)$ , such functions are continuous with respect to time. This latter assumption is needed to express the residual in terms of space-time inner products in  $L^2(Q)$ .

Let  $1 \leq n \leq N$  and let  $T \in \mathcal{T}^n$  or  $T \in \overline{\mathcal{T}}^{n-1,n}$  or  $T \in \underline{\mathcal{T}}^{n-1,n}$ . We use the notation  $h_T$  for the diameter of  $T$  and  $\|\cdot\|_T$  for the norm on  $L^2(T)$ . Similarly, the norm on  $L^2(T \times I_n)$  is denoted by  $\|\cdot\|_{T \times I_n}$ . The corresponding inner products are denoted by  $(\cdot, \cdot)_T$  and  $(\cdot, \cdot)_{T \times I_n}$ , respectively. We collect in  $\mathcal{F}_T$  all the faces of  $T$  and in  $\mathcal{F}_T^{\text{int}}$  those that are subsets of  $\Omega$ . The set of all faces of the mesh  $\mathcal{T}^n$  is denoted by  $\mathcal{F}^n$ , a generic mesh face by  $F$ , and we use a similar notation for norms and inner products as above. For an interface  $F$ ,  $\mathbf{n}_F$  denotes its unit normal oriented in the sense the jump is evaluated, while, for a mesh element  $T$ ,  $\mathbf{n}_T$  denotes its outward unit normal. Finally,  $\mathcal{T}_T$  stands for all mesh elements sharing a face with the element  $T$ , whereas  $\mathcal{T}_F$  denotes the mesh elements sharing the face  $F$ .

The following scaled space-time Poincaré inequality is instrumental in deriving our error upper bound.

**Lemma 2.2** (Scaled space-time Poincaré inequality). *Let  $1 \leq n \leq N$  and let  $T \in \underline{\mathcal{T}}^{n-1,n}$ . Let  $\Pi_0$  denote the  $L^2(Q)$ -orthogonal projection onto constants in each space-time element  $T \times I_n$ . Set  $C_P := 1/\pi$ . Then, for all  $\varphi \in H^1(T \times I_n)$ ,*

$$\|\varphi - \Pi_0 \varphi\|_{T \times I_n} \leq C_P (h_T^2 \|\nabla \varphi\|_{T \times I_n}^2 + (\tau^n)^2 \|\partial_t \varphi\|_{T \times I_n}^2)^{\frac{1}{2}}. \quad (2.4)$$

*Proof.* The proof is straightforward; we present it for completeness. Let  $\varphi \in H^1(T \times I_n)$  and, for all  $t \in I_n$ , set  $\tilde{\varphi}(t) := |T|^{-1} \int_T \varphi(\mathbf{x}, t) \, d\mathbf{x}$ . Observing that  $(\varphi - \tilde{\varphi})$  and  $(\tilde{\varphi} - \Pi_0 \varphi)$  are  $L^2(T \times I_n)$ -orthogonal, we infer  $\|\varphi - \Pi_0 \varphi\|_{T \times I_n}^2 = \|\varphi - \tilde{\varphi}\|_{T \times I_n}^2 + \|\tilde{\varphi} - \Pi_0 \varphi\|_{T \times I_n}^2$ . For the first term on the right-hand side, the usual Poincaré inequality on  $T$  (which is convex) yields, cf. Payne and Weinberger [31] and Bebendorf [5], that

$$\begin{aligned} \|\varphi - \tilde{\varphi}\|_{T \times I_n}^2 &= \int_{I_n} \left( \int_T |\varphi - \tilde{\varphi}|^2(\mathbf{x}, t) \, d\mathbf{x} \right) dt \\ &\leq \int_{I_n} \left( C_P^2 h_T^2 \int_T |\nabla \varphi|^2(\mathbf{x}, t) \, d\mathbf{x} \right) dt = C_P^2 h_T^2 \|\nabla \varphi\|_{T \times I_n}^2. \end{aligned}$$

For the second term, observing that  $\Pi_0\varphi = (\tau^n)^{-1} \int_{I_n} \tilde{\varphi}(t) dt$ , the one-dimensional Poincaré inequality on  $I_n$  and the Cauchy–Schwarz inequality yield

$$\begin{aligned} \|\tilde{\varphi} - \Pi_0\varphi\|_{T \times I_n}^2 &= |T| \int_{I_n} |\tilde{\varphi} - \Pi_0\varphi|^2(t) dt \leq |T| \pi^{-2} (\tau^n)^2 \int_{I_n} |d_t \tilde{\varphi}|^2(t) dt \\ &= |T|^{-1} \pi^{-2} (\tau^n)^2 \int_{I_n} \left( \int_T \partial_t \varphi(\mathbf{x}, t) d\mathbf{x} \right)^2 dt \\ &\leq \pi^{-2} (\tau^n)^2 \int_{I_n} \left( \int_T |\partial_t \varphi|^2(\mathbf{x}, t) d\mathbf{x} \right) dt = \pi^{-2} (\tau^n)^2 \|\partial_t \varphi\|_{T \times I_n}^2. \end{aligned}$$

Collecting the above bounds yields (2.4).  $\square$

Finally, to alleviate the notation while allowing for nonconforming functions, we denote by  $\nabla$  the broken space gradient. Other choices are possible.

## 2.3 Error measure

The error measure described in this section combines a space-time mesh-dependent dual norm plus a non-conformity term.

### 2.3.1 Space-time mesh-dependent dual norm

Recalling definition (2.1b) of the space  $Y$ , we equip it with the norm

$$\|\varphi\|_{Y, T \times I_n}^2 := C_{T,n} (h_T^2 \|\nabla \varphi\|_{T \times I_n}^2 + (\tau^n)^2 \|\partial_t \varphi\|_{T \times I_n}^2), \quad \forall 1 \leq n \leq N, T \in \mathcal{T}^{n-1,n}, \quad (2.5a)$$

$$\|\varphi\|_Y^2 := \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} \|\varphi\|_{Y, T \times I_n}^2. \quad (2.5b)$$

$\|\cdot\|_Y$  is indeed a norm since functions in  $Y$  vanish on  $\Omega \times \{t_F\}$  and on  $\partial\Omega \times (0, t_F)$ . We observe that this norm depends on the space-time meshes. The positive quantities  $C_{T,n}$  are user-dependent weights which are further discussed in Section 2.3.4. Our robustness results are *not influenced* by the value assigned to these weights, see Remark 3.5. Later on, the same weight  $C_{T,n}$  and notation will be used for any subelement of  $T$ .

The first building block of our error measure is the quantity

$$\mathcal{J}_{u, \text{FR}}(u_{h\tau}) := \sup_{\varphi \in Y, \|\varphi\|_Y=1} \int_0^{t_F} \{(u_{h\tau} - u, \partial_t \varphi) + (\boldsymbol{\sigma}(u, \nabla u) - \boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}), \nabla \varphi)\}(t) dt. \quad (2.6)$$

Let  $v \in L^2(0, t_F; H^1(\mathcal{T}))$ . Define the residual  $R(v) \in Y'$  such that, for all  $\varphi \in Y$ ,

$$\langle R(v), \varphi \rangle_{Y', Y} := \int_0^{t_F} \{(f, \varphi) + (v, \partial_t \varphi) - (\boldsymbol{\sigma}(v, \nabla v), \nabla \varphi)\}(t) dt + (u_0, \varphi(\cdot, 0)). \quad (2.7)$$

Then we infer, owing to (2.2), that

$$\mathcal{J}_{u, \text{FR}}(u_{h\tau}) = \sup_{\varphi \in Y, \|\varphi\|_Y=1} \langle R(u_{h\tau}), \varphi \rangle_{Y', Y}, \quad (2.8)$$

showing that  $\mathcal{J}_{u, \text{FR}}(u_{h\tau})$  is a *dual norm of the residual of the approximate solution*.

The error measure  $\mathcal{J}_{u, \text{FR}}(u_{h\tau})$  cannot be computed easily in practice (in the test cases where the exact solution  $u$  is available). Indeed, its evaluation requires solving the following (infinite-dimensional space-time) problem: Find  $\psi \in Y$  such that  $(\psi, \varphi)_Y = \langle R(u_{h\tau}), \varphi \rangle_{Y', Y}$  for all  $\varphi \in Y$ , where  $(\cdot, \cdot)_Y$  denotes the inner product corresponding to the  $\|\cdot\|_Y$ -norm. Then, it is immediate that  $\mathcal{J}_{u, \text{FR}}(u_{h\tau}) = \|\psi\|_Y$ . However,

a computable upper bound on  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$  can be readily derived using the Cauchy–Schwarz inequality, leading to

$$\mathcal{J}_{u,\text{FR}}(u_{h\tau}) \leq \sup_{\varphi \in Y, \|\varphi\|_Y=1} \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} e_{\text{FR},T}^n \|\varphi\|_{Y,T \times I_n} \leq \left\{ \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} (e_{\text{FR},T}^n)^2 \right\}^{\frac{1}{2}} \quad (2.9)$$

that we denote by  $e_{\text{FR}}$ , with for all  $1 \leq n \leq N$  and all  $T \in \mathcal{T}^{n-1,n}$ ,

$$e_{\text{FR},T}^n := C_{T,n}^{-\frac{1}{2}} \{ (\tau^n)^{-2} \|u_{h\tau} - u\|_{T \times I_n}^2 + h_T^{-2} \|\sigma(u, \nabla u) - \sigma(u_{h\tau}, \nabla u_{h\tau})\|_{T \times I_n}^2 \}^{\frac{1}{2}}. \quad (2.10)$$

**Remark 2.3** (Dual norm of the residual and energy-type norms). *Under appropriate assumptions, a functional framework can be introduced for the nonlinear differential operator in (1.1a), and the error between  $u$  and  $u_{h\tau}$  can be measured in the corresponding energy-type norms. For conforming approximations ( $u_{h\tau} \in X$ ), following Verfürth (see, e.g., [37, Proposition 2.1]), the energy error can be bounded from above and from below by a dual norm of the residual. Such equivalence results hinge on suitable a priori bounds of the linearized differential operator, where it is in particular difficult to trace the influence of the size of the nonlinearities or of the advection dominance. Here we measure the error directly by (2.6) and avoid energy-type norms and the question of their equivalence with a dual norm of the residual.*

### 2.3.2 Nonconformity

As we allow for nonconformities ( $X_h$  is not a subspace of  $X$ ), we need to introduce a second building block of our error measure,

$$\mathcal{J}_{u,\text{NC}}(u_{h\tau}) := \left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} \sum_{F \in \mathcal{F}_T} C_{T,n}^{-1} h_T^{-2} C_{\underline{\mathbf{K}},\phi,T,F,n} \|\llbracket u - u_{h\tau} \rrbracket\|_{F \times I_n}^2 \right\}^{\frac{1}{2}}, \quad (2.11)$$

where  $\llbracket \cdot \rrbracket$  means the jump across interfaces and the actual value at boundary faces. The measure (2.11) is inspired by the penalty terms used in discontinuous Galerkin methods, cf., e.g., the early work of Arnold [3] and more recent work [18, 17] on heterogeneous diffusion problems with advection dominance. The positive constants  $C_{\underline{\mathbf{K}},\phi,T,F,n}$  are weights (physical unit is  $\text{L}^3\text{T}^{-2}$ ). Contrary to the weights  $C_{T,n}$  of (2.5a), the weights  $C_{\underline{\mathbf{K}},\phi,T,F,n}$  cannot be chosen arbitrarily because they are used to formulate the flux approximation property in Assumption 4.1. We refer to (7.7) in Section 7 for an example in the discontinuous Galerkin method. A straightforward consequence of the fact that the weak solution  $u$  is in  $X$  is the following lemma:

**Lemma 2.4** (Jumps of the weak solution). *For all  $1 \leq n \leq N$ , all  $T \in \overline{\mathcal{T}}^{n-1,n}$ , and all  $F \in \mathcal{F}_T$ ,  $\llbracket u \rrbracket = 0$  in  $L^2(F \times I_n)$ .*

Lemma 2.4 readily yields

$$\mathcal{J}_{u,\text{NC}}(u_{h\tau}) = \left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} \sum_{F \in \mathcal{F}_T} C_{T,n}^{-1} h_T^{-2} C_{\underline{\mathbf{K}},\phi,T,F,n} \|\llbracket u_{h\tau} \rrbracket\|_{F \times I_n}^2 \right\}^{\frac{1}{2}}. \quad (2.12)$$

Thus, in contrast to  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$ ,  $\mathcal{J}_{u,\text{NC}}(u_{h\tau})$  is easily computable. Finally, note that  $\mathcal{J}_{u,\text{NC}}(u_{h\tau}) = 0$  if and only if  $u_{h\tau} \in X$ , that is, if and only if  $u_{h\tau}$  is  $X$ -conforming.

### 2.3.3 Error measure

Our error measure is the sum of (2.6) and (2.11), i.e.

$$\mathcal{J}_u(u_{h\tau}) := \mathcal{J}_{u,\text{FR}}(u_{h\tau}) + \mathcal{J}_{u,\text{NC}}(u_{h\tau}). \quad (2.13)$$

There holds  $\mathcal{J}_u(u_{h\tau}) = 0$  if and only if  $u_{h\tau} = u$ . Indeed, if  $u_{h\tau} = u$ ,  $\mathcal{J}_u(u_{h\tau})$  clearly equals 0. Conversely, if  $\mathcal{J}_u(u_{h\tau}) = 0$ , then  $\mathcal{J}_{u,\text{FR}}(u_{h\tau}) = \mathcal{J}_{u,\text{NC}}(u_{h\tau}) = 0$ . Thus  $u_{h\tau}$  is in  $X$  and hence, since the weak solution is uniquely characterized by the property  $\langle R(u), \varphi \rangle_{Y',Y} = 0$  for all  $\varphi \in Y$ , see (2.2) and (2.7), we infer  $u_{h\tau} = u$ .

### 2.3.4 Choice of the weights $C_{T,n}$

One relevant choice for the weights  $C_{T,n}$  is such that their physical unit is  $T^{-1}$  since, with this choice, the error measure has the same physical unit as the classical energy norm in the context of advection-diffusion problems. A possible example is

$$C_{T,n} := \left( \frac{t_F}{(\tau^n)^2} + \frac{C_{\phi,T,n}}{h_T^2} + \frac{C_{\mathbf{K},T,n}}{h_T^2} \right), \quad (2.14)$$

with the advection-dependent weight  $C_{\phi,T,n} := h_\Omega \|\phi'(u_{h\tau})\|_{\infty, T \times I_n}$  and the diffusion-dependent weight  $C_{\mathbf{K},T,n} := (h_\Omega/h_T) \|\mathbf{K}(u_{h\tau})\|_{\infty, T \times I_n}$  (here,  $h_\Omega$  denotes the diameter of the domain  $\Omega$ ). For advection-dominated unsteady problems, the first two addends on the right-hand side of (2.14) can be expected to be of similar size if the Courant numbers  $\tau^n \|\phi'(u_{h\tau})\|_{\infty, T \times I_n} / h_T$  are of order unity, as well as the ratios  $t_F \|\phi'(u_{h\tau})\|_{\infty, T \times I_n} / h_\Omega$  (meaning that the final time allows particles advected by the flow to cross a relevant part of the domain). Moreover, the ratio of the second to the third addend is of the order of the local Péclet numbers  $h_T \|\phi'(u_{h\tau})\|_{\infty, T \times I_n} / \|\mathbf{K}(u_{h\tau})\|_{\infty, T \times I_n}$ , so that the third addend is negligible for dominant advection.

**Remark 2.5** (Steady case). *For the steady variant of (1.1a)–(1.1c), a similar reasoning can actually be applied to motivate that robustness with respect to the Péclet number, say  $Pe$ , can be expected for the present setting. Indeed, dropping the index  $n$  and choosing  $C_T := \left( \frac{C_{\phi,T}}{h_T^2} + \frac{C_{\mathbf{K},T}}{h_T^2} \right)$ , with  $C_{\phi,T} := h_T \|\phi'(u_{h\tau})\|_{\infty, T}$ ,  $C_{\mathbf{K},T} := \|\mathbf{K}(u_{h\tau})\|_{\infty, T}$ , and local Péclet numbers  $Pe_T := C_{\phi,T} / C_{\mathbf{K},T}$ , the upper bound (2.9) on the error measure is multiplied by the factor  $Pe^{\frac{1}{2}}$  with respect to the usual energy norm, while the error indicators (cf. (3.2)–(3.5) below) are multiplied by the factor  $Pe^{-\frac{1}{2}}$  with respect to the classical energy indicators (without cutoff factors), see, e.g., Verfürth [41] and [17]. In this point, a similarity to the approach of Sangalli [36] can be observed.*

## 3 A posteriori error estimate

This section collects the main results of this paper concerning the error upper bound. The approximation error is measured by (2.13), and the estimators are defined using an equilibrated flux reconstruction.

### 3.1 Equilibrated flux reconstruction

In order to proceed as generally as possible, in particular without the definition of any numerical scheme for approximating (1.1a)–(1.1c), we make the following assumption (recall that  $\mathbf{H}(\text{div}, \Omega)$  is spanned by vector fields in  $[L^2(\Omega)]^d$  which admit a weak divergence in  $L^2(\Omega)$ , cf. [8]):

**Assumption 3.1** (Space-time equilibrated flux reconstruction). *There exists a flux reconstruction  $\mathbf{t}_{h\tau}$  such that  $\mathbf{t}_{h\tau} \in \mathbf{L}^2(0, t_F; \mathbf{H}(\text{div}, \Omega))$  and such that this flux is equilibrated in the sense that*

$$(f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, 1)_{T \times I_n} = 0 \quad \forall 1 \leq n \leq N, T \in \mathcal{T}^{n-1, n}. \quad (3.1)$$

Specific constructions of the flux  $\mathbf{t}_{h\tau}$  depend on the spatial discretization at hand; various examples are discussed in [20]. We present an example of such construction in the context of discontinuous Galerkin methods for problem (1.1a)–(1.1c) in Section 7.

**Remark 3.2** (Local mass conservation). *Equation (3.1) expresses local mass conservation over the space-time element  $T \times I_n$ . A similar assumption has been made in [20], see equation (3.4) therein. In [20], however, the local mass conservation had to be satisfied over a given mesh element  $T$  for all times  $t \in I_n$ . The present assumption, being more general, allows for more flexibility. In particular, it allows us to use the scaled Poincaré inequality (2.4) in the proof of the error upper bound. This local mass conservation is supposed in (3.1) on the elements of the common coarsening of  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$  only, which in particular fits to the Crank-Nicolson time stepping of Section 7. Local mass conservation on the elements of  $\mathcal{T}^n$  may be used in other cases.*



## 3.2 Guaranteed a posteriori error estimate

We are now in a position to state our main result concerning the error upper bound. For all  $1 \leq n \leq N$  and all  $T \in \mathcal{T}^{n-1,n}$ , we define the *residual*, *flux*, and *nonconformity estimators* respectively

$$\eta_{\mathbf{R},T}^n := C_{T,n}^{-\frac{1}{2}} C_{\mathbf{P}} \|f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}\|_{T \times I_n}, \quad (3.2)$$

$$\eta_{\mathbf{F},T}^n := C_{T,n}^{-\frac{1}{2}} h_T^{-1} \|\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}\|_{T \times I_n}, \quad (3.3)$$

$$\eta_{\mathbf{NC},T}^n := \left\{ \sum_{T' \in \overline{\mathcal{T}}^{n-1,n}, T' \subset T} \sum_{F \in \mathcal{F}_{T'}} C_{T',n}^{-1} h_{T'}^{-2} C_{\mathbf{K},\phi,T',F,n} \|\llbracket u_{h\tau} \rrbracket\|_{F \times I_n}^2 \right\}^{\frac{1}{2}}, \quad (3.4)$$

and set  $\eta_{\mathbf{FR},T}^n := \eta_{\mathbf{F},T}^n + \eta_{\mathbf{R},T}^n$ . These estimators are local-in-space and in-time. We define their global space-time versions as  $\eta_{\bullet} := \left\{ \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} (\eta_{\bullet,T}^n)^2 \right\}^{\frac{1}{2}}$  for  $\bullet \in \{\mathbf{F}, \mathbf{R}, \mathbf{FR}, \mathbf{NC}\}$ . Finally, we define the *initial condition estimator* as

$$\eta_{\mathbf{IC}} := \left\{ \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} (\eta_{\mathbf{IC},T}^n)^2 \right\}^{\frac{1}{2}}, \quad \eta_{\mathbf{IC},T}^n := C_{T,n}^{-\frac{1}{2}} (\tau^n)^{-\frac{1}{2}} \|u_0 - u_{h\tau}(\cdot, 0)\|_T. \quad (3.5)$$

**Theorem 3.3** (Guaranteed a posteriori error estimate). *Let  $u \in X$  be the weak solution given by (2.2) and let  $u_{h\tau} \in X_h$  be arbitrary. Let Assumption 3.1 hold true. Let  $\eta_{\mathbf{R},T}^n$ ,  $\eta_{\mathbf{F},T}^n$ ,  $\eta_{\mathbf{NC},T}^n$ , and  $\eta_{\mathbf{IC}}$  be defined by (3.2)–(3.5). Then,*

$$\mathcal{J}_u(u_{h\tau}) \leq \eta_{\mathbf{FR}} + \eta_{\mathbf{NC}} + \eta_{\mathbf{IC}}. \quad (3.6)$$

The proof is given in Section 5. We now present several remarks.

**Remark 3.4** (Interpretation of the error estimators). *The flux estimator  $\eta_{\mathbf{F},T}^n$  is related to the possible violation of the constitutive relation (1.2) at the discrete level, the residual estimator  $\eta_{\mathbf{R},T}^n$  to the possible violation of the equilibrium condition (1.1a) at the discrete level, while  $\eta_{\mathbf{NC},T}^n$  and  $\eta_{\mathbf{IC},T}^n$  are related to the possible violation of the constraints ( $u \in X$  and  $u(\cdot, 0) = u_0$ ) at the discrete level.*

**Remark 3.5** (Weights  $C_{T,n}$ ). *We stress that our robustness results are not influenced by the value assigned to the weights  $C_{T,n}$  of (2.5a). Indeed, multiplying them by a positive factor  $\lambda$  scales both error measures  $\mathcal{J}_{u,\mathbf{FR}}(u_{h\tau})$  and  $\mathcal{J}_{u,\mathbf{NC}}(u_{h\tau})$  defined by (2.6) and (2.11) by the factor  $\lambda^{-\frac{1}{2}}$ , while the error estimators  $\eta_{\mathbf{R},T}^n$ ,  $\eta_{\mathbf{F},T}^n$ ,  $\eta_{\mathbf{NC},T}^n$ , and  $\eta_{\mathbf{IC},T}^n$  defined by (3.2)–(3.5) are scaled by the same factor. Hence, the ratio of error measure  $\mathcal{J}_u(u_{h\tau})$  to error indicators is independent of the scaling factor  $\lambda$ .*

**Remark 3.6** (Initial condition estimator). *Whenever the weights  $C_{T,n}$  are chosen according to (2.14), there holds  $C_{T,n} \geq t_{\mathbf{F}}(\tau^n)^{-2}$ , whence we readily infer that  $\eta_{\mathbf{IC}} \leq \|u_{h\tau}(\cdot, 0) - u_0\|_{\Omega}$ , that is,  $\eta_{\mathbf{IC}}$  is upper-bounded by the usual  $L^2(\Omega)$ -norm of the approximation error on the initial condition.*

**Remark 3.7** (Comparison with [20]). *In [20], a posteriori (augmented) energy norm estimates for the heat equation are derived, which leads to considering a potential reconstruction for nonconforming discretization schemes along with the flux reconstruction. Using the nonconformity error measure (2.11), we circumvent here this potential reconstruction which can be intricate on time-varying meshes and leads to restrictions on the discretization parameters in efficiency proofs, see [20, Section 3.2.1].*

## 4 Efficiency and robustness

This section deals with the efficiency and robustness of our estimates. For simplicity, we assume that the source term  $f$  is a piecewise space-time polynomial (on  $\overline{\mathcal{T}}^{n-1,n}$ ); otherwise, a classical data oscillation term has to be included in the error lower bound. We likewise need to assume that both  $u_{h\tau}$  and  $\mathbf{t}_{h\tau}$  are piecewise space-time polynomials (on  $\overline{\mathcal{T}}^{n-1,n}$ ). We observe that, because of the nonlinear functions  $\mathbf{K}(\cdot)$  and  $\phi(\cdot)$ , the

flux  $\sigma(u_{h\tau}, \nabla u_{h\tau})$  is not necessarily a piecewise polynomial, even if  $u_{h\tau}$  is. We take into account quadrature errors resulting from these nonlinearities by letting

$$\bar{\sigma}(u_{h\tau}, \nabla u_{h\tau})|_T := \mathcal{P}_T(\sigma(u_{h\tau}, \nabla u_{h\tau})|_T), \quad \forall 1 \leq n \leq N, T \in \mathcal{T}^{n-1,n}, \quad (4.1)$$

where  $\mathcal{P}_T$  is a projection-type operator mapping onto piecewise space-time polynomials, see (7.5)–(7.6) below for an example. We define the *quadrature error estimator*

$$\eta_{\text{qd},T}^n := C_{T,n}^{-\frac{1}{2}} h_T^{-1} \|\bar{\sigma}(u_{h\tau}, \nabla u_{h\tau}) - \sigma(u_{h\tau}, \nabla u_{h\tau})\|_{T \times I_n}, \quad (4.2)$$

together with its global space-time version  $\eta_{\text{qd}} := \left\{ \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1,n}} (\eta_{\text{qd},T}^n)^2 \right\}^{\frac{1}{2}}$ .

Henceforth,  $A \lesssim B$  means that there exists a constant  $C$  such that  $A \leq CB$ , with  $C$  only depending on space dimension, maximal polynomial degree, shape-regularity of the meshes  $\mathcal{T}^n$  for all  $1 \leq n \leq N$ , the maximal ratio  $h_T/h_{T'}$  over all  $1 \leq n \leq N$ ,  $T \in \mathcal{T}^{n-1,n}$  (the common coarsening of  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$ ) and  $T' \in \bar{\mathcal{T}}^{n-1,n}$  (the common refinement of  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$ ),  $T'$  a subelement of  $T$ , and the maximal ratios  $(C_{T,n}/C_{T',n})^{\frac{1}{2}}$ , for all  $T \in \mathcal{T}^{n-1,n}$  and  $T' \in \bar{\mathcal{T}}^{n-1,n}$  sharing a face and all  $1 \leq n \leq N$ . We thus in particular suppose that both refinement and coarsening are not too abrupt.

## 4.1 Approximation property of the flux reconstruction

Let, for  $1 \leq n \leq N$  and  $T \in \bar{\mathcal{T}}^{n-1,n}$ ,  $\eta_{\text{clas},T}^n$  be given by

$$\begin{aligned} & h_T \|f - \partial_t u_{h\tau} + \nabla \cdot (\bar{\sigma}(u_{h\tau}, \nabla u_{h\tau}))\|_{T \times I_n} \\ & + \left\{ \sum_{F \in \mathcal{F}_T^{\text{int}}} h_F \|[\![\bar{\sigma}(u_{h\tau}, \nabla u_{h\tau})]\!] \cdot \mathbf{n}_F\|_{F \times I_n}^2 \right\}^{\frac{1}{2}} + \left\{ \sum_{F \in \mathcal{F}_T} C_{\mathbf{k},\phi,T,F,n} \|[\![u_{h\tau}]\!]\|_{F \times I_n}^2 \right\}^{\frac{1}{2}}. \end{aligned} \quad (4.3)$$

The quantity  $\eta_{\text{clas},T}^n$  can be viewed as a classical residual-based a posteriori error estimator. In order to carry the analysis without specifying a numerical scheme, we make the following assumption on the reconstructed flux:

**Assumption 4.1** (Flux approximation property). *There holds*

$$\|\bar{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}\|_{T \times I_n}^2 \lesssim \sum_{T' \in \bar{\mathcal{T}}^{n-1,n}, T' \subset T} (\eta_{\text{clas},T'}^n)^2 \quad \forall 1 \leq n \leq N, T \in \mathcal{T}^{n-1,n}.$$

An example on how to achieve Assumption 4.1 in the context of discontinuous Galerkin methods in space and the Crank–Nicolson scheme in time is presented in Section 7.

## 4.2 Robust local-in-space and in-time efficiency

To control quadrature errors, we introduce coefficients  $0 < \gamma_{\text{qd},T}^n$  and require

$$\max_{T' \in \bar{\mathcal{T}}_T} \eta_{\text{qd},T'}^n \leq \gamma_{\text{qd},T}^n (\eta_{\text{FR},T}^n + \eta_{\text{NC},T}^n), \quad \forall 1 \leq n \leq N, T \in \mathcal{T}^{n-1,n}, \quad (4.4)$$

recalling that  $\bar{\mathcal{T}}_T$  collects the mesh elements sharing a face with  $T$ .

**Theorem 4.2** (Robust local-in-space and in-time efficiency). *Let a time step  $1 \leq n \leq N$  and a mesh element  $T \in \mathcal{T}^{n-1,n}$  be fixed. Let (4.4), with  $\gamma_{\text{qd},T}^n$  small enough, hold. Let Assumption 4.1 hold true. Recall the definition (2.10) of  $e_{\text{FR},T}^n$  and define  $\mathcal{J}_{u,\text{NC},T}(u_{h\tau})$  as  $\eta_{\text{NC},T}^n$  in (3.4), with  $[\![u_{h\tau}]\!]$  replaced by  $[\![u - u_{h\tau}]\!]$ . Then*

$$\eta_{\text{FR},T}^n + \eta_{\text{NC},T}^n \lesssim \left\{ \sum_{T' \in \bar{\mathcal{T}}_T} (e_{\text{FR},T'}^n)^2 \right\}^{\frac{1}{2}} + \mathcal{J}_{u,\text{NC},T}(u_{h\tau}). \quad (4.5)$$

**Remark 4.3** (Comment on Theorem 4.2). *Estimate (4.5) says that our estimators represent a robust local lower bound for the error measures  $e_{\text{FR},T}^n$ , whose Hilbertian sum provides an upper bound on the error measure  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$ , see (2.9), augmented by the jump seminorm.*

### 4.3 Robust global efficiency

As above, in order to control quadrature errors, we introduce a coefficient  $0 < \gamma_{\text{qd}}$  and require

$$\eta_{\text{qd}} \leq \gamma_{\text{qd}}(\eta_{\text{FR}} + \eta_{\text{NC}}). \quad (4.6)$$

Then we have the following full equivalence result between our estimators and  $\mathcal{J}_u(u_{h\tau})$ .

**Theorem 4.4** (Robust global efficiency). *Let (4.6), with  $\gamma_{\text{qd}}$  small enough, be satisfied. Let Assumption 4.1 hold true. Then*

$$\eta_{\text{FR}} + \eta_{\text{NC}} \lesssim \mathcal{J}_u(u_{h\tau}). \quad (4.7)$$

## 5 Proof of the error upper bound

The goal of this section is to prove Theorem 3.3. The proof is decomposed in two steps.

**Lemma 5.1** (Bound on  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$ ). *There holds  $\mathcal{J}_{u,\text{FR}}(u_{h\tau}) \leq \eta_{\text{FR}} + \eta_{\text{IC}}$ .*

*Proof.* Let  $\varphi \in Y$  with  $\|\varphi\|_Y = 1$  be given. Using (2.7), subtracting the quantity  $\int_0^{t_{\text{F}}} (\mathbf{t}_{h\tau}, \nabla \varphi)(t) dt + \int_0^{t_{\text{F}}} (\nabla \cdot \mathbf{t}_{h\tau}, \varphi)(t) dt$ , which is equal to zero owing to the Green theorem since  $\mathbf{t}_{h\tau} \in \mathbf{L}^2(0, t_{\text{F}}; \mathbf{H}(\text{div}, \Omega))$  by Assumption 3.1, integrating by parts in time since  $\partial_t u_{h\tau} \in L^2(Q)$  by assumption, and using that  $\varphi(\cdot, 0) = -\int_0^{t_{\text{F}}} \partial_t \varphi(t) dt$  since  $\varphi(\cdot, t_{\text{F}}) = 0$  by assumption, we infer

$$\begin{aligned} & \langle R(u_{h\tau}), \varphi \rangle_{Y', Y} \\ &= \int_0^{t_{\text{F}}} \{ (f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, \varphi) - (\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}, \nabla \varphi) \}(t) dt - (u_{h\tau}(\cdot, 0) - u_0, \varphi(\cdot, 0)) \\ &= \int_0^{t_{\text{F}}} \{ (f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, \varphi) + (u_{h\tau}(\cdot, 0) - u_0, \partial_t \varphi) - (\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}, \nabla \varphi) \}(t) dt. \end{aligned}$$

We now employ the second part of Assumption 3.1, namely (3.1). This enables us to rewrite equivalently

$$\begin{aligned} \langle R(u_{h\tau}), \varphi \rangle_{Y', Y} &= \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1, n}} \{ (f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, \varphi - \Pi_0 \varphi)_{T \times I_n} \\ &\quad + (u_{h\tau}(\cdot, 0) - u_0, \partial_t \varphi)_{T \times I_n} - (\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}, \nabla \varphi)_{T \times I_n} \}. \end{aligned}$$

The Cauchy–Schwarz inequality, the scaled space-time Poincaré inequality (2.4), and the definition (2.5a) lead to the bound

$$\begin{aligned} & \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1, n}} \{ \|f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}\|_{T \times I_n} C_{\text{P}} C_{T,n}^{-\frac{1}{2}} \|\varphi\|_{Y, T \times I_n} \\ & \quad + C_{T,n}^{-\frac{1}{2}} ((\tau^n)^{-2} \|u_0 - u_{h\tau}(\cdot, 0)\|_{T \times I_n}^2 + h_T^{-2} \|\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}\|_{T \times I_n}^2)^{\frac{1}{2}} \|\varphi\|_{Y, T \times I_n} \} \\ &= \sum_{n=1}^N \sum_{T \in \mathcal{T}^{n-1, n}} (\eta_{\text{R}, T}^n + ((\eta_{\text{F}, T}^n)^2 + (\eta_{\text{IC}, T}^n)^2)^{\frac{1}{2}}) \|\varphi\|_{Y, T \times I_n} \end{aligned}$$

on  $\langle R(u_{h\tau}), \varphi \rangle_{Y', Y}$ , whence the assertion follows from  $((\eta_{\text{F}, T}^n)^2 + (\eta_{\text{IC}, T}^n)^2)^{\frac{1}{2}} \leq \eta_{\text{F}, T}^n + \eta_{\text{IC}, T}^n$ , the Cauchy–Schwarz inequality, the definition (2.5b) of the  $\|\cdot\|_Y$ -norm, the fact that  $\|\varphi\|_Y = 1$ , and the relation (2.8).  $\square$

**Lemma 5.2** (Bound on  $\mathcal{J}_{u,\text{NC}}(u_{h\tau})$ ). *There holds  $\mathcal{J}_{u,\text{NC}}(u_{h\tau}) = \eta_{\text{NC}}$ .*

*Proof.* Immediate owing to (2.12) and (3.4).  $\square$

## 6 Proof of the error lower bound

We present here the proofs of Theorems 4.2 and 4.4. An important ingredient is the bubble function technique (see [39, 41]) that we extend here to space-time bubbles.

### 6.1 Proof of Theorem 4.2

*Proof.* Let a time step  $1 \leq n \leq N$  and a mesh element  $T \in \overline{\mathcal{T}}^{n-1,n}$  be fixed. The first two steps of the proof are devoted to bounding the quantity  $\eta_{\text{clas},T}^n$  defined by (4.3), while staying on the common refinement  $\overline{\mathcal{T}}^{n-1,n}$ .

*Step 1, bound on  $C_{T,n}^{-\frac{1}{2}} \|f - \partial_t u_{h\tau} + \nabla \cdot \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{T \times I_n}$ .* Let us prove that

$$C_{T,n}^{-\frac{1}{2}} \|f - \partial_t u_{h\tau} + \nabla \cdot \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{T \times I_n} \lesssim e_{\text{FR},T}^n + \eta_{\text{qd},T}^n, \quad (6.1)$$

with  $e_{\text{FR},T}^n$  defined by (2.10). Set  $v_{T,n} := (f - \partial_t u_{h\tau} + \nabla \cdot \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}))|_{T \times I_n}$ . Let  $\psi_{T,n}$  be the *space-time bubble function* on  $T \times I_n$  given by the product of the barycentric coordinates on  $T$  and of the barycentric coordinates on  $I_n$ . Note that both  $\psi_{T,n}$  and  $v_{T,n}$  are (space-time) polynomials since we are on refinement of both  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$ . Then, by norm equivalence in finite-dimensional spaces, there holds

$$(v_{T,n}, v_{T,n})_{T \times I_n} \lesssim (v_{T,n}, \psi_{T,n} v_{T,n})_{T \times I_n}. \quad (6.2)$$

Using an inverse inequality *separately* in space and in time, we obtain

$$h_T \|\nabla(\psi_{T,n} v_{T,n})\|_{T \times I_n} \lesssim \|\psi_{T,n} v_{T,n}\|_{T \times I_n}, \quad (6.3a)$$

$$\tau^n \|\partial_t(\psi_{T,n} v_{T,n})\|_{T \times I_n} \lesssim \|\psi_{T,n} v_{T,n}\|_{T \times I_n}. \quad (6.3b)$$

The norm  $\|\cdot\|_{Y,T \times I_n}$  defined by (2.5a) was precisely designed in order to use these inequalities at the present stage. Using (6.3), (2.5a), and  $\|\psi_{T,n}\|_{\infty,T,n} \leq 1$ , we infer

$$\begin{aligned} C_{T,n}^{-1} \|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n}^2 &= (h_T^2 \|\nabla(\psi_{T,n} v_{T,n})\|_{T \times I_n}^2 + (\tau^n)^2 \|\partial_t(\psi_{T,n} v_{T,n})\|_{T \times I_n}^2) \\ &\lesssim \|\psi_{T,n} v_{T,n}\|_{T \times I_n}^2 \leq \|v_{T,n}\|_{T \times I_n}^2. \end{aligned} \quad (6.4)$$

Thus, using the definition of  $v_{T,n}$ , (6.2), (2.2), the Green theorem, and (6.4) yields, with the notation  $\boldsymbol{\sigma}_d := \boldsymbol{\sigma}(u, \nabla u) - \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})$ ,

$$\begin{aligned} &C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 \\ &\lesssim C_{T,n}^{-1} (f - \partial_t u_{h\tau} + \nabla \cdot \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}), \psi_{T,n} v_{T,n})_{T \times I_n} \\ &= C_{T,n}^{-1} \frac{(u_{h\tau} - u, \partial_t(\psi_{T,n} v_{T,n}))_{T \times I_n} + (\boldsymbol{\sigma}_d, \nabla(\psi_{T,n} v_{T,n}))_{T \times I_n}}{\|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n}} \|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n} \\ &\lesssim \frac{(u_{h\tau} - u, \partial_t(\psi_{T,n} v_{T,n}))_{T \times I_n} + (\boldsymbol{\sigma}_d, \nabla(\psi_{T,n} v_{T,n}))_{T \times I_n}}{\|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n}} C_{T,n}^{-\frac{1}{2}} \|v_{T,n}\|_{T \times I_n}. \end{aligned} \quad (6.5)$$

Thus,

$$C_{T,n}^{-\frac{1}{2}} \|v_{T,n}\|_{T \times I_n} \lesssim \frac{(u_{h\tau} - u, \partial_t(\psi_{T,n} v_{T,n}))_{T \times I_n} + (\boldsymbol{\sigma}_d, \nabla(\psi_{T,n} v_{T,n}))_{T \times I_n}}{\|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n}},$$

whence (6.1) follows from the Cauchy–Schwarz inequality, the definitions (2.10) of  $e_{\text{FR},T}^n$  and (4.2) of  $\eta_{\text{qd},T}^n$ , and the triangle inequality.

*Step 2, bound on  $C_{T,n}^{-\frac{1}{2}} h_T^{-\frac{1}{2}} \|\overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{\mathbf{n}_F}$  for all  $F \in \mathcal{F}_T^{\text{int}}$ .* Let  $F \in \mathcal{F}_T^{\text{int}}$  be fixed. Recalling that  $\mathcal{T}_F$  denotes the simplices sharing the face  $F$ , let us prove that

$$C_{T,n}^{-\frac{1}{2}} h_T^{-\frac{1}{2}} \|\overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{\mathbf{n}_F} \lesssim \sum_{T' \in \mathcal{T}_F} (e_{\text{FR},T'}^n + \eta_{\text{qd},T'}^n). \quad (6.6)$$

Set  $v_{F,n} := \llbracket \bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}) \rrbracket \cdot \mathbf{n}_F|_{F \times I_n}$ . Let  $\psi_{F,n}$  be the space-time bubble function on  $F \times I_n$  given by the product of the barycentric coordinates with vertices in  $F$  and of the barycentric coordinates on  $I_n$ . Then, by norm equivalence in finite-dimensional spaces, there holds

$$(v_{F,n}, v_{F,n})_{F \times I_n} \lesssim (v_{F,n}, \psi_{F,n} v_{F,n})_{F \times I_n}. \quad (6.7)$$

Using the same notation for the extension of the function  $v_{F,n}$  onto  $\mathcal{T}_F$  by constant values in the direction of the normal  $\mathbf{n}_F$  of  $F$ , we also infer the estimate

$$\|v_{F,n}\|_{\mathcal{T}_F \times I_n} \lesssim h_F^{\frac{1}{2}} \|v_{F,n}\|_{F \times I_n}. \quad (6.8)$$

Using (6.7), the Green theorem, (2.2), the Cauchy–Schwarz inequality, (6.8), the fact that  $\|\psi_{F,n}\|_{\infty, \mathcal{T}_F \times I_n} \leq 1$ , and (6.4) where we replace  $\psi_{T,n} v_{T,n}$  by  $\psi_{F,n} v_{F,n}$  leads to

$$\begin{aligned} & C_{T,n}^{-1} h_T^{-1} \|v_{F,n}\|_{F \times I_n}^2 \\ & \lesssim C_{T,n}^{-1} h_T^{-1} (\llbracket \bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}) \rrbracket \cdot \mathbf{n}_F, \psi_{F,n} v_{F,n})_{F \times I_n} \\ & = C_{T,n}^{-1} h_T^{-1} (f - \partial_t u_{h\tau} + \nabla \cdot \bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}), \psi_{F,n} v_{F,n})_{\mathcal{T}_F \times I_n} \\ & \quad + C_{T,n}^{-1} h_T^{-1} \frac{(u - u_{h\tau}, \partial_t(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n} - (\boldsymbol{\sigma}_d, \nabla(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n}}{\|\psi_{F,n} v_{F,n}\|_{Y, \mathcal{T}_F, n}} \|\psi_{F,n} v_{F,n}\|_{Y, \mathcal{T}_F, n} \\ & \lesssim C_{T,n}^{-\frac{1}{2}} \|f - \partial_t u_{h\tau} + \nabla \cdot \bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{\mathcal{T}_F \times I_n} C_{T,n}^{-\frac{1}{2}} h_T^{-\frac{1}{2}} \|v_{F,n}\|_{F \times I_n} \\ & \quad + \frac{(u - u_{h\tau}, \partial_t(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n} - (\boldsymbol{\sigma}_d, \nabla(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n}}{\|\psi_{F,n} v_{F,n}\|_{Y, \mathcal{T}_F, n}} C_{T,n}^{-\frac{1}{2}} h_T^{-\frac{1}{2}} \|v_{F,n}\|_{F \times I_n}. \end{aligned}$$

Thus,

$$\begin{aligned} C_{T,n}^{-\frac{1}{2}} h_T^{-\frac{1}{2}} \|v_{F,n}\|_{F \times I_n} & \lesssim C_{T,n}^{-\frac{1}{2}} \|f - \partial_t u_{h\tau} + \nabla \cdot \bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})\|_{\mathcal{T}_F \times I_n} \\ & \quad + \frac{(u - u_{h\tau}, \partial_t(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n} - (\boldsymbol{\sigma}_d, \nabla(\psi_{F,n} v_{F,n}))_{\mathcal{T}_F \times I_n}}{\|\psi_{F,n} v_{F,n}\|_{Y, \mathcal{T}_F, n}}. \end{aligned}$$

Finally, the first term on the right-hand side is bounded using (6.1) for each  $T' \in \mathcal{T}_F$ , while proceeding as in step 1 for the second term yields (6.6).

*Step 3, conclusion.* Let now  $1 \leq n \leq N$  and an element  $T$  from the common coarsening  $\mathcal{T}^{n-1,n}$  be fixed. Combining (6.1) and (6.6) yields

$$\sum_{T' \in \overline{\mathcal{T}}^{n-1,n}, T' \subset T} (C_{T',n}^{-\frac{1}{2}} h_{T'}^{-1} \eta_{\text{clas}, T'}^n)^2 \lesssim \sum_{T' \in \mathcal{T}_T} (e_{\text{FR}, T'}^n + \eta_{\text{qd}, T'}^n)^2 + (\eta_{\text{NC}, T}^n)^2.$$

Using the triangle inequality and Assumption 4.1, we infer

$$(\eta_{\text{F}, T}^n)^2 \lesssim (\eta_{\text{qd}, T}^n)^2 + C_{T,n}^{-1} h_T^{-2} \sum_{T' \in \overline{\mathcal{T}}^{n-1,n}, T' \subset T} (\eta_{\text{clas}, T'}^n)^2.$$

Similarly, using the triangle inequality, an inverse space inequality, and Assumption 4.1 leads to

$$(\eta_{\text{R}, T}^n)^2 \lesssim C_{T,n}^{-1} h_T^{-2} \sum_{T' \in \overline{\mathcal{T}}^{n-1,n}, T' \subset T} (\eta_{\text{clas}, T'}^n)^2.$$

Combining the above inequalities we obtain

$$\eta_{\text{FR}, T}^n \lesssim \sum_{T' \in \mathcal{T}_T} (e_{\text{FR}, T'}^n + \eta_{\text{qd}, T'}^n) + \eta_{\text{NC}, T}^n.$$

Finally, invoking the requirement (4.4) on quadrature errors to discard the quantities  $\eta_{\text{qd}, T'}^n$  from the right-hand side yields the conclusion.  $\square$

## 6.2 Proof of Theorem 4.4

*Proof.* We use the notations  $v_{T,n}$  and  $v_{F,n}$  from the previous section. It is sufficient to show that

$$\left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 \right\}^{\frac{1}{2}} \lesssim \mathcal{J}_u(u_{h\tau}) + \eta_{\text{qd}}, \quad (6.9a)$$

$$\left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} \sum_{F \in \mathcal{F}_T^{\text{int}}} C_{T,n}^{-1} h_T^{-1} \|v_{F,n}\|_{F \times I_n}^2 \right\}^{\frac{1}{2}} \lesssim \mathcal{J}_u(u_{h\tau}) + \eta_{\text{qd}}, \quad (6.9b)$$

since choosing  $\gamma_{\text{qd}}$  in (4.6) small enough then yields (4.7). We only show (6.9a); (6.9b) follows by similar arguments. Let  $1 \leq n \leq N$  and  $T \in \overline{\mathcal{T}}^{n-1,n}$ . It follows from the first estimate in (6.5) and the Green theorem that

$$\begin{aligned} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 &\lesssim C_{T,n}^{-1} \{ (f - \partial_t u_{h\tau}, \psi_{T,n} v_{T,n})_{T \times I_n} - (\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}), \nabla(\psi_{T,n} v_{T,n}))_{T \times I_n} \\ &\quad + (\boldsymbol{\sigma}(u_{h\tau}, \nabla u_{h\tau}) - \overline{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}), \nabla(\psi_{T,n} v_{T,n}))_{T \times I_n} \}. \end{aligned}$$

Set  $\lambda|_{T \times I_n} := C_{T,n}^{-1} \psi_{T,n} v_{T,n}$  and observe that  $\lambda \in Y$ . Recall the notation  $R(u_{h\tau})$  from (2.7),  $\mathcal{J}_{u,\text{FR}}(u_{h\tau})$  from (2.6), and  $\eta_{\text{qd},T}^n$  from (4.2). Summing the above inequality over all  $1 \leq n \leq N$  and all  $T \in \overline{\mathcal{T}}^{n-1,n}$  and using the Green theorem, the Cauchy–Schwarz inequality, (6.3a), and  $\|\psi_{T,n}\|_{\infty,T,n} \leq 1$  yields

$$\begin{aligned} &\sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 \lesssim \langle R(u_{h\tau}), \lambda \rangle_{Y',Y} + \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} \eta_{\text{qd},T}^n C_{T,n}^{-\frac{1}{2}} \|v_{T,n}\|_{T \times I_n} \\ &\leq \frac{\langle R(u_{h\tau}), \lambda \rangle_{Y',Y}}{\|\lambda\|_Y} \|\lambda\|_Y + \eta_{\text{qd}} \left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 \right\}^{\frac{1}{2}} \\ &\leq \mathcal{J}_{u,\text{FR}}(u_{h\tau}) \|\lambda\|_Y + \eta_{\text{qd}} \left\{ \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2 \right\}^{\frac{1}{2}}. \end{aligned} \quad (6.10)$$

The definition (2.5b) of the  $\|\cdot\|_Y$ -norm, that of  $\lambda$ , the fact that we suppose that the ratio  $h_T/h_{T'}$ ,  $T \in \overline{\mathcal{T}}^{n-1,n}$ ,  $T' \in \overline{\mathcal{T}}^{n-1,n}$ ,  $T' \subset T$ , bounded and (6.4) lead to

$$\|\lambda\|_Y^2 = \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-2} \|\psi_{T,n} v_{T,n}\|_{Y,T \times I_n}^2 \lesssim \sum_{n=1}^N \sum_{T \in \overline{\mathcal{T}}^{n-1,n}} C_{T,n}^{-1} \|v_{T,n}\|_{T \times I_n}^2. \quad (6.11)$$

Combining (6.10) and (6.11) proves (6.9a).  $\square$

## 7 Application to the discontinuous Galerkin method

We apply here the abstract framework of Sections 3 and 4 to the discontinuous Galerkin method as an example of nonconforming space discretization scheme, with Crank–Nicolson time stepping. This consists in specifying the flux reconstruction  $\mathbf{t}_{h\tau}$  and in verifying Assumptions 3.1 and 4.1.

### 7.1 The discontinuous Galerkin method

For all  $0 \leq n \leq N$ , let  $V_h^n := \mathbb{P}_p(\mathcal{T}^n)$ ,  $p \geq 1$ , be spanned by piecewise polynomials on the mesh  $\mathcal{T}^n$  with total degree  $\leq p$ . Recall that  $\mathcal{F}^n$  collects all the mesh faces of  $\mathcal{T}^n$ . For each  $F \in \mathcal{F}^n$ , along with the jump operator  $[\![\cdot]\!]_F$ , we consider the (arithmetic) average operator  $\{\!\!\{ \cdot \}\!\!\}$  (conventionally yielding the actual value at boundary faces).

The space-time approximation  $u_{h\tau}$  is continuous and piecewise affine in time, and is defined by its values  $u_h^n$  at  $t^n$  for all  $0 \leq n \leq N$ . We take  $u_h^0$  as the  $L^2$ -orthogonal projection of  $u_0$  onto  $V_h^0$ . Then, for all  $1 \leq n \leq N$ , we look for  $u_h^n \in V_h^n$  such that

$$\begin{aligned} & (\partial_t u_h^n, v_h) + \frac{1}{2} \sum_{m=n-1}^n \left\{ (\boldsymbol{\sigma}(u_h^m, \nabla u_h^m), \nabla v_h) + \sum_{F \in \mathcal{F}^m} \alpha_{\underline{\mathbf{K}}, F}^m h_F^{-1} (\llbracket u_h^m \rrbracket, \llbracket v_h \rrbracket)_F \right. \\ & + \sum_{F \in \mathcal{F}^m} (H_F(u_h^m), \llbracket v_h \rrbracket)_F - \sum_{F \in \mathcal{F}^m} (\{\!\!\{ \underline{\mathbf{K}}(u_h^m) \nabla u_h^m \}\!\!\}, \mathbf{n}_F, \llbracket v_h \rrbracket)_F \\ & \left. - \theta \sum_{F \in \mathcal{F}^m} (\{\!\!\{ \underline{\mathbf{K}}(u_h^m) \nabla v_h \}\!\!\}, \mathbf{n}_F, \llbracket u_h^m \rrbracket)_F - (f^m, v_h) \right\} = 0 \quad \forall v_h \in V_h^n, \end{aligned} \quad (7.1)$$

where the parameter  $\theta$  is chosen in  $\{-1, 0, 1\}$  according to the variant of interior penalty Galerkin method and  $\partial_t u_h^n := (\tau^n)^{-1}(u_h^n - u_h^{n-1})$ ; we suppose for simplicity that  $f$  is continuous and piecewise affine in time and denote  $f^n := f(t^n)$ . The penalty coefficient  $\alpha_{\underline{\mathbf{K}}, F}^n$  can be taken as  $\sigma \|\underline{\mathbf{K}}(u_h^n)\|_{\infty, \mathcal{T}_F}$ , where  $\sigma$  only depends on mesh-regularity and the polynomial degree  $p$ . For problems with internal layers caused by locally small diffusion, using diffusion-dependent weighted averages and scaling  $\alpha_{\underline{\mathbf{K}}, F}^n$  with the harmonic average of the normal component of  $\underline{\mathbf{K}}(u_h^n)$  at  $F$  can be more effective, see [18]. Finally, the advection term in (7.1) has been discretized using a numerical flux  $H_F(u_h^n)$  satisfying the following reasonable assumption:

**Assumption 7.1** (Numerical flux for advection). *For all  $1 \leq n \leq N$ , all mesh faces  $F \in \mathcal{F}^n$ , and all  $v_h \in V_h^n$ , there holds*

$$\|H_F(v_h) - \{\!\!\{ \phi(v_h) \}\!\!\}, \mathbf{n}_F\|_F \lesssim \|\phi'(v_h)\|_{\infty, \mathcal{T}_F} \|\llbracket v_h \rrbracket\|_F.$$

An example is the numerical flux of Lax–Friedrichs type, which consists of the centered flux  $\{\!\!\{ \phi(\cdot) \}\!\!\}, \mathbf{n}_F$  supplemented by a stabilization term penalizing the interface jumps. In the numerical experiments of Section 8, we employ a numerical upwinding flux, see (8.1).

## 7.2 Flux reconstruction

We now specify the flux reconstruction  $\mathbf{t}_{h\tau}$  of Assumption 3.1. Let  $l \geq 0$ . We construct  $\mathbf{t}_{h\tau}$  continuous and piecewise affine in time, with  $\mathbf{t}_h^n := \mathbf{t}_{h\tau}(t^n)$ ,  $0 \leq n \leq N$ , in the Raviart–Thomas–Nédélec finite element space  $\mathbf{RTN}_l(\mathcal{T}^n)$  (on the mesh  $\mathcal{T}^n$ ). This space is defined as, cf. Brezzi and Fortin [8],  $\mathbf{RTN}_l(\mathcal{T}^n) := \{\mathbf{v}_h \in \mathbf{H}(\text{div}, \Omega); \mathbf{v}_h|_T \in \mathbf{RTN}_l(T) \text{ for all } T \in \mathcal{T}^n\}$ , where  $\mathbf{RTN}_l(T) := [\mathbb{P}_l(T)]^d + \mathbf{x}\mathbb{P}_l(T)$ . In particular,  $\mathbf{v}_h \in \mathbf{RTN}_l(\mathcal{T}^n)$  is such that  $\nabla \cdot \mathbf{v}_h \in \mathbb{P}_l(T)$  for all  $T \in \mathcal{T}^n$ ,  $\mathbf{v}_h \cdot \mathbf{n}_F \in \mathbb{P}_l(F)$  for all  $F \in \mathcal{F}^n$ , and such that its normal trace is continuous at all interfaces. Following Kim [23] and [16], we set:

**Definition 7.2** (Reconstructed flux  $\mathbf{t}_{h\tau}$ ). *Let  $l \geq 0$ . For all  $0 \leq n \leq N$ , we specify the degrees of freedom of  $\mathbf{t}_h^n \in \mathbf{RTN}_l(\mathcal{T}^n)$  by setting, for all  $T \in \mathcal{T}^n$ , all  $F \in \mathcal{F}_T$ , and all  $q_h \in \mathbb{P}_l(F)$ ,*

$$(\mathbf{t}_h^n \cdot \mathbf{n}_F, q_h)_F = (-\{\!\!\{ \underline{\mathbf{K}}(u_h^n) \nabla u_h^n \}\!\!\}, \mathbf{n}_F + \alpha_{\underline{\mathbf{K}}, F}^n h_F^{-1} \llbracket u_h^n \rrbracket, q_h)_F + (H_F(u_h^n), q_h)_F, \quad (7.2)$$

and all  $\mathbf{r}_h \in [\mathbb{P}_{l-1}(T)]^d$ ,

$$(\mathbf{t}_h^n, \mathbf{r}_h)_T = -(\underline{\mathbf{K}}(u_h^n) \nabla u_h^n, \mathbf{r}_h)_T + \theta \sum_{F \in \mathcal{F}_T} w_F (\underline{\mathbf{K}}(u_h^n) \mathbf{r}_h \cdot \mathbf{n}_F, \llbracket u_h^n \rrbracket)_F + (\phi(u_h^n), \mathbf{r}_h)_T, \quad (7.3)$$

where  $w_F := \frac{1}{2}$  for interfaces and  $w_F := 1$  for boundary faces.

**Remark 7.3** (Diffusive and advective fluxes). *We can split  $\mathbf{t}_h^n$  into  $\mathbf{t}_{D,h}^n + \mathbf{t}_{A,h}^n$ , where the diffusive flux  $\mathbf{t}_{D,h}^n$  is defined by the first two terms on the right-hand sides of (7.2) and (7.3), whereas the advective flux  $\mathbf{t}_{A,h}^n$  is defined by the last terms. Then,  $\mathbf{t}_{D,h}^n$  is a reconstruction of the diffusive flux  $\mathbf{t}_D := -\underline{\mathbf{K}}(u) \nabla u$  and  $\mathbf{t}_{A,h}^n$  is a reconstruction of the advective flux  $\mathbf{t}_A := \phi(u)$ .*

**Remark 7.4** (Alternative construction). *Instead of prescribing directly the degrees of freedom for  $\mathbf{t}_h^n$ , it is also possible to reconstruct the flux by solving local Neumann problems by mixed finite elements, see [19]. This approach can achieve a tighter relationship between the error and the estimator but is more expensive.*

**Remark 7.5** (Spatial and temporal errors). *In the present setting, from their behavior with respect to time,  $\eta_{F,T}^n$  and  $\eta_{\text{NC},T}^n$  can be used as spatial error estimators and  $\eta_{R,T}^n$  as a temporal error estimator.*

### 7.3 Verification of Assumption 3.1

**Lemma 7.6** (Local conservation). *Let  $u_h^n$ ,  $1 \leq n \leq N$ , solve (7.1) and let  $\mathbf{t}_{h\tau}$  be defined by (7.2)–(7.3) with  $l \geq 0$ . Then,  $\mathbf{t}_{h\tau}$  satisfies  $\mathbf{t}_{h\tau} \in \mathbf{L}^2(0, t_F; \mathbf{H}(\text{div}, \Omega))$  and, for all  $1 \leq n \leq N$  and all  $T \in \mathcal{T}^{n-1, n}$ ,*

$$(f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, v_h)_{T \times I_n} = 0 \quad \forall v_h \in \mathbb{P}_{\min(p, l)}(T), \quad (7.4)$$

so that, in particular, the local space-time conservation property (3.1) holds true.

*Proof.* By construction, we have, for  $0 \leq n \leq N$ ,  $\mathbf{t}_h^n \in \mathbf{H}(\text{div}, \Omega)$ , so that  $\mathbf{t}_{h\tau} \in \mathbf{L}^2(0, t_F; \mathbf{H}(\text{div}, \Omega))$ . Let  $1 \leq n \leq N$ ,  $T \in \mathcal{T}^{n-1, n}$ , and  $v_h \in \mathbb{P}_{\min(p, l)}(T)$  be given. Using that  $u_{h\tau}$ ,  $f$ , and  $\mathbf{t}_{h\tau}$  are affine in time on  $I_n$ , that  $T$  is from the common coarsening of the meshes  $\mathcal{T}^{n-1}$  and  $\mathcal{T}^n$ , and the Green theorem, we obtain

$$\begin{aligned} (f - \partial_t u_{h\tau} - \nabla \cdot \mathbf{t}_{h\tau}, v_h)_{T \times I_n} &= -\tau^n (\partial_t u_{h\tau}, v_h)_T + \frac{\tau^n}{2} \sum_{m=n-1}^n (f^m - \nabla \cdot \mathbf{t}_h^m, v_h)_T \\ &= -\tau^n \left( \partial_t u_{h\tau} - \frac{1}{2} \sum_{m=n-1}^n f^m, v_h \right)_T + \frac{\tau^n}{2} \sum_{m=n-1}^n \sum_{T' \in \mathcal{T}^m, T' \subset T} \{ (\mathbf{t}_h^m, \nabla v_h)_{T'} - (\mathbf{t}_h^m \cdot \mathbf{n}_{T'}, v_h)_{\partial T'} \}. \end{aligned}$$

Let  $m = n$  or  $m = n - 1$ . Since  $\nabla v_h \in [\mathbb{P}_{l-1}(T')]^d$  for any  $T'$ , (7.3) yields that  $(\mathbf{t}_h^m, \nabla v_h)_{T'}$  equals

$$-(\mathbf{K}(u_h^m) \nabla u_h^m, \nabla v_h)_{T'} + \theta \sum_{F \in \mathcal{F}_{T'}} w_F (\mathbf{K}(u_h^m) \nabla v_h \cdot \mathbf{n}_F, \llbracket u_h^m \rrbracket)_F + (\phi(u_h^m), \nabla v_h)_{T'}.$$

Furthermore, the fact that  $v_h|_F \in \mathbb{P}_l(F)$  for all  $F \in \mathcal{F}_{T'}$  and (7.2) yield that  $-(\mathbf{t}_h^m \cdot \mathbf{n}_{T'}, v_h)_{\partial T'}$  equals

$$\sum_{F \in \mathcal{F}_{T'}} \mathbf{n}_{T'} \cdot \mathbf{n}_F \{ -\llbracket \mathbf{K}(u_h^m) \nabla u_h^m \rrbracket \cdot \mathbf{n}_F + \alpha_{\mathbf{K}, F}^m h_F^{-1} \llbracket u_h^m \rrbracket, v_h \}_F + (H_F(u_h^m), v_h)_F.$$

Extending  $v_h$  by 0 outside  $T$  so that a function in  $V_h^n$  is obtained and using the above identities and the definition (7.1) of the discontinuous Galerkin scheme yields (7.4).  $\square$

### 7.4 Verification of Assumption 4.1

To verify Assumption 4.1, we need first to specify the space-time piecewise polynomial function  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)$ . Let  $0 \leq n \leq N$  and let  $T \in \mathcal{T}^n$ . We define  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)|_T \in \mathbf{RTN}_l(T)$  such that, for all  $F \in \mathcal{F}_T$  and all  $q_h \in \mathbb{P}_l(F)$ ,

$$(\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)|_T \cdot \mathbf{n}_F, q_h)_F = (\boldsymbol{\sigma}(u_h^n, \nabla u_h^n)|_T \cdot \mathbf{n}_F, q_h)_F, \quad (7.5)$$

and all  $\mathbf{r}_h \in [\mathbb{P}_{l-1}(T)]^d$ ,

$$(\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n), \mathbf{r}_h)_T = (\boldsymbol{\sigma}(u_h^n, \nabla u_h^n), \mathbf{r}_h)_T. \quad (7.6)$$

Here,  $l$  is the polynomial degree used for reconstructing the flux  $\mathbf{t}_h^n$  in Section 7.2. We observe that, locally in each mesh element  $T \in \mathcal{T}^n$ ,  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)|_T$  belongs to  $\mathbf{RTN}_l(T)$  (as  $\mathbf{t}_h^n$  does), but, globally,  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) \notin \mathbf{RTN}_l(\mathcal{T}^n)$  in general because the normal component of  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)$  is in general discontinuous across interfaces. Finally, the space-time function  $\bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})$  is taken to be continuous and piecewise affine in time, matching the values  $\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)$  at  $t^n$  for all  $0 \leq n \leq N$ .

**Lemma 7.7** (Flux approximation). *Let  $\mathbf{t}_{h\tau}$  be defined by (7.2)–(7.3) with  $l \geq 0$  and let  $\bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})$  be defined by (7.5)–(7.6). For  $1 \leq n \leq N$  and  $T \in \overline{\mathcal{T}}^{n-1, n}$ , define the constants  $C_{\mathbf{K}, \phi, T, F, n}$  of (2.11) by*

$$C_{\mathbf{K}, \phi, T, F, n} := (\alpha_{\mathbf{K}, F}^n)^2 h_F^{-1} + \|\mathbf{K}(u_h^n)\|_{\infty, T}^2 h_F^{-1} + \|\phi'(u_h^n)\|_{\infty, T_F}^2 h_F. \quad (7.7)$$

Then, Assumption 4.1 holds true.



*Proof.* Let  $0 \leq n \leq N$  and let  $T \in \mathcal{T}^n$ . Let  $F \in \mathcal{F}_T$ . Since  $\boldsymbol{\sigma}(u_h^n, \nabla u_h^n) = \mathbf{K}(u_h^n) \nabla u_h^n - \boldsymbol{\phi}(u_h^n)$  and using (7.5), we rewrite (7.2) as

$$\begin{aligned} (\mathbf{t}_h^n \cdot \mathbf{n}_F, q_h)_F &= (-\{\{\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n)\}\} \cdot \mathbf{n}_F + \alpha_{\mathbf{K}, F}^n h_F^{-1} \llbracket u_h^n \rrbracket, q_h)_F \\ &\quad + (H_F(u_h^n) - \{\{\boldsymbol{\phi}(u_h^n)\}\} \cdot \mathbf{n}_F, q_h)_F, \end{aligned} \quad (7.8)$$

for all  $q_h \in \mathbb{P}_l(F)$ , and using (7.6), we rewrite (7.3) as

$$(\mathbf{t}_h^n, \mathbf{r}_h)_T = -(\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n), \mathbf{r}_h)_T + \theta \sum_{F \in \mathcal{F}_T} w_F (\mathbf{K}(u_h^n) \mathbf{r}_h \cdot \mathbf{n}_F, \llbracket u_h^n \rrbracket)_F, \quad (7.9)$$

for all  $\mathbf{r}_h \in [\mathbb{P}_{l-1}(T)]^d$ . Set  $\mathbf{v}_h := \bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) + \mathbf{t}_h^n$  so that  $\mathbf{v}_h|_T \in \mathbf{RTN}_l(T)$ . Owing to (7.8),  $\mathbf{v}_h|_T \cdot \mathbf{n}_F$  can be rewritten as

$$(1 - w_F) \llbracket \bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) \rrbracket \cdot \mathbf{n}_F + \alpha_{\mathbf{K}, F}^n h_F^{-1} \Pi_{l,F}(\llbracket u_h^n \rrbracket) + \Pi_{l,F}(H_F(u_h^n) - \{\{\boldsymbol{\phi}(u_h^n)\}\} \cdot \mathbf{n}_F),$$

where  $\Pi_{l,F}$  denotes the  $L^2(F)$ -orthogonal projection onto  $\mathbb{P}_l(F)$ . Thus, using Assumption 7.1 on the numerical flux for advection, we infer

$$\|\mathbf{v}_h \cdot \mathbf{n}_F\|_F \lesssim (1 - w_F) \|\llbracket \bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) \rrbracket \cdot \mathbf{n}_F\|_F + (\alpha_{\mathbf{K}, F}^n h_F^{-1} + \|\boldsymbol{\phi}'(u_h^n)\|_{\infty, \mathcal{T}_F}) \|\llbracket u_h^n \rrbracket\|_F. \quad (7.10)$$

Moreover, owing to (7.9) and using an inverse inequality in space, we infer

$$(\mathbf{v}_h, \mathbf{r}_h)_T = \theta \sum_{F \in \mathcal{F}_T} w_F (\mathbf{K}(u_h^n) \mathbf{r}_h \cdot \mathbf{n}_F, \llbracket u_h^n \rrbracket)_F \lesssim \|\mathbf{K}(u_h^n)\|_{\infty, T} \|\mathbf{r}_h\|_T \sum_{F \in \mathcal{F}_T} h_F^{-\frac{1}{2}} \|\llbracket u_h^n \rrbracket\|_F. \quad (7.11)$$

Using (7.10) and (7.11) in the classical bound  $\|\mathbf{v}_h\|_T^2 \lesssim \sum_{F \in \mathcal{F}_T} h_F \|\mathbf{v}_h \cdot \mathbf{n}_F\|_F^2 + \left( \sup_{\mathbf{r}_h \in [\mathbb{P}_{l-1}(T)]^d} \frac{(\mathbf{v}_h, \mathbf{r}_h)_T}{\|\mathbf{r}_h\|_T} \right)^2$ , valid for all  $\mathbf{v}_h \in \mathbf{RTN}_l(T)$ , and owing to (7.7),

$$\|\bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) + \mathbf{t}_h^n\|_T^2 \lesssim \sum_{F \in \mathcal{F}_T^{\text{int}}} h_F \|\llbracket \bar{\boldsymbol{\sigma}}(u_h^n, \nabla u_h^n) \rrbracket \cdot \mathbf{n}_F\|_F^2 + \sum_{F \in \mathcal{F}_T} C_{\mathbf{K}, \boldsymbol{\phi}, T, F, n} \|\llbracket u_h^n \rrbracket\|_F^2. \quad (7.12)$$

Let now  $1 \leq n \leq N$  and  $T \in \mathcal{T}^{n-1, n}$ . Using that both  $\bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau})$  and  $u_{h\tau}$  are piecewise affine in time, we have, cf. [20, Lemma 6.1],

$$\begin{aligned} \|\bar{\boldsymbol{\sigma}}(u_{h\tau}, \nabla u_{h\tau}) + \mathbf{t}_{h\tau}\|_{T \times I_n}^2 &\lesssim \tau^n \sum_{m=n-1}^n \|\bar{\boldsymbol{\sigma}}(u_h^m, \nabla u_h^m) + \mathbf{t}_h^m\|_T^2 \\ &= \tau^n \sum_{m=n-1}^n \sum_{T' \in \mathcal{T}^m, T' \subset T} \|\bar{\boldsymbol{\sigma}}(u_h^m, \nabla u_h^m) + \mathbf{t}_h^m\|_{T'}^2 \lesssim \sum_{T' \in \overline{\mathcal{T}}^{n-1, n}, T' \subset T} (\eta_{\text{clas}, T'}^n)^2, \end{aligned}$$

employing (7.12) on  $T' \in \mathcal{T}^m$ ; actually, the first terms of  $\eta_{\text{clas}, T'}^n$  do not appear.  $\square$

**Remark 7.8** (Choice of the reconstruction degree  $l$ ). *A typical choice for the polynomial degree in the flux reconstruction  $\mathbf{t}_{h\tau}$  is  $l \in \{p-1, p\}$ . Choosing larger values for  $l$  might, however, be needed to satisfy the balancing criteria (4.4) or (4.6) with  $\gamma_{\text{qd}, T}$  or  $\gamma_{\text{qd}}$  small enough, as required, respectively, in Theorems 4.2 and 4.4.*

## 8 Numerical experiments

In this section, we present several numerical experiments illustrating the a posteriori error estimates of this paper. We consider the application to the discontinuous Galerkin method of Section 7.

## 8.1 Setting

We consider the problem (1.1a)–(1.1c) where we replace the homogeneous Dirichlet boundary condition (1.1b) by an inhomogeneous one; the additional error for nonpolynomial Dirichlet boundary data is neglected. We employ the discontinuous Galerkin method (7.1) with  $\theta = 0$  and the upwind numerical flux

$$H_F(u_h^n)|_F = \begin{cases} \phi(u_h^n)|_F^L \cdot \mathbf{n}_F & \text{if } A > 0, \\ \phi(u_h^n)|_F^R \cdot \mathbf{n}_F & \text{if } A \leq 0, \end{cases} \quad F \in \mathcal{F}^n, \quad n = 1, \dots, N, \quad (8.1)$$

where  $A := \{\{\phi'(u_h^n)\}\}_F \cdot \mathbf{n}_F$  and  $v_h^n|_F^L$  and  $v_h^n|_F^R$  denote traces of a function  $v_h^n$  on  $F \in \mathcal{F}^n$  from the direction and the opposite direction of  $\mathbf{n}_F$ , respectively. The penalty parameter  $\alpha_{\mathbf{K},F}^n$  in (7.1) is chosen as  $\alpha_{\mathbf{K},F}^n := 10p^2 \|\mathbf{K}(u_h^n)\|_{\infty, \mathcal{T}_F}$ , following [11] and Houston *et al.* [21]. We test  $p = 1, 2, 3$  (we employ the notation  $P_1, P_2, P_3$ ). For the linearization of (7.1), we use the Newton-like method of [13] where the (approximate) construction of the Jacobian matrix is avoided via the idea of easy-to-evaluate flux matrix. The volume integrals are evaluated with the aid of the Dunavant quadrature rule [14] of order  $3p + 2$ , whereas the face integrals with the aid of the Gauss quadrature rule with  $2p + 2$  nodes.

Since the error measure  $\mathcal{J}_{u, \text{FR}}(u_{h\tau})$  cannot be computed easily in practice even if the exact solution  $u$  is known, see the discussion in Section 2.3.1, we deal with its upper bound  $e_{\text{FR}}$ , see (2.9)–(2.10). The fully computable part  $\mathcal{J}_{u, \text{NC}}(u_{h\tau})$  is denoted herein by  $e_{\text{NC}}$  and we consider the error  $e := e_{\text{FR}} + e_{\text{NC}}$  (observe that  $e_{\text{NC}} = \eta_{\text{NC}}$ ). We evaluate the effectivity index  $i_{\text{eff}} := \frac{e}{e}$ , where  $\eta := \eta_{\text{FR}} + \eta_{\text{NC}} + \eta_{\text{IC}}$ . Since  $e$  is an upper bound on the actual error measure  $\mathcal{J}_u(u_{h\tau})$ ,  $i_{\text{eff}}$  can become smaller than one.

We consider square domains with uniform discretizations characterized by the space and time steps  $h_m$  and  $\tau_m$ ,  $m = 1, 2, 3$ , respectively. We choose  $\{h_1, \tau_1\}$  and then set  $h_{m+1} = h_m/2$ ,  $\tau_{m+1} = \tau_m/2^p$  for  $m = 2, 3$ . The space grids are triangulations with right-angled triangles resulting from diagonal cuttings of squares with edges equal to  $h_l = h_T/\sqrt{2}$ . We evaluate the experimental order of convergence  $\text{EOC} := \frac{\log(E_m/E_{m-1})}{\log(h_m/h_{m-1})}$ ,  $m = 2, 3$ , where  $E_m$  is either an error or an error estimator on the space-time discretization  $\{h_m, \tau_m\}$ . In order to verify that the error and its estimate are distributed in the same way in the computational domain, we also define, for all  $1 \leq n \leq N$  and all  $T \in \mathcal{T}^n$ , the errors  $e_T^n := e_{\text{FR}, T}^n + \eta_{\text{NC}, T}^n$  and the error estimators  $\eta_T^n := \eta_{\text{FR}, T}^n + \eta_{\text{NC}, T}^n$ . All the computations are performed in double precision on a PC with Intel Pentium 4-M 2.5 GHz processor and Linux operating system. Machine precision is  $10^{-15}$ .

## 8.2 Robustness with respect to advection dominance, final time, and discretization parameters

We consider first a model advection-diffusion problem with linear diffusive but nonlinear advective term

$$\partial_t u - \nabla \cdot (\varepsilon \nabla u - \phi(u)) = 0 \quad \text{in } Q, \quad (8.2)$$

where  $\varepsilon > 0$ ,  $\phi(u) = (u^2/2, u^2/2)^T$ ,  $\Omega = (-1, 1) \times (-1, 1)$ , and  $t_F = 1$ . The initial and boundary conditions are chosen in such a way that the exact solution is, see [12],

$$u(x, y, t) = \left( 1 + \exp\left(\frac{x + y + 1 - t}{2\varepsilon}\right) \right)^{-1}. \quad (8.3)$$

This problem exhibits an inner layer moving in the diagonal direction. The steepness of the layer increases with decreasing  $\varepsilon$ .

### 8.2.1 Robustness with respect to advection dominance

We first employ the initial space-time step  $\{h_1, \tau_1\} = \{1/6, 0.01\}$  and compare two values  $\varepsilon = 10^{-2}$  and  $\varepsilon = 10^{-4}$ . The errors  $e$ , error estimators  $\eta$ , effectivity indices  $i_{\text{eff}}$ , and EOC are summarized in Tables 1–2. We observe that  $i_{\text{eff}}$  behaves similarly for  $\varepsilon = 10^{-2}$  and  $\varepsilon = 10^{-4}$ , which indicates robustness with respect to advection dominance. Figures 2–3 show the corresponding diagonal cuts (bottom left to top right) of  $u_{h\tau}$  at  $t = t_F$ . We observe a sharp capturing of the solution for increasing  $p$  and decreasing  $h$  and  $\tau$ .

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	1	5.03E-02	9.25E-03	2.44E-03	9.37E-04	3.14E-02	4.39E-02	0.7377
8.33E-02	7.07E-03	1	2.21E-02	5.46E-03	1.89E-03	6.40E-04	1.47E-02	2.26E-02	0.8199
	(EOC)		( 1.18)	( 0.76)	( 0.37)	( 0.55)	( 1.10)	( 0.96)	
4.17E-02	5.00E-03	1	1.13E-02	3.13E-03	1.31E-03	3.76E-04	3.36E-03	8.17E-03	0.5646
	(EOC)		( 0.97)	( 0.80)	( 0.53)	( 0.77)	( 2.13)	( 1.47)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	2	3.06E-02	5.20E-03	1.57E-03	1.32E-03	1.34E-02	2.13E-02	0.5953
8.33E-02	5.00E-03	2	1.16E-02	1.77E-03	5.59E-04	4.07E-04	1.79E-03	4.48E-03	0.3350
	(EOC)		( 1.40)	( 1.56)	( 1.49)	( 1.70)	( 2.90)	( 2.25)	
4.17E-02	2.50E-03	2	5.19E-03	6.73E-04	1.63E-04	1.19E-04	3.05E-04	1.24E-03	0.2114
	(EOC)		( 1.16)	( 1.39)	( 1.78)	( 1.77)	( 2.55)	( 1.85)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	3	2.31E-02	3.91E-03	9.11E-04	1.69E-03	3.01E-03	9.39E-03	0.3472
8.33E-02	3.54E-03	3	7.41E-03	9.22E-04	1.57E-04	3.26E-04	6.20E-04	1.98E-03	0.2380
	(EOC)		( 1.64)	( 2.09)	( 2.54)	( 2.38)	( 2.28)	( 2.24)	
4.17E-02	1.25E-03	3	2.55E-03	2.54E-04	1.80E-05	3.81E-05	8.88E-05	3.95E-04	0.1406
	(EOC)		( 1.54)	( 1.86)	( 3.12)	( 3.10)	( 2.80)	( 2.33)	

Table 1: Example (8.2)–(8.3),  $\varepsilon = 10^{-2}$ : errors, error estimators, and effectivity indices

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	1	2.05E-01	2.99E-02	4.75E-03	3.44E-03	6.64E-17	3.79E-02	0.1612
8.33E-02	7.07E-03	1	1.51E-01	2.96E-02	4.70E-03	4.96E-03	1.36E-12	3.90E-02	0.2160
	(EOC)		( 0.44)	( 0.01)	( 0.02)	( -0.53)	( -14.32)	( -0.04)	
4.17E-02	5.00E-03	1	9.07E-02	2.82E-02	4.50E-03	6.71E-03	2.41E-07	3.90E-02	0.3284
	(EOC)		( 0.73)	( 0.07)	( 0.06)	( -0.43)	( -17.44)	( -0.00)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	2	1.29E-01	3.07E-02	3.84E-03	8.19E-03	7.15E-05	4.22E-02	0.2649
8.33E-02	5.00E-03	2	8.85E-02	2.04E-02	2.63E-03	5.63E-03	1.62E-03	2.99E-02	0.2744
	(EOC)		( 0.54)	( 0.59)	( 0.54)	( 0.54)	( -4.50)	( 0.50)	
4.17E-02	2.50E-03	2	5.69E-02	1.30E-02	1.76E-03	3.60E-03	5.79E-03	2.38E-02	0.3411
	(EOC)		( 0.64)	( 0.65)	( 0.58)	( 0.64)	( -1.83)	( 0.33)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.67E-01	1.00E-02	3	1.38E-01	3.46E-02	3.41E-03	1.40E-02	7.67E-10	5.14E-02	0.2978
8.33E-02	3.54E-03	3	8.54E-02	1.49E-02	1.60E-03	4.93E-03	4.38E-06	2.11E-02	0.2108
	(EOC)		( 0.69)	( 1.22)	( 1.09)	( 1.50)	( -12.48)	( 1.28)	
4.17E-02	1.25E-03	3	4.86E-02	5.91E-03	7.38E-04	1.46E-03	2.90E-04	8.26E-03	0.1514
	(EOC)		( 0.81)	( 1.34)	( 1.12)	( 1.76)	( -6.05)	( 1.36)	

Table 2: Example (8.2)–(8.3),  $\varepsilon = 10^{-4}$ : errors, error estimators, and effectivity indices

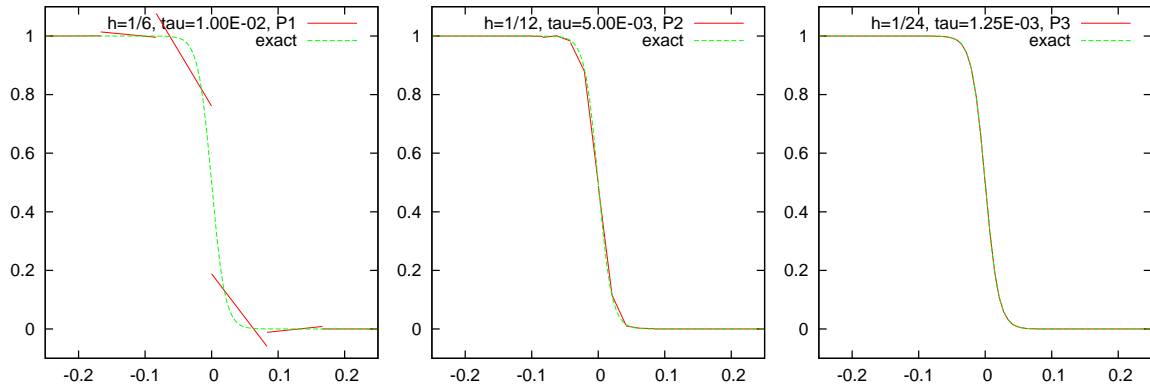


Figure 2: Example (8.2)–(8.3),  $\varepsilon = 10^{-2}$ : comparison of the exact and approximate solutions along diagonal cut,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

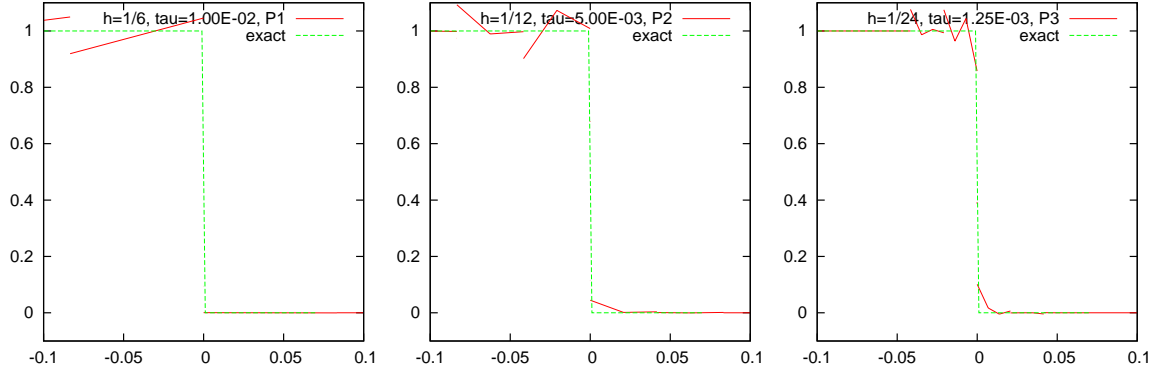


Figure 3: Example (8.2)–(8.3),  $\varepsilon = 10^{-4}$ : comparison of the exact and approximate solutions along diagonal cut,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

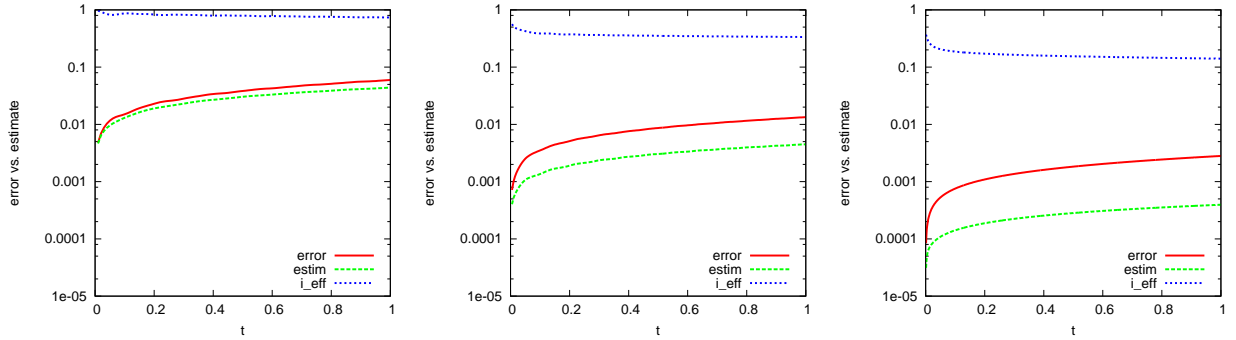


Figure 4: Example (8.2)–(8.3),  $\varepsilon = 10^{-2}$ : dependence of the error  $e$ , the error estimate  $\eta$ , and the effectivity index  $i_{\text{eff}}$  on the final time  $t_F$ ,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

### 8.2.2 Robustness with respect to the final time

We still employ the initial space-time step  $\{h_1, \tau_1\} = \{1/6, 0.01\}$ , but let the final time  $t_F$  vary. Figures 4–5 show the dependence of the error  $e$ , the error estimator  $\eta$ , and the effectivity index  $i_{\text{eff}}$  on  $t_F$ . We observe that  $i_{\text{eff}}$  is almost constant (more precisely slightly decreasing) over the interval  $(0, t_F)$ , which indicates that the present approach is robust with respect to the final time  $t_F$ .

### 8.2.3 Robustness with respect to discretization parameters

We employ here a different initial space-time step  $\{h_1, \tau_1\} = \{1/6, 0.1\}$ , with  $\varepsilon = 10^{-2}$ . We present the results in Table 3. We observe a very similar behavior to in Table 1, which confirms that our estimates do not require any type of matching between the space and time steps, i.e., robustness with respect to discretization parameters.

## 8.3 Degenerate parabolic problems and robustness with respect to size of non-linearity

We consider here two more model problems.

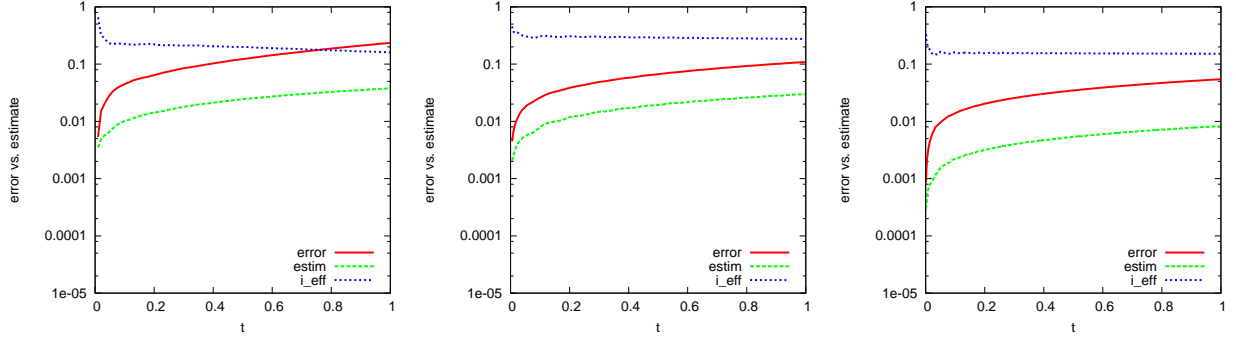


Figure 5: Example (8.2)–(8.3),  $\varepsilon = 10^{-4}$ : dependence of the error  $e$ , the error estimate  $\eta$ , and the effectivity index  $i_{\text{eff}}$  on the final time  $t_F$ ,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{\text{FR}}$	$\ e\ _{\text{NC}} = \eta_{\text{NC}}$	$\eta_{\text{F}}$	$\eta_{\text{R}}$	$\eta_{\text{IC}}$	$\eta$	$i_{\text{eff}}$
1.67E-01	1.00E-01	1	1.59E-01	6.85E-02	1.82E-02	8.12E-02	2.52E-02	1.93E-01	0.8483
8.33E-02	7.07E-02	1	1.17E-01	4.03E-02	1.33E-02	5.36E-02	1.04E-02	1.17E-01	0.7444
	(EOC)		(0.44)	(0.77)	(0.45)	(0.60)	(1.28)	(0.72)	
4.17E-02	5.00E-02	1	7.77E-02	2.96E-02	7.98E-03	2.50E-02	1.85E-03	6.41E-02	0.5973
	(EOC)		(0.59)	(0.45)	(0.74)	(1.10)	(2.50)	(0.87)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{\text{FR}}$	$\ e\ _{\text{NC}} = \eta_{\text{NC}}$	$\eta_{\text{F}}$	$\eta_{\text{R}}$	$\eta_{\text{IC}}$	$\eta$	$i_{\text{eff}}$
1.67E-01	1.00E-01	2	1.60E-01	7.23E-02	1.17E-02	1.13E-01	1.07E-02	2.07E-01	0.8926
8.33E-02	5.00E-02	2	8.90E-02	5.03E-02	4.67E-03	3.71E-02	1.44E-03	9.31E-02	0.6688
	(EOC)		(0.84)	(0.52)	(1.32)	(1.61)	(2.90)	(1.15)	
4.17E-02	2.50E-02	2	4.58E-02	3.60E-02	1.55E-03	1.00E-02	2.43E-04	4.76E-02	0.5818
	(EOC)		(0.96)	(0.48)	(1.59)	(1.89)	(2.57)	(0.97)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{\text{FR}}$	$\ e\ _{\text{NC}} = \eta_{\text{NC}}$	$\eta_{\text{F}}$	$\eta_{\text{R}}$	$\eta_{\text{IC}}$	$\eta$	$i_{\text{eff}}$
1.67E-01	1.00E-01	3	1.61E-01	1.41E-01	8.67E-03	1.21E-01	2.42E-03	2.73E-01	0.9009
8.33E-02	3.54E-02	3	6.84E-02	6.23E-02	2.02E-03	2.83E-02	5.56E-04	9.31E-02	0.7120
	(EOC)		(1.23)	(1.18)	(2.10)	(2.10)	(2.12)	(1.55)	
4.17E-02	1.25E-02	3	2.47E-02	2.32E-02	2.65E-04	3.54E-03	8.30E-05	2.70E-02	0.5641
	(EOC)		(1.47)	(1.43)	(2.93)	(3.00)	(2.74)	(1.78)	

Table 3: Example (8.2)–(8.3),  $\varepsilon = 10^{-2}$  with  $\{h_1, \tau_1\} = \{1/6, 0.1\}$ : errors, error estimators, and effectivity indices

### 8.3.1 Degenerate parabolic problem

First, we consider a nonlinear degenerate advection-diffusion problem from Kačur [22]

$$\partial_t u - \nabla \cdot (\mathbf{K}(u) \nabla u - \phi(u)) = 0 \quad \text{in } Q, \quad (8.4)$$

where  $\mathbf{K}(u) = 2\varepsilon u \mathbb{I}$  ( $\mathbb{I}$  being the identity matrix),  $\varepsilon = 10^{-2}$ ,  $\phi(u) = v(u^2, 0)^T$ ,  $v = 1/2$ ,  $\Omega = (0, 1) \times (0, 1)$ , and  $t_F = 1$ . The initial and boundary conditions are chosen in such a way that the exact solution is

$$u(x, y, t) = \begin{cases} 1 - \exp\left(\frac{v(x-vt-x_0)}{2\varepsilon}\right) & \text{for } x \leq vt + x_0, \\ 0 & \text{for } x > vt + x_0, \end{cases} \quad (8.5)$$

where  $x_0 = 1/4$  is the initial position of the front. If  $u = 0$  (for  $x \geq vt + x_0$ ), the diffusive term degenerates. The initial space-time step is  $\{h_1, \tau_1\} = \{1/8, 0.01\}$ .

The errors  $e$ , error estimators  $\eta$ , effectivity indices  $i_{\text{eff}}$ , and EOC are summarized in Table 4. We again observe that the effectivity index  $i_{\text{eff}}$  does not change too much for all computations. Figure 6 shows cuts of  $u_{h\tau}$  at  $y = 0.5$  and  $t = t_F$ ; accurate capturing of the exact solution is observed. Finally, Figure 7 shows the distribution of the local error  $e_T^n$  and the local error estimator  $\eta_T^n$  at  $t = t_F$  for  $P_2$  approximation with the space-time step  $\{h_2, \tau_2\}$ . We observe a close agreement in the distributions of error and error estimator.

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.25E-01	1.00E-02	1	4.25E-02	6.01E-03	2.18E-03	2.59E-04	2.45E-02	3.29E-02	0.6787
6.25E-02	7.07E-03	1	2.18E-02	3.75E-03	1.77E-03	2.53E-04	8.57E-03	1.43E-02	0.5577
	(EOC)		( 0.96)	( 0.68)	( 0.30)	( 0.03)	( 1.52)	( 1.21)	
3.12E-02	5.00E-03	1	1.03E-02	2.35E-03	1.22E-03	2.49E-04	2.37E-03	6.09E-03	0.4815
	(EOC)		( 1.08)	( 0.68)	( 0.54)	( 0.03)	( 1.86)	( 1.23)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.25E-01	1.00E-02	2	2.49E-02	3.98E-03	1.37E-03	8.74E-04	6.35E-03	1.23E-02	0.4258
6.25E-02	5.00E-03	2	1.07E-02	1.87E-03	5.97E-04	3.12E-04	1.11E-03	3.79E-03	0.3025
	(EOC)		( 1.23)	( 1.09)	( 1.19)	( 1.48)	( 2.52)	( 1.70)	
3.12E-02	2.50E-03	2	4.36E-03	8.82E-04	2.37E-04	9.92E-05	1.54E-04	1.34E-03	0.2561
	(EOC)		( 1.29)	( 1.09)	( 1.33)	( 1.65)	( 2.85)	( 1.50)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.25E-01	1.00E-02	3	1.72E-02	4.00E-03	8.83E-04	1.03E-03	1.19E-03	6.90E-03	0.3249
6.25E-02	3.54E-03	3	6.57E-03	1.37E-03	2.64E-04	1.45E-04	1.07E-04	1.86E-03	0.2337
	(EOC)		( 1.39)	( 1.54)	( 1.74)	( 2.84)	( 3.48)	( 1.89)	
3.12E-02	1.25E-03	3	2.40E-03	4.65E-04	7.58E-05	2.58E-05	7.47E-06	5.68E-04	0.1981
	(EOC)		( 1.45)	( 1.56)	( 1.80)	( 2.49)	( 3.84)	( 1.71)	

Table 4: Example (8.4)–(8.5): errors, error estimators, and effectivity indices

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E-00	1.00E-02	1	1.53E-01	4.67E-03	3.05E-03	1.07E-03	8.87E-02	9.70E-02	0.6155
5.00E-01	7.07E-03	1	5.62E-02	2.71E-03	1.98E-03	1.68E-03	3.30E-02	3.88E-02	0.6586
	(EOC)		( 1.44)	( 0.79)	( 0.62)	( -0.65)	( 1.43)	( 1.32)	
2.50E-01	5.00E-03	1	2.08E-02	1.83E-03	1.36E-03	2.31E-03	1.19E-02	1.69E-02	0.7463
	(EOC)		( 1.43)	( 0.57)	( 0.54)	( -0.46)	( 1.47)	( 1.20)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E-00	1.00E-02	2	5.29E-02	3.76E-03	1.34E-03	1.01E-02	4.12E-02	5.63E-02	0.9950
5.00E-01	5.00E-03	2	1.94E-02	2.16E-03	7.46E-04	6.81E-03	1.38E-02	2.34E-02	1.0854
	(EOC)		( 1.44)	( 0.80)	( 0.84)	( 0.57)	( 1.58)	( 1.27)	
2.50E-01	2.50E-03	2	7.07E-03	1.57E-03	5.36E-04	5.30E-03	5.36E-03	1.28E-02	1.4760
	(EOC)		( 1.46)	( 0.46)	( 0.48)	( 0.36)	( 1.36)	( 0.88)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E-00	1.00E-02	3	3.11E-02	7.79E-03	2.25E-03	3.36E-02	2.40E-02	6.76E-02	1.7411
5.00E-01	3.54E-03	3	1.10E-02	2.80E-03	8.26E-04	1.23E-02	8.55E-03	2.44E-02	1.7722
	(EOC)		( 1.50)	( 1.48)	( 1.45)	( 1.46)	( 1.49)	( 1.47)	
2.50E-01	1.25E-03	3	3.94E-03	1.04E-03	3.07E-04	4.60E-03	3.19E-03	9.14E-03	1.8356
	(EOC)		( 1.48)	( 1.43)	( 1.43)	( 1.41)	( 1.42)	( 1.42)	

Table 5: Example (8.6)–(8.7),  $m = 2$ : errors, error estimators, and effectivity indices

$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E+00	1.00E-02	1	5.62E-01	1.65E-02	8.22E-03	4.89E-03	4.11E-01	4.40E-01	0.7597
5.00E-01	7.07E-03	1	3.51E-01	1.50E-02	7.25E-03	1.29E-02	2.20E-01	2.53E-01	0.6920
	(EOC)		( 0.68)	( 0.13)	( 0.18)	( -1.40)	( 0.91)	( 0.80)	
2.50E-01	5.00E-03	1	2.01E-01	1.54E-02	7.01E-03	2.69E-02	1.32E-01	1.80E-01	0.8317
	(EOC)		( 0.80)	( -0.03)	( 0.05)	( -1.06)	( 0.73)	( 0.49)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E+00	1.00E-02	2	3.62E-01	2.51E-02	8.69E-03	7.35E-02	3.08E-01	4.14E-01	1.0703
5.00E-01	5.00E-03	2	2.08E-01	2.12E-02	7.04E-03	6.98E-02	1.41E-01	2.39E-01	1.0426
	(EOC)		( 0.80)	( 0.25)	( 0.30)	( 0.07)	( 1.13)	( 0.80)	
2.50E-01	2.50E-03	2	1.17E-01	2.26E-02	7.40E-03	7.94E-02	8.41E-02	1.93E-01	1.3828
	(EOC)		( 0.82)	( -0.10)	( -0.07)	( -0.19)	( 0.74)	( 0.30)	
$h$	$\tau$	$\mathbb{P}_k$	$\ e\ _{FR}$	$\ e\ _{NC = \eta_{NC}}$	$\eta_F$	$\eta_R$	$\eta_{IC}$	$\eta$	$i_{eff}$
1.00E+00	1.00E-02	3	2.76E-01	6.26E-02	1.84E-02	2.79E-01	2.20E-01	5.80E-01	1.7134
5.00E-01	3.54E-03	3	1.55E-01	3.28E-02	9.41E-03	1.44E-01	1.13E-01	2.98E-01	1.5916
	(EOC)		( 0.83)	( 0.93)	( 0.97)	( 0.95)	( 0.97)	( 0.96)	
2.50E-01	1.25E-03	3	8.60E-02	1.83E-02	5.38E-03	8.08E-02	6.91E-02	1.73E-01	1.6626
	(EOC)		( 0.85)	( 0.84)	( 0.81)	( 0.83)	( 0.70)	( 0.78)	

Table 6: Example (8.6)–(8.7),  $m = 4$ : errors, error estimators, and effectivity indices

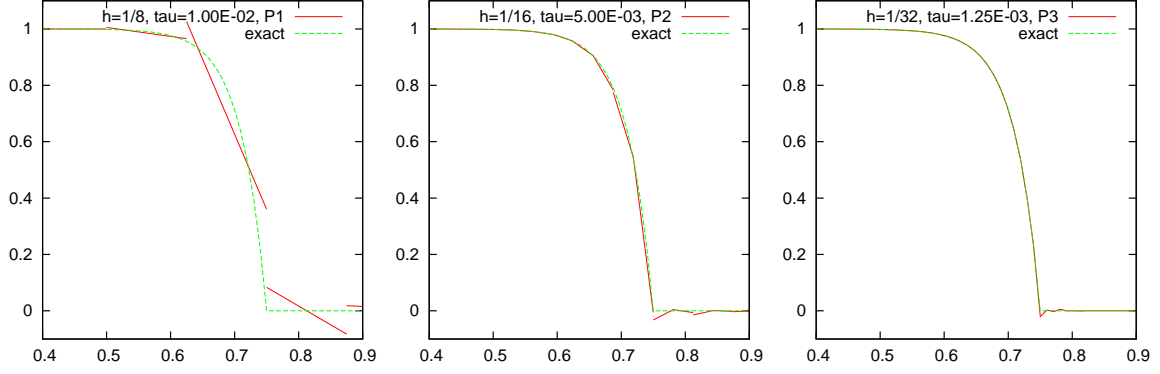


Figure 6: Example (8.4)–(8.5): comparison of the exact and approximate solutions along cut  $y = 0.5$ ,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

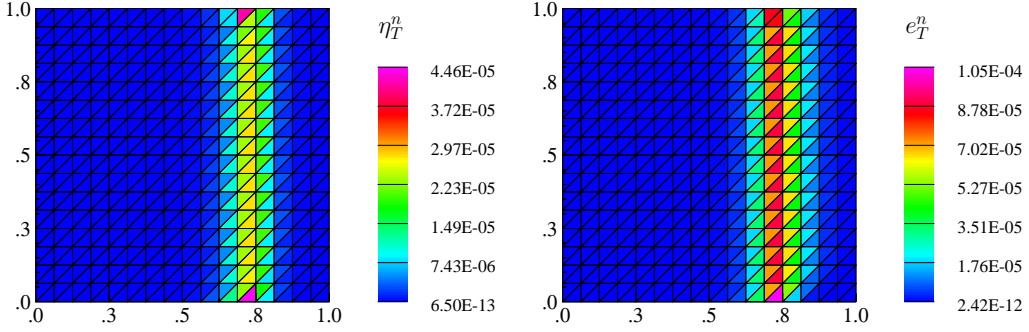


Figure 7: Example (8.4)–(8.5): distribution of the element error estimate  $\eta_T^n$  (left) and of the local error  $e_T^n$  (right) at  $t = t_F$  for  $P_2$  approximation on  $\{h_2, \tau_2\}$

### 8.3.2 Robustness with respect to size of nonlinearity

The last example is the porous medium equation

$$\partial_t u - \nabla \cdot (\underline{\mathbf{K}}(u) \nabla u) = 0 \quad \text{in } Q, \quad (8.6)$$

where  $\underline{\mathbf{K}}(u) = m|u|^{m-1}\mathbb{I}$ ,  $\Omega = (-6, 6) \times (-6, 6)$ ,  $t_F = 1$ , and  $m > 1$ . The initial and boundary conditions are chosen in such a way that the exact solution is, see Barenblatt [4] or Radu *et al.* [34],

$$u(x, y, t) = \left\{ \frac{1}{t+1} \left[ 1 - \frac{m-1}{4m^2} \frac{x^2 + y^2}{(t+1)^{1/m}} \right]_+^{\frac{m}{m-1}} \right\}^{\frac{1}{m}}, \quad (8.7)$$

where  $[a]_+ = \max(a, 0)$ ,  $a \in \mathbb{R}$ . In order to demonstrate robustness with respect to size of nonlinearity, we compare results for the values  $m = 2$  and  $m = 4$  (noticing that the case  $m = 4$  falls in fact beyond the adopted setting for the continuous problem). We employ the initial space-time step  $\{h_1, \tau_1\} = \{1, 0.01\}$ .

The errors  $e$ , error estimators  $\eta$ , effectivity indices  $i_{\text{eff}}$ , and EOC are summarized in Tables 5–6. We observe that the effectivity index  $i_{\text{eff}}$  almost does not change between  $m = 2$  and  $m = 4$ , indicating robustness with respect to size of nonlinearity. Figures 8 and 9 show cuts of  $u_{h\tau}$  at  $y = 0$  and  $t = t_F$ ; accurate capturing of the exact solution is observed. Finally, Figures 10 and 11 show the distribution of the local error  $e_T^n$  and the local error estimator  $\eta_T^n$  at  $t = t_F$  for  $P_2$  approximation with the space-time step  $\{h_2, \tau_2\}$ . We observe a close agreement in the distributions of error and error estimator.

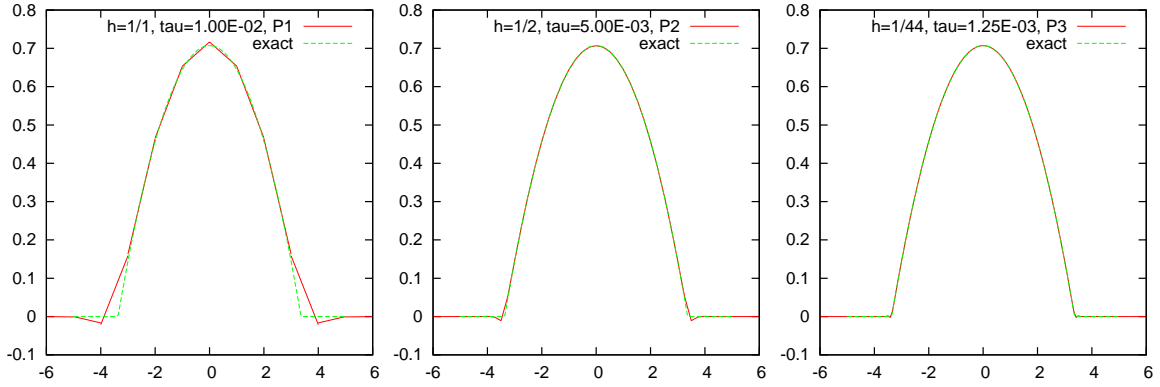


Figure 8: Example (8.6)–(8.7),  $m = 2$ : comparison of the exact and approximate solutions along cut  $y = 0$ ,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

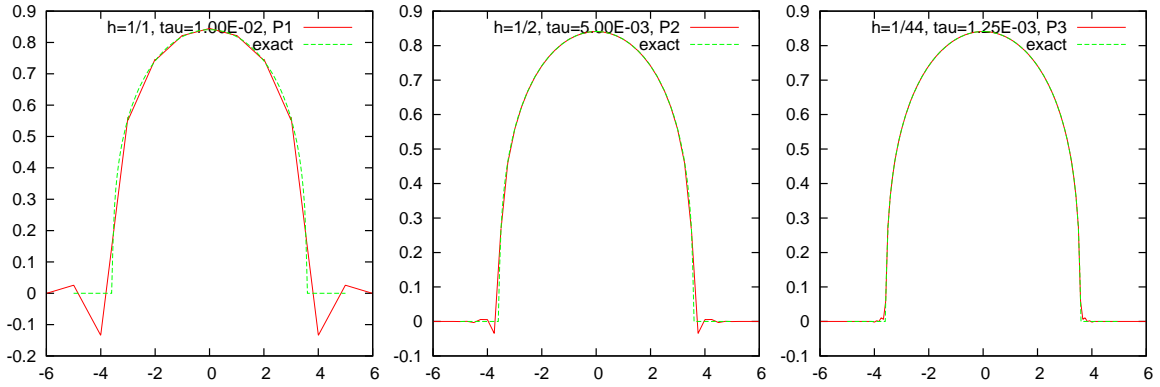


Figure 9: Example (8.6)–(8.7),  $m = 4$ : comparison of the exact and approximate solutions along cut  $y = 0$ ,  $P_1$  approximation on  $\{h_1, \tau_1\}$  (left),  $P_2$  approximation on  $\{h_2, \tau_2\}$  (center), and  $P_3$  approximation on  $\{h_3, \tau_3\}$  (right)

## References

- [1] M. AINSWORTH, *A framework for obtaining guaranteed error bounds for finite element approximations*, J. Comput. Appl. Math., 234 (2010), pp. 2618–2632.
- [2] H. W. ALT AND S. LUCKHAUS, *Quasilinear elliptic-parabolic differential equations*, Math. Z., 183 (1983), pp. 311–341.
- [3] D. N. ARNOLD, *An interior penalty finite element method with discontinuous elements*, SIAM J. Numer. Anal., 19 (1982), pp. 742–760.
- [4] G. I. BARENBLATT, *On some unsteady motions of a liquid and gas in a porous medium*, Akad. Nauk SSSR. Prikl. Mat. Meh., 16 (1952), pp. 67–78.
- [5] M. BEBENDORF, *A note on the Poincaré inequality for convex domains*, Z. Anal. Anwendungen, 22 (2003), pp. 751–756.
- [6] A. BERGAM, C. BERNARDI, AND Z. MGHAZLI, *A posteriori analysis of the finite element discretization of some parabolic equations*, Math. Comp., 74 (2005), pp. 1117–1138.



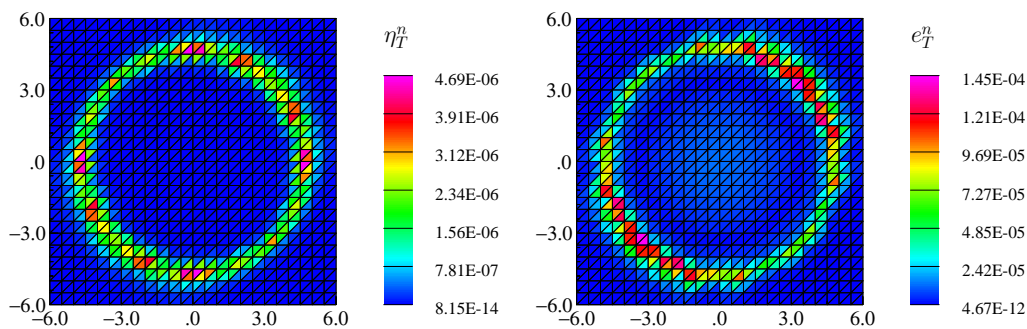


Figure 10: Example (8.6)–(8.7),  $m = 2$ : distribution of the element error estimate  $\eta_T^n$  (left) and of the local error  $e_T^n$  (right) at  $t = t_F$  for  $P_2$  approximation on  $\{h_2, \tau_2\}$

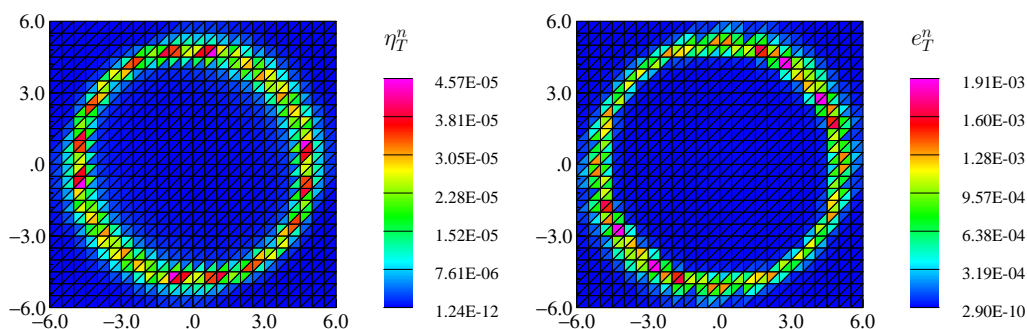


Figure 11: Example (8.6)–(8.7),  $m = 4$ : distribution of the element error estimate  $\eta_T^n$  (left) and of the local error  $e_T^n$  (right) at  $t = t_F$  for  $P_2$  approximation on  $\{h_2, \tau_2\}$

- [7] D. BRAESS, V. PILLWEIN, AND J. SCHÖBERL, *Equilibrated residual error estimates are p-robust*, Comput. Methods Appl. Mech. Engrg., 198 (2009), pp. 1189–1197.
- [8] F. BREZZI AND M. FORTIN, *Mixed and hybrid finite element methods*, vol. 15 of Springer Series in Computational Mathematics, Springer-Verlag, New York, 1991.
- [9] A. L. CHAILLOU AND M. SURI, *Computable error estimators for the approximation of nonlinear problems by linearized models*, Comput. Methods Appl. Mech. Engrg., 196 (2006), pp. 210–224.
- [10] ———, *A posteriori estimation of the linearization error for strongly monotone nonlinear operators*, J. Comput. Appl. Math., 205 (2007), pp. 72–87.
- [11] V. DOLEJŠÍ, *Analysis and application of the IIPG method to quasilinear nonstationary convection-diffusion problems*, J. Comput. Appl. Math., 222 (2008), pp. 251–273.
- [12] V. DOLEJŠÍ, M. FEISTAUER, V. KUČERA, AND V. SOBOTÍKOVÁ, *An optimal  $L^\infty(L^2)$ -error estimate of the discontinuous Galerkin method for a nonlinear nonstationary convection-diffusion problem*, IMA J. Numer. Anal., 28 (2008), pp. 496–521.
- [13] V. DOLEJŠÍ, M. HOLÍK, AND J. HOZMAN, *Efficient solution strategy for the semi-implicit discontinuous Galerkin discretization of the Navier–Stokes equations*, J. Comput. Phys., 230 (2011), pp. 4176–4200.
- [14] D. A. DUNAVANT, *High degree efficient symmetrical Gaussian quadrature rules for the triangle*, Internat. J. Numer. Methods Engrg., 21 (1985), pp. 1129–1148.
- [15] L. EL ALAOU, A. ERN, AND M. VOHRALÍK, *Guaranteed and robust a posteriori error estimates and balancing discretization and linearization errors for monotone nonlinear problems*, Comput. Methods Appl. Mech. Engrg., 200 (2011), pp. 2782–2795.

- [16] A. ERN, S. NICAISE, AND M. VOHRALÍK, *An accurate  $\mathbf{H}(\text{div})$  flux reconstruction for discontinuous Galerkin approximations of elliptic problems*, C. R. Math. Acad. Sci. Paris, 345 (2007), pp. 709–712.
- [17] A. ERN, A. F. STEPHANSEN, AND M. VOHRALÍK, *Guaranteed and robust discontinuous Galerkin a posteriori error estimates for convection–diffusion–reaction problems*, J. Comput. Appl. Math., 234 (2010), pp. 114–130.
- [18] A. ERN, A. F. STEPHANSEN, AND P. ZUNINO, *A discontinuous Galerkin method with weighted averages for advection-diffusion equations with locally small and anisotropic diffusivity*, IMA J. Numer. Anal., 29 (2009), pp. 235–256.
- [19] A. ERN AND M. VOHRALÍK, *Flux reconstruction and a posteriori error estimation for discontinuous Galerkin methods on general nonmatching grids*, C. R. Math. Acad. Sci. Paris, 347 (2009), pp. 441–444.
- [20] ———, *A posteriori error estimation based on potential and flux reconstruction for the heat equation*, SIAM J. Numer. Anal., 48 (2010), pp. 198–223.
- [21] P. HOUSTON, J. ROBSON, AND E. SÜLI, *Discontinuous Galerkin finite element approximation of quasilinear elliptic boundary value problems I: The scalar case*, IMA J. Numer. Anal., 25 (2005), pp. 726–749.
- [22] J. KAČUR, *Solution of degenerate convection–diffusion problems by the method of characteristics*, SIAM J. Numer. Anal., 39 (2001), pp. 858–879.
- [23] K. Y. KIM, *A posteriori error estimators for locally conservative methods of nonlinear elliptic problems*, Appl. Numer. Math., 57 (2007), pp. 1065–1080.
- [24] P. LADEVÈZE, *Comparaison de modèles de milieux continus*, Ph.D. thesis, Université Pierre et Marie Curie (Paris 6), 1975.
- [25] R. LUCE AND B. I. WOHLMUTH, *A local a posteriori error estimator based on equilibrated fluxes*, SIAM J. Numer. Anal., 42 (2004), pp. 1394–1414.
- [26] C. MAKRIDAKIS AND R. H. NOCHETTO, *Elliptic reconstruction and a posteriori error estimates for parabolic problems*, SIAM J. Numer. Anal., 41 (2003), pp. 1585–1594.
- [27] P. NEITTAANMÄKI AND S. REPIN, *Reliable methods for computer simulation*, vol. 33 of Studies in Mathematics and its Applications, Elsevier Science B.V., Amsterdam, 2004. Error control and a posteriori estimates.
- [28] R. H. NOCHETTO, A. SCHMIDT, AND C. VERDI, *A posteriori error estimation and adaptivity for degenerate parabolic problems*, Math. Comp., 69 (2000), pp. 1–24.
- [29] M. OHLBERGER, *A posteriori error estimates for vertex centered finite volume approximations of convection–diffusion–reaction equations*, M2AN Math. Model. Numer. Anal., 35 (2001), pp. 355–387.
- [30] F. OTTO,  *$L^1$ -contraction and uniqueness for quasilinear elliptic-parabolic equations*, J. Differential Equations, 131 (1996), pp. 20–38.
- [31] L. E. PAYNE AND H. F. WEINBERGER, *An optimal Poincaré inequality for convex domains*, Arch. Rational Mech. Anal., 5 (1960), pp. 286–292.
- [32] M. PICASSO, *Adaptive finite elements for a linear parabolic problem*, Comput. Methods Appl. Mech. Engrg., 167 (1998), pp. 223–237.
- [33] W. PRAGER AND J. L. SYNGE, *Approximations in elasticity based on the concept of function space*, Quart. Appl. Math., 5 (1947), pp. 241–269.
- [34] F. A. RADU, I. S. POP, AND P. KNABNER, *Error estimates for a mixed finite element discretization of some degenerate parabolic equations*, Numer. Math., 109 (2008), pp. 285–311.

- [35] S. I. REPIN, *Estimates of deviations from exact solutions of initial-boundary value problem for the heat equation*, Atti Accad. Naz. Lincei Cl. Sci. Fis. Mat. Natur. Rend. Lincei (9) Mat. Appl., 13 (2002), pp. 121–133.
- [36] G. SANGALLI, *Robust a-posteriori estimator for advection-diffusion-reaction problems*, Math. Comp., 77 (2008), pp. 41–70.
- [37] R. VERFÜRTH, *A posteriori error estimates for nonlinear problems.  $L^r(0, T; L^p(\Omega))$ -error estimates for finite element discretizations of parabolic equations*, Math. Comp., 67 (1998), pp. 1335–1360.
- [38] ———, *A posteriori error estimates for nonlinear problems:  $L^r(0, T; W^{1,p}(\Omega))$ -error estimates for finite element discretizations of parabolic equations*, Numer. Methods Partial Differential Equations, 14 (1998), pp. 487–518.
- [39] ———, *A posteriori error estimates for finite element discretizations of the heat equation*, Calcolo, 40 (2003), pp. 195–212.
- [40] ———, *A posteriori error estimates for non-linear parabolic equations*. Tech. report, Ruhr-Universität Bochum, 2004.
- [41] ———, *Robust a posteriori error estimates for nonstationary convection-diffusion equations*, SIAM J. Numer. Anal., 43 (2005), pp. 1783–1802.