

# DUALITY BASED ERROR ESTIMATION IN THE PRESENCE OF DISCONTINUITIES

SUSANNE BECKERS, JÖRN BEHRENS, AND WINNIFRIED WOLLNER

ABSTRACT. Goal-oriented mesh adaptation, in particular using the dual-weighted residual (DWR) method, is known in many cases to produce very efficient meshes. For obtaining such meshes the (numerical) solution of an adjoint problem is needed to weight the residuals appropriately with respect to their relevance for the overall error. For hyperbolic problems already the weak primal problem requires in general an additional entropy condition to assert uniqueness of solutions; this difficulty is also reflected when considering adjoints to hyperbolic problems involving discontinuities where again an additional requirement (reversibility) is needed to select appropriate solutions.

Within this article, an approach to the DWR method for hyperbolic problems based on an artificial viscosity approximation is proposed. It is discussed why the proposed method provides a well-posed dual problem, while a direct, formal, application of the dual problem does not. Moreover, we will discuss a further, novel, approach in which the forward problem need not be modified, thus allowing for an unchanged forward solution. The latter procedure introduces an additional residual term in the error estimation, accounting for the inconsistency between primal and dual problem.

Finally, the effectivity of the extended error estimator, assessing the global error by a suitable functional of interest, is tested numerically; and the advantage over a formal estimator approach is demonstrated.

## 1. INTRODUCTION

Within this article, we are concerned with deriving a posteriori error estimates for certain functional values of (entropy) solutions to a first order hyperbolic partial differential equation (PDE)

$$(1) \quad \begin{aligned} \partial_t u + \nabla_x \cdot f(u) &= 0 && \text{on } \Omega \times I, \\ u(x, t) &= u_\Gamma(x, t) && \text{on } \Gamma \times I, \\ u(x, 0) &= u_{\text{ini}}(x) && \text{on } \Omega, \end{aligned}$$

where  $\Omega \subset \mathbb{R}^d$  is a given domain,  $I = (0, T)$  some time-interval,  $f: \mathbb{R} \rightarrow \mathbb{R}^d$  is a smooth function, and  $u_\Gamma, u_{\text{ini}}$  some appropriate given data. It is well known, that the entropy-solutions of (1) will not assume the given boundary values  $u_\Gamma$  on all of  $\partial\Omega$ . In fact, for smooth  $u_\Gamma, u_{\text{ini}}$  the boundary conditions are given by the condition

$$\min_{k \in (u, u_\gamma)} \left( \text{sign}(u - u_\gamma) (f(u) - f(u_\gamma)) \cdot n \right)$$

for all  $(x, t) \in \partial\Omega \times I$ ; where  $(u, u_\gamma)$  denotes the interval spanned by the extreme values  $u$  and  $u_\gamma$  and  $n$  the outward unit normal on  $\Gamma$ , see, [6, Theorem 2]. It shows the delicacy of the calculations that the analogous results for  $L^\infty$ -boundary and initial conditions has only been obtained recently by [19].

Moreover, it is well known, that (1), in general, has no smooth solutions, even if the data are arbitrarily smooth. This is due to the potential creation of so called

---

*Date:* November 20, 2017.

*Key words and phrases.* dual weighted residual, hyperbolic problems, discontinuous Galerkin, artificial viscosity.

shocks at which the entropy solution  $u$  of (1) is no longer continuous. Indeed, for smooth data the solution  $u$  of (1) will be in  $BV(\Omega \times I)$ , see, e.g., [6] and thus the PDE is only satisfied in the distributional sense. For  $L^\infty$ -boundary and initial data the solution will generally be in  $L^\infty(\Omega \times I)$ , only, however such a (weak entropy) solution will still assume its initial and boundary values, see, [19].

A posteriori error estimation for, in particular, elliptic and parabolic PDEs has a long history, see, e.g., [4, 2, 40], and many others. In particular, for elliptic problems it is possible to show that, suitable, adaptive algorithms based on residual error estimates provide meshes of quasi-optimal cardinality for the energy  $H^1$ -norm, see, e.g., the pioneering work [10] and the reviews of the current state of the art [23]. For hyperbolic problems, a posteriori upper bounds of residual type can be obtained for the  $L^1$ -error of the solution, see, e.g., the survey [24].

In contrast to the previously mentioned error estimates for global error norms, goal-oriented error estimation is concerned with estimating a, post-processed, quantity of interest, see, e.g., [2, 7, 11]. In the goal-oriented context, the DWR method, cf., [7, 5], provides an error estimator consisting of weighted residuals of the primal and dual equation. In case of hyperbolic equations, the mathematically sound formulation of a suitable dual solution is not as straightforward. Indeed, adjoint calculus relies on differentiability of the corresponding problem, which is a subtle issue, see for instance [37], [38], and [13, 14], where the problem was tackled by ‘shift differentiability’, suitably modified adjoint based derivative computations, and application of artificial viscosity to the primal and adjoint equations, respectively. The difficulty becomes apparent, as the formal dual problem is a transport equation with potentially discontinuous coefficients for which a useful solution concept requires the notion of reversible solutions, see, e.g., [8].

The difficulties given above for defining the adjoint problem have led to several approaches to goal oriented error estimation for hyperbolic problems. First, the standard approach consists in replacing the hyperbolic problem by an elliptic/parabolic one by adding viscous regularization with a ‘sufficiently small’ viscosity parameter, see, e.g., [18, 35]. Further, [12] argue, without proof, that the functional value is differentiable with respect to the viscosity parameter, and thus the error due to the regularization can be estimated by a Taylor-expansion of the goal-functional. [26] utilize adjoint information to post-process goal-functionals where first the shock positions are recovered to higher accuracy and then the shock is treated as an additional, internal, boundary for the adjoint; this is similar to the derivative representation utilizing ‘shift-differentiability’. Finally, an error representation can be derived directly for the hyperbolic problem assuming that the solution is sufficiently smooth, see, [17]. The thus derived estimator can then formally be applied to hyperbolic problems where the assumptions made in the derivation are no longer satisfied. From the results shown there it appears that this technique works extremely well if the support of the goal-functional and the shock do not interact, while a large overestimation of the error is reported in [17] for functionals interacting with the shock, although subsequent refinement can cure this defect.

In this paper, we are concerned with deriving a hybrid of the two approaches above. Namely, we will consider a DWR error representation that only requires a modification in the dual problem, but is well-posed also for problems with non-smooth solutions thereby combining the advantages of the previously mentioned methods. Further research into the question of how to estimate the error in goal functionals is warranted; since despite the difficulties of providing a mathematically sound framework towards adjoint based goal-oriented in hyperbolic problems, approaches based on the formal ideas are utilized in applications, see, e.g., [29, 27, 39].

The paper is organized as follows: In Section 2, we provide a detailed motivation for the considered setting. To illustrate the substantial difficulty associated with discontinuous solutions to the PDE, we will consider a simple linear test-case in which all relevant difficulties become apparent in Section 3. Further, this section will provide given analytical solutions for the numerical discussion in the end.

In Section 4, we briefly specify a particular discretization, and then, in Section 5, the main part of this paper starts. We propose a new approach towards adjoint based error estimation in which we do not change the forward problem, thereby allowing for an arbitrary discretization of the PDE under consideration. The only change we apply is that we perturb the dual problem, which we account for by an additional consistency term in the derived error estimation.

In the final Section 6, we first provide a numerical experiment in a 1d setting where we can in fact evaluate the error identity, up to quadrature error. This example highlights the difficulty of the formal approach compared to our new hybrid estimator without the additional errors coming from the approximation of the weights as it is always necessary in the DWR method. Finally, a 2d example is shown in which the additional approximation of the weights is considered. The numerical experiments confirm that our new modified adjoint approach is advantageous compared to the formal DWR method when one is interested not only in suitable error indicator but also in an estimate of the error in the goal-functional.

## 2. MOTIVATION

For the use of the DWR method, it is suitable to rewrite the conservation law (1) in weak form. In this context, let  $A(\cdot, \cdot): W \times V \rightarrow \mathbb{R}$  be a semi-linear form, i.e., it is linear in its second component, and let  $F(\cdot)$  be a linear functional on  $V$ .

Suppose that  $u \in W$  is a weak entropy solution of

$$(2) \quad A(u, \psi) = F(\psi) \quad \forall \psi \in V,$$

where  $W$  is a suitable function space and  $V$  is the proper test function space. As discussed in the introduction suitable choices are  $W = BV(\Omega \times I)$  or  $W = L^\infty(\Omega \times I)$  depending on the regularity of the boundary and initial data. In any case, functions in  $W$  in general have only distributional, but not weak, derivatives, and thus test-functions must be sufficiently regular, e.g.,  $V = C^1(\bar{\Omega} \times \bar{I})$  to allow derivatives to be applied to the test-function.

W.l.o.g, we assume that the same form may be used for the discretization, by the finite element method. Then to estimate the discretization error in a goal-functional  $J$ , the standard approach, see, e.g., [7, 5] requires a dual-problem of finding  $z \in V$  solving

$$(3) \quad A'(u; \varphi, z) = J'(u; \varphi), \quad \forall \varphi \in W$$

where  $A'(u, \varphi, z)$  denotes the derivative of  $A(u; z)$  with respect to  $u$  in direction  $\varphi$ . Unfortunately, since (3) is a transport equation, with potentially discontinuous coefficient, the regularity  $z \in V$  can not be expected in general and we will discuss this in more detail in Section 3. The problem of non-fitting solution and test spaces for the primal and dual problem is also mentioned in [7, Remark 2.3].

There are two obvious possibilities to match the solution spaces and test spaces in this setting: For a linear problem, modification of the goal functional, and consequently the data of the dual problem, can increase the regularity of the dual solution  $z$  and thus allowing  $z$  to be used as a test function in (2). Second, and more generally applicable to nonlinear problems, artificial viscosity can be used to prevent shocks and obtain sufficiently smooth adjoint solutions, see for instance [13, 14, 34]. Additionally, we will present subsequently a third alternative, giving a well-posed

dual problem without modification in the primal equation which may be advantageous if changes in the discretization of the forward problem are not suitable.

### 3. DISCONTINUOUS TEST CASE

To illustrate, that the difficulty of unsuitable ansatz- and test-spaces is not only a formality, we consider the most simple setting in which this becomes apparent. Namely, we consider a scalar one dimensional transport equation with discontinuous initial data.

We specify the domain  $\Omega = (-2, 2)$  and the time interval  $I = (0, 1)$  and consider the advection problem of finding a function  $u_0: \Omega \times I \rightarrow \mathbb{R}$  solving

$$\begin{aligned} \partial_t u_0(x, t) + \nabla_x u_0(x, t) &= 0 && \text{on } \Omega \times I, \\ u_0(-2, t) &= 0 && \text{on } I, \\ u_0(x, 0) &= u_{\text{ini}}(x) && \text{on } \Omega, \end{aligned}$$

with the discontinuous initial condition

$$(4) \quad u_0(x, 0) = u_{\text{ini}}(x) = \begin{cases} 1, & -1 < x < 0, \\ 0.5, & x = -1 \text{ or } x = 0, \\ 0, & \text{otherwise.} \end{cases}$$

To avoid confusion, the index 0 to the solution indicates that no viscosity is present in the defining equation. Its weak, entropy, solution for  $t \leq 1$  is given by

$$u_0(x, t) = u_{\text{ini}}(x - t) = \begin{cases} 1, & -1 + t < x < t, \\ 0.5, & x = -1 + t \text{ or } x = t, \\ 0, & \text{otherwise.} \end{cases}$$

which is simply a translation of the initial condition along the characteristic curves.

It is clear, that the above function does not possess weak derivatives. However, for all smooth test-functions  $\psi \in C^1(\overline{\Omega \times I})$  it satisfies the weak problem

$$\begin{aligned} (5) \quad A_0(u_0, \psi) &:= - \int_0^T \int_{\Omega} u_0(x, t) (\partial_t + \partial_x) \psi(x, t) \, dx \, dt \\ &\quad - \int_{\Omega} u_{\text{ini}}(x) \psi(x, 0) \, dx + \int_{\Omega} u_0(x, 1) \psi(x, 1) \, dx \\ &\quad + \int_I u(2, t) \psi(2, t) \, dt \\ &= 0 \end{aligned}$$

note that the boundary integrals are indeed meaningful since weak entropy solutions admit boundary traces in the  $L^1$  sense. Now, consider the goal functional

$$(6) \quad J(u_0) = \int_{\mathbb{R}} u_0(x, 1) z_T(x) \, dx,$$

with the weight  $z_T$  indicating an area of interest

$$z_T(x) := \begin{cases} 1, & 0 \leq x \leq 1, \\ 0, & \text{otherwise.} \end{cases}$$

Using formal duality for (5) gives that the corresponding dual problem consists of finding  $z_0$  such that for all  $\psi \in C^1(\overline{\Omega} \times \overline{I})$  satisfying

$$\begin{aligned}
 A_0(\psi, z_0) &= \int_I \int_{\Omega} (\partial_t + \partial_x)\psi(x, t) z_0(x, t) \, dx \, dt \\
 &+ \int_{\Omega} z(x, 0)\psi(x, 0) \, dx - \int_{\Omega} z_T(x)\psi(x, 1) \, dx \\
 (7) \qquad &+ \int_I z_0(-2, t)\psi(-2, t) \, dt \\
 &= 0.
 \end{aligned}$$

Exchanging the temporal variable by  $t \mapsto 1-t$  shows that this is exactly the forward problem (5) with the initial condition  $z_T(x)$ .

Unfortunately, now neither  $u_0$  nor  $z_0$  have a weak derivative, but in fact  $\partial_x$  and  $\partial_t$  are only defined in the sense of distributions. Consequently neither the weak form (5) tested with  $z_0$  nor the dual (7) tested with  $u_0$  is well defined.

*Remark 1.* Clearly, taking  $\partial_t u_0 + \partial_x u_0 = 0$  the volume integral could be set to zero, however, since our discretization will provide separate approximations to  $\partial_t u_0$  and  $\partial_x u_0$  it is desirable to have that each individual product, i.e.,  $\partial_t u_0 z_0$  is well defined!

#### 4. DISCRETIZATION SCHEMES

Hyperbolic problems can develop or maintain discontinuities in the solutions, as seen in the advection example in Section 3 and is generally well known. One approach for an accurate and efficient method to solve advection dominated problems numerically are the discontinuous Galerkin (dG) methods. These methods combined with slope limiters are able to capture the physically relevant discontinuities without producing spurious oscillations, [9].

Some of the first to apply the dG method were W. Reed and T. Hill, [30], in 1973. DG methods are generalizations of finite volume methods but possess also properties of finite element methods, as for instance the simple handling of complex geometries and of boundary conditions. The advantage of dG lies in the discontinuities at the element boundaries and the thereby resulting simple routines for parallelization and adaptivity. These advantages, however, have to be bought by the price of a higher number of degrees of freedom than for the continuous finite element schemes.

In this section, we will discuss the spatial discretization of the primal problem (1) and its viscous regularization

$$\begin{aligned}
 \partial_t u + \nabla_x \cdot f(u) &= \varepsilon \Delta u && \text{on } \Omega \times I, \\
 (8) \qquad u(x, t) &= u_{\Gamma}(x, t) && \text{on } \Gamma \times I, \\
 u(x, 0) &= u_{\text{ini}}(x) && \text{on } \Omega,
 \end{aligned}$$

and the corresponding dual by a discontinuous Galerkin method. For the temporal discretization, we confine ourselves to a method-of-lines setting in which only one spatial mesh is considered at each time-point.

For the spatial discretization the domain  $\Omega$  is decomposed into a set  $\mathcal{E}$  of non-overlapping open elements  $E$  (intervals if  $d = 1$ , triangles or quadrilaterals if  $d = 2$ , and so on) of diameter  $h_E$  such that

$$\overline{\Omega} = \bigcup_{E \in \mathcal{E}} \overline{E}.$$

For each element  $E \in \mathcal{E}$  the flux in the element is defined as

$$F_\varepsilon(v)(x, t) := f(v(x, t)) - \varepsilon \partial_x v(x, t), \quad (x, t) \in E \times (0, 1).$$

For the boundaries, we select a suitable, consistent and conservative, flux-function  $\mathcal{H}: \mathbb{R} \times \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}$ , cf., [9, 16], i.e., it holds

$$\mathcal{H}(v, w, n) = -\mathcal{H}(w, v, -n), \quad \mathcal{H}(v, v, n) = n \cdot F(v).$$

Then following the SIPG method of [41, 32] for the diffusion part, we obtain the weak form for the spatial discretization

$$a_\varepsilon(u, v) = \sum_{E \in \mathcal{E}} \left( - \int_E F_\varepsilon(u) \nabla v \, dx + \int_{\partial E} \mathcal{H}(u^+, u^-, n) v \, ds \right. \\ \left. + \int_{\partial E} \frac{c_\varepsilon}{h_E} [u][v] - \varepsilon \{n \cdot \nabla u\}[v] - \varepsilon \{n \cdot \nabla v\}[u] \, ds \right)$$

for any functions  $u, v$  being piecewise smooth with respect to the subdivision  $\mathcal{E}$ . Here the superscripts  $+$  and  $-$  denote the interior and exterior trace of the piecewise smooth functions,  $n$  the outer unit normal vector, and  $\{\}$  and  $[\ ]$  denote the average and jump across the given boundary; together with an appropriate definition of the exterior trace when  $\partial\Omega \cap \partial E \neq \emptyset$ . Following [41], we set  $c_\varepsilon = p^2$  where  $p$  is the polynomial degree of the dG-ansatz when  $\varepsilon > 0$  and  $c_\varepsilon = 0$  when  $\varepsilon = 0$ . Note for  $\varepsilon = 0$  all terms involving  $\varepsilon$  vanish, i.e.,

$$a_0(u, v) = \sum_{E \in \mathcal{E}} \left( - \int_E f(u) \nabla v \, dx + \int_{\partial E} \mathcal{H}(u^+, u^-, n) v \, ds \right).$$

Choosing the finite dimensional space  $V^h = \{\psi \in L^2(\Omega) : \forall E \ \psi|_E \in P_p(E)\}$  where  $P_p(E)$  is the space of polynomials of degree  $p$  on element  $E$ .

With these definitions, a semi-discrete version of (1) is to find a function  $u_\varepsilon^h: \Omega \times I \rightarrow \mathbb{R}$  with  $u_\varepsilon^h(t) \in V^h$  solving

$$- \int_I \int_\Omega u_\varepsilon^h \partial_t \varphi \, dx \, dt + \int_I a_\varepsilon(u_\varepsilon^h, \varphi) \, dt + \int_\Omega u_{\text{ini}}(x) \varphi(x, 0) + u_\varepsilon^h(x, T) \varphi(x, T) \, dx = 0$$

for all  $\varphi \in C^1(\bar{I}; V_h)$ . This corresponds to the weak form (5) which treats the discontinuous test case of Section 3.

For the temporal-discretization, we utilize the simplest discretization by a  $\theta$ -scheme using its  $dG(0)$  interpretation. The setting can be extended to varying spatial meshes in time with some additional technicalities, see, e.g., [33, 15]. We consider  $0 = t_0 < t_1 < \dots < t_M = T$ , define  $k_i = t_i - t_{i-1}$  for  $i = 1, \dots, M$  and consider the space-time discrete space  $X = \{v: \Omega \times I \mid v|_{(t_{i-1}, t_i]} \in V^h, \ i = 1, \dots, M\}$  of piecewise constants in time. For abbreviation, we write  $v^i := v|_{(t_{i-1}, t_i]}$  for a function  $v \in X$ . Then the dG(0) discretization in time reads: find  $u_\varepsilon^{kh} \in X$  solving

$$(9) \quad \int_{t_{i-1}}^{t_i} \int_\Omega \partial_t u_\varepsilon^{kh} \varphi \, dx \, dt + \int_{t_{i-1}}^{t_i} a_\varepsilon(u_\varepsilon^{kh}, \varphi) \, dt + \int_\Omega (u_\varepsilon^i - u_\varepsilon^{i-1}) \varphi(x) \, dx = 0$$

for all  $i = 1, \dots, M$  and  $\varphi \in V^h$  where  $u_\varepsilon^0 = u_{\text{ini}}$ . Noting that  $\partial_t u_\varepsilon^{kh} \equiv 0$  and that  $a_\varepsilon(u_\varepsilon^{kh}, \varphi) = a_\varepsilon(u_\varepsilon^i, \varphi)$  is constant on  $(t_{i-1}, t_i)$  gives the implicit Euler-scheme

$$\int_\Omega u_\varepsilon^i \varphi(x) \, dx + k_i a_\varepsilon(u_\varepsilon^i, \varphi) = \int_\Omega u_\varepsilon^{i-1} \varphi(x) \, dx \quad \forall \varphi \in V^h.$$

In particular, the above shows, that any weak solution to (1) solves (9), with  $\varepsilon = 0$ , as well, i.e.,

$$(10) \quad \int_{t_{i-1}}^{t_i} \int_\Omega \partial_t u_0 \varphi \, dx \, dt + \int_{t_{i-1}}^{t_i} a_0(u_0, \varphi) \, dt + \int_\Omega (u_0(t_i^+) - u_0(t_i^-)) \varphi(x) \, dx = 0$$

for all  $\varphi \in V^h \cup C^1(\bar{\Omega})$ . We summarize the time-steps of equation (9) by defining

$$\begin{aligned} A_\varepsilon(u, \varphi) &:= \sum_{i=1}^M \int_{t_{i-1}}^{t_i} \int_{\Omega} \partial_t u_\varepsilon^{kh} \varphi \, dx + a_\varepsilon(u_\varepsilon^{kh}, \varphi) \, dt \\ &\quad + \sum_{i=2}^M \int_{\Omega} (u_\varepsilon^i - u_\varepsilon^{i-1}) \varphi^i(x) \, dx + \int_{\Omega} u_\varepsilon^1 \varphi^1(x) \, dx \end{aligned}$$

for  $u, \varphi \in X$ . Then  $u_\varepsilon^{kh} \in X$  solves (9) if and only if

$$(11) \quad A_\varepsilon(u_\varepsilon^{kh}, \varphi) = \int_{\Omega} u_{\text{ini}} \varphi^1(x) \, dx \quad \forall \varphi \in X.$$

Similar procedures can give the explicit Euler-scheme, and many other known time-stepping procedures. Note, however, that not in all cases is  $A_0(u_0, \varphi) = 0$  satisfied for the solution of (1) since quadrature errors may occur in the derivation.

A simple example is the explicit Euler-scheme combined with a left box rule for integration: The forward Euler method gives a piecewise linear solution,

$$u_0^{kh}|_{[t_{i-1}, t_i]} = \left(1 - \frac{t - t_{i-1}}{t_i - t_{i-1}}\right) u_0^{i-1} + \frac{t - t_{i-1}}{t_i - t_{i-1}} u_0^i,$$

for  $i = 1, \dots, M$ , which is continuous on  $I$ . Integration as

$$\sum_{i=1}^M \int_{t_{i-1}}^{t_i} \int_{\Omega} \partial_t u_0^{kh} \varphi \, dx \, dt$$

can be done exactly by the box rule, since the piecewise derivative and  $\varphi$  are piecewise constant in time. But

$$\sum_{i=1}^M \int_{t_{i-1}}^{t_i} a_0(u_0^{kh}, \varphi) \, dt \approx \sum_{i=1}^M k_i a_0\left((u_0^{kh})^{i-1}, \varphi\right) \, dt, \quad k_i := t_i - t_{i-1},$$

causes integration errors, because  $u_0^{kh}$  is still piecewise linear in time, but evaluation of  $a_0(u_0^{kh}, \varphi)$  will in general not result in constant functions in time. The error of the box rule is of order  $O(k)$ . Thus,  $A_0(u_0, \varphi) = 0$  does not have to be satisfied exactly.

Now, we define the linearized operator in some  $u \in X$  by

$$\begin{aligned} A'_\varepsilon(u; \psi, \varphi) &= \sum_{i=1}^M \int_{t_{i-1}}^{t_i} \int_{\Omega} \partial_t \psi \varphi \, dx + a'_\varepsilon(u; \psi, \varphi) \, dt \\ &\quad + \sum_{i=2}^M \int_{\Omega} (\psi^i - \psi^{i-1}) \varphi^i(x) \, dx + \int_{\Omega} \psi^1 \varphi^1(x) \, dx \end{aligned}$$

for all  $\psi, \varphi \in X$ , where

$$\begin{aligned} a'_\varepsilon(u; \psi, \varphi) &= \sum_{E \in \mathcal{E}} \left( - \int_E F'_\varepsilon(u; \psi) \nabla \varphi \, dx + \int_{\partial E} \mathcal{H}'(u^+, u^-, n) \begin{pmatrix} \psi^+ \\ \psi^- \end{pmatrix} \varphi \, ds \right. \\ &\quad \left. + \int_{\partial E} \frac{c_\varepsilon}{h_E} [\psi][\varphi] - \varepsilon \{n \cdot \nabla \psi\}[\varphi] - \varepsilon \{n \cdot \nabla \varphi\}[\psi] \, ds \right) \end{aligned}$$

where  $\mathcal{H}'(u^+, u^-, n)$  is the corresponding derivative of the flux with respect to  $u^+, u^-$ .

When concerned with the evaluation of the discretization error in a linear functional  $J: BV(\Omega \times I) \rightarrow \mathbb{R}$  the corresponding discrete dual problem reads as: find  $z \in X$  solving

$$(12) \quad A'_\varepsilon(u; \varphi, z) = J(\varphi) \quad \forall \varphi \in X$$

To assert, that this is indeed a suitable discretization of the corresponding formal continuous dual problem suitable numerical fluxes  $\mathcal{H}$  and functionals  $J$  are required, see, e.g., [16]. In what follows, we assume that the numerical fluxes are chosen suitably asserting that the continuous dual solution  $z_\varepsilon$  satisfies

$$(13) \quad A'_\varepsilon(z_\varepsilon; u_\varepsilon - u_\varepsilon^{kh}, z_\varepsilon) = J(u_\varepsilon - u_\varepsilon^{kh})$$

## 5. ERROR ESTIMATOR WITH CORRECTION TERM

We will first provide a generic, and well known error estimate for the discretization error in the functional  $J$ , see, e.g., [7]. As usual, we define the Lagrangian  $\mathcal{L}_\varepsilon(u, z) = J(u) - A_\varepsilon(u, z) + \int_\Omega u_{\text{ini}}(x)z(x, 0) dx$  to simplify the exposition.

**Theorem 1.** *Assume that the consistency relation  $A_\varepsilon(u_\varepsilon, z_\varepsilon) = \int_\Omega u_{\text{ini}}(x)z(x, 0) dx$  holds for the solution  $u_\varepsilon$  of (8) and its adjoint  $z_\varepsilon$  and that the adjoint is consistent in the sense that (13) holds. Define  $\mathbf{x} = (u_\varepsilon, z_\varepsilon)$  and  $\mathbf{x}^{kh} = (u_\varepsilon^{kh}, z_\varepsilon^{kh})$  as solution of (11) and (12) and assume that  $\mathcal{L}_\varepsilon$  is three times continuously differentiable on the line  $\text{conv}(\mathbf{x}, \mathbf{x}^{kh})$ .*

*Then for its discrete approximation  $u_\varepsilon^{kh}$  given by (11) it holds the error representation*

$$J(u_\varepsilon) - J(u_\varepsilon^{kh}) = \frac{1}{2} \mathcal{L}'_\varepsilon(\mathbf{x}^{kh})(\mathbf{x} - \tilde{\mathbf{x}}^{kh}) + \mathcal{R}$$

*with an arbitrary  $\tilde{\mathbf{x}}^{kh} \in X^2$  and a remainder  $\mathcal{R}$ ; cubic in the error  $e = \mathbf{x} - \mathbf{x}^{kh}$ .*

*Proof.* The proof is the standard calculation for the DWR-method, see, e.g., [7]. It is only detailed to clarify where the consistency assumption enters.

By our consistency assumption and (11) it holds

$$J(u_\varepsilon) - J(u_\varepsilon^{kh}) = \mathcal{L}_\varepsilon(u_\varepsilon, z_\varepsilon) - \mathcal{L}_\varepsilon(u_\varepsilon^{kh}, z_\varepsilon^{kh}).$$

Then by the assumed differentiability, and the remainder term of the trapezoidal rule it is

$$\begin{aligned} J(u_\varepsilon) - J(u_\varepsilon^{kh}) &= \frac{1}{2} \int_0^1 \mathcal{L}'_\varepsilon(\mathbf{x}^{kh} + s(e))(e) ds \\ &= \frac{1}{2} \left( \mathcal{L}'_\varepsilon(\mathbf{x})(e) + \mathcal{L}'_\varepsilon(\mathbf{x}^{kh})(e) \right) \\ &\quad + \frac{1}{2} \int_0^1 \mathcal{L}'''_\varepsilon(\mathbf{x}^{kh} + se)(e, e, e) s(s-1) ds. \end{aligned}$$

Utilizing Galerkin-orthogonality for the primal and adjoint-consistency (13) we get  $\mathcal{L}'_\varepsilon(\mathbf{x})(e) = 0$  and the definition of  $\mathbf{x}_\varepsilon^{kh}$  gives

$$\mathcal{L}'_\varepsilon(\mathbf{x}^{kh})(e) = \mathcal{L}'_\varepsilon(\mathbf{x}^{kh})(\mathbf{x} - \tilde{\mathbf{x}}^{kh})$$

showing the assertion.  $\square$

*Remark 2.* Before we continue a few remarks are in order.

- (1) If  $\varepsilon > 0$ , the problem (8) is a semi-linear parabolic problem, for which the assumptions of the theorem regarding consistency and differentiability can be verified in many cases.
- (2) However, if  $\varepsilon = 0$  several of the assumptions are questionable, e.g., we have seen in Section 3 that the definition of  $A_\varepsilon(u_\varepsilon, z_\varepsilon)$  is not straightforward. Furthermore, differentiability of the Lagrangian is a delicate issue as it can not rely on a simple chain rule, since differences of discontinuous



functions are generally not even directionally differentiable in  $L^1$ . To see this, compare [36, Example 3.1.1.], let

$$v_\varepsilon = \begin{cases} 1 & x < \varepsilon, \\ 0 & \text{otherwise.} \end{cases}$$

Then clearly  $v_\varepsilon \rightarrow v_0$  pointwise a.e. as  $\varepsilon \rightarrow 0$  but

$$\int_0^1 \frac{v_\varepsilon - v_0}{\varepsilon} ds = \int_0^\varepsilon \frac{1}{\varepsilon} ds = 1$$

and consequently the limit of the difference quotient can not be a function. This makes calculations of the type

$$\frac{d}{ds} \mathcal{L}_\varepsilon(\mathbf{x}^{kh} + s(e)) = \mathcal{L}'_\varepsilon(\mathbf{x}^{kh} + s(e))e,$$

needed in the proof, quite intricate. Notice that in the example above the mapping  $\varepsilon \mapsto \int_0^1 v_\varepsilon(s) ds = \varepsilon$  is clearly differentiable with the expected derivative but this can not be obtained by application of the chain rule to the mappings  $\mathbb{R} \rightarrow L^1$  with  $\varepsilon \mapsto v_\varepsilon$  and  $L^1 \rightarrow \mathbb{R}$  given by  $v \mapsto \int_0^1 v(s) ds$  as it is commonly argued.

Consequently, a standard approach to the goal-oriented error estimation consists of taking  $\varepsilon > 0$  (small enough) and then estimate the discretization error for the respective given  $\varepsilon$ . Unfortunately, this introduces an additional viscosity to the problem (1) for the sole purpose of estimating the error. This may be undesirable if numerics are tuned to add little viscosity in order to provide a better resolution of the discontinuities.

To this end, we now propose a new approach that avoids adding additional viscosity to the primal problem. We now introduce the residuals of the inviscid problem and its dual by

$$(14) \quad \rho(u_0, z_\varepsilon - z_\varepsilon^{hk}) := \int_\Omega u_{\text{ini}}(z_\varepsilon - z_\varepsilon^{hk}) dx - A_0(u_0, z_\varepsilon - z_\varepsilon^{hk}),$$

$$(15) \quad \rho^*(w; z_\varepsilon, u_0 - u_0^{hk}) := J(u_0 - u_0^{hk}) - A'_0(w; u_0 - u_0^{hk}, z_\varepsilon).$$

**Theorem 2.** *Let  $\Omega \subset \mathbb{R}$  and  $J(u) = \int_\Omega \omega(x)u(x, T) dx$  with a weight  $\omega \in L^\infty(\Omega)$ . Further, let  $u_0$  be the entropy solution to (1) and  $u_0^{kh}$  its discrete approximation by (11) with  $\varepsilon = 0$  and let  $z_\varepsilon^{kh}$  be the adjoint defined by (12) and assume that the numerical fluxes are such that the continuous adjoint satisfies (13). Define  $\mathbf{x} = (u_0, z_\varepsilon)$  and  $\mathbf{x}^{kh} = (u_0^{kh}, z_\varepsilon^{kh})$  and assume that  $\mathcal{L}_0$  is three times continuously differentiable on the line  $\text{conv}(\mathbf{x}, \mathbf{x}^{kh})$ .*

*Then the error representation*

$$J(u_0) - J(u_0^{kh}) = \frac{1}{2} \left( \rho^*(u_0^{kh}; z_\varepsilon^{kh}, u_0^{kh} - \tilde{u}_0^{hk}) + \rho(u_0^{kh}, z_\varepsilon - \tilde{z}_\varepsilon^{hk}) + \rho^*(u_0; z_\varepsilon, u_0 - u_0^{hk}) \right) + \mathcal{R}$$

*holds with arbitrary  $\tilde{u}_0^{kh}, \tilde{z}_\varepsilon^{kh} \in X$  and a remainder  $\mathcal{R}$ ; cubic in the error  $e = \mathbf{x} - \mathbf{x}^{kh}$ .*

*Proof.* We know, that the solution  $u_0 \in BV(\Omega \times I)$  and thus  $f'(u_0) \in L^\infty(\Omega)$ . Thus the solution  $z_\varepsilon$  of the adjoint problem

$$\begin{aligned} -\partial_t z_\varepsilon - f'(u_0) \nabla_x z_\varepsilon &= \varepsilon \Delta_\varepsilon && \text{on } \Omega \times I, \\ z(x, t) &= 0 && \text{on } \Gamma \times I, \\ z(x, T) &= \omega(x) && \text{on } \Omega, \end{aligned}$$

satisfies the regularity  $z_\varepsilon \in C(\bar{I}, L^2(\Omega)) \cap L^2(I, H^2(\Omega) \cap H_0^1(\Omega))$ , see, e.g., [22]. By standard embedding theorems, see, e.g., [1], it is therefore  $z_\varepsilon \in L^2(I, C^1(\bar{\Omega}))$ .

As a consequence,  $z_\varepsilon$  is a suitable test function in  $A_0$ , see (10). Analogous to Theorem 1 we get

$$\begin{aligned} J(u_0) - J(u_0^{kh}) &= \mathcal{L}_0(u_0, z_\varepsilon) - L(u_0^{kh}, z_\varepsilon^{kh}) \\ &= \frac{1}{2} \int_0^1 \mathcal{L}'_0(\mathbf{x}^{kh} + s(e))(e) \, ds \\ &= \frac{1}{2} \left( \mathcal{L}'_0(\mathbf{x})(e) + \mathcal{L}'_0(\mathbf{x}^{kh})(e) \right) \\ &\quad + \frac{1}{2} \int_0^1 \mathcal{L}'''_0(\mathbf{x}^{kh} + se)(e, e, e) s(s-1) \, ds. \end{aligned}$$

In contrast to Theorem (1)

$$(16) \quad \mathcal{L}'_0(\mathbf{x})(e) = \frac{1}{2} \left[ \rho(u_0, z_\varepsilon - z_\varepsilon^{kh}) + \rho^*(u_0; z_\varepsilon, u_0 - u_0^{kh}) \right]$$

$$(17) \quad = \frac{1}{2} \rho^*(u_0; z_\varepsilon, u_0 - u_0^{kh}) \neq 0$$

in general, since  $z_\varepsilon$  does not solve the adjoint for  $\varepsilon = 0$ . With the remaining arguments as in the proof of Theorem 1 we obtain the assertion.  $\square$

*Remark 3.* The simplification  $J(u) = \int_\Omega \omega(x)u(x, T) \, dx$  and  $\Omega \subset \mathbb{R}$  is only needed to assert that  $z_\varepsilon$  is regular enough. This can be done in many other situations as well but for simplicity of the exposition we avoid the additional notation overhead.

In comparison to the error given in [7] the first two residuals now contain  $z_\varepsilon$  instead of  $z_0$ , as was expected, but an additional dual residual,  $\rho^*(z_\varepsilon, u_0 - u_0^{kh})$  has to be taken into account.

In fact, to obtain a computable error estimator, the errors  $u_0 - \tilde{u}_0^{kh}$  and  $z_\varepsilon - \tilde{z}_\varepsilon^{kh}$  need to be approximated. For this several heuristic techniques are known, see, e.g., [2, 7, 40]. For brevity, we will denote the resulting weights by  $w_z \approx z_\varepsilon - \tilde{z}_\varepsilon^{kh}$  and  $w_u \approx u_0 - \tilde{u}_0^{kh}$ . Further, if element-wise indicators are used for mesh-refinement, integration by parts needs to be performed such that the derivatives in the residual are all on the discrete solution, and not on the weight. This is required to assert the correct localization of the indicators. An new alternative approach, [31], can directly work on the weak form of the residual by testing with a suitable partition of unity in the weak form.

The resulting time averaged element indicator is given by

$$\begin{aligned} 2\eta_E &:= \rho_E(u_0^{kh}, w_z) \\ &= \sum_{i=1}^M \left[ \int_{t_{i-1}}^{t_i} \int_E \partial_t u_0^{kh} w_z + \nabla \cdot f(u_0^{kh}) w_z \, dx \right. \\ &\quad \left. + \int_{\partial E} \{n \cdot f(u_0^{kh})\} [w_z] + \mathcal{H}((u_0^{kh})^+, (u_0^{kh})^-, n) w_z \, ds \right] dt \\ &\quad + \sum_{i=2}^M \left[ \int_E (u_0^{kh}(t_{i-1}^+) - u_0^{kh}(t_{i-1}^-)) w_z(t_{i-1}^+) \, dx \right] \\ &\quad + \int_E u_0^{kh}(t_0^+) w_z(t_0^+) \, dx. \end{aligned}$$

Analogously, the dual time averaged element indicators

$$2\eta_E^* = \rho_E^*(u_0^{kh}; z_\varepsilon^{kh}, w_u)$$

and the adjoint-consistency indicator

$$2\hat{\eta}_E^* = \rho_E^*(u_0; z_\varepsilon, w_u)$$

can be defined. We notice, that  $\hat{\eta}_E^*$  can not be evaluated exactly, hence unless  $u_0$  and  $z_\varepsilon$  are known need to be approximated. For the purpose of this article, we simply take  $\hat{\eta}_E^* = \eta_E^*$  as this approximation, although more elaborate techniques might be advantageous as we will see.

We thus obtain the element indicator

$$(18) \quad \eta_E := \eta_E + \eta_E^* + \hat{\eta}_E^*$$

and the approximation

$$(19) \quad \eta := \sum_{E \in \mathcal{E}} \eta_E \approx J(u_0) - J(u_0^{hk})$$

to the global error.

In many other cases, the error identity is not the simple sum of weighted discrete residuals as in Theorem 1 but contains additional residuals on the continuous level as in Theorem 2. In many such cases, the size of the additional residual can be used to estimate the error due to the parameter inducing the inconsistency, see, e.g., [42, 21]. In other cases such residuals can be used to steer the accuracy of non-linear solvers, see, e.g., [20, 28]. In particular, the appearance of such a term may be utilized in the future to couple the size of the additional viscosity for the dual problem to the size of the functional error.

Concluding, in this section it was shown that a modification in the dual equation introduces an additional dual residual. Given the goal functional of the above mentioned advection problem the error in the goal functional can be computed easily. This representation of the error in the goal functional which is due to the introduction of diffusion in the dual equation is going to be evaluated numerically in the next section.

## 6. NUMERICAL EXPERIMENTS

In this section, we study the dependence of the absolute value of the additional residual on the spatial grid size numerically on a 1D and a 2D test case. Then, we investigate the behavior of the local error estimators with and without the additional dual residual and, in the end, we find the global error estimator including the additional residual to gain a better effectivity index as the global estimator without the artificial viscosity. We chose our first test case, such that analytical solutions  $u_0$  and  $z_0$  are known. Therefore we can evaluate  $w_z = z_\varepsilon - \tilde{z}_\varepsilon^{kh}$  and  $w_u = u_0 - \tilde{u}_u^{kh}$  at quadrature points. Since the analytic solutions are not known in general, we use interpolation techniques in our second test case to approximate the weights. In both examples at hand, we used the explicit Euler method for time discretization and Lagrange polynomials of degree two for spatial discretization by the above introduced dG method without any limiter. With the dG method the box-shaped initial condition for the primal,  $u_{\text{ini}}$ , and the dual,  $z_{\text{T}}$ , can be initialized without any initialization error. But the numerical advection causes some over and under shootings in front of steep gradients since we did not apply any limiter. In contrast, the dual problem with diffusion advects the initial condition  $z_{\text{T}}$  to the left hand side and smooths the steep gradients. In this case, the numerical solution is close to the analytical one. For sufficiently smooth solutions the SIPG method provides convergence of  $L^2$ -errors of the order  $p + 1$ , where  $p$  is the order of the polynomial, compare [32]. For discontinuous initial conditions, the order of convergence is lower.

**6.1. 1D Example.** In the following, the spatial discretization on  $\Omega = (-2, 2)$  in dG manner uses basis and test function polynomials of order 2 and the discrete solutions are evaluated such that a numerical quadrature, using a composite trapezoidal rule on each element, can be performed. The value of the goal functional for the discrete solution,  $J(u_0^{hk})$ , is also determined exactly by the trapezoidal rule. The goal value of the analytic solution,  $J(u_0)$ , is one.

In this setting, the global value of the additional residual was computed for different spatial resolutions and a fixed time step size of  $k = 0.0001$ .

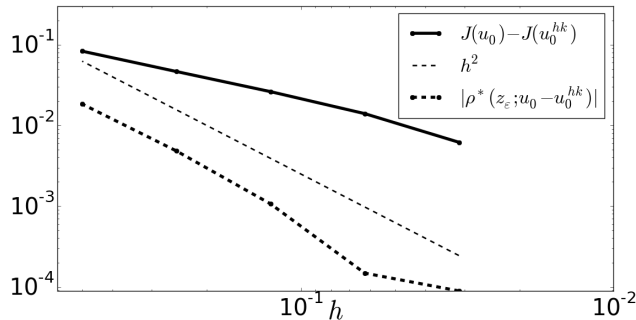


FIGURE 1. Absolute value of the additional residual,  $k = 0.0001$ ,  $\varepsilon = 0.1$ .

Fig. 1 shows the absolute global value of the additional residual for  $\varepsilon = 0.1$ . This extra term converges with second order to zero and is thus faster than the actual error in the goal functional, implying that the additional residual is negligible on sufficiently fine meshes.

However, the difference to the classical formulation is not only the additional residual, but also the replacement of the discontinuous dual function by the solution of the dual advection diffusion equation.

All three residuals together, element-wise evaluated, give the local error estimators, see (18). For a uniform grid the local error estimator  $\eta_{E,uni}$  indicates the area of influence for the goal functional, (6).

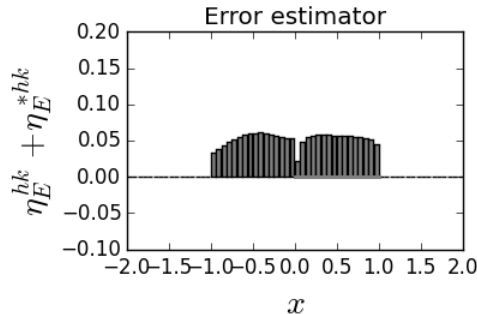


FIGURE 2. Absolute value of local error estimators on a uniform grid with  $h = 0.0625$  and artificial viscosity  $\varepsilon = 0.05$

The area of interest, on which the goal functional is evaluated, is the interval  $[0, 1]$ . Thus, for the discontinuous test case presented in this paper, each element over which the box shaped function moves, has theoretically an equally high local error estimator, while the regions outside are of minor influence to the value of the

goal functional. This is reflected in the numerical results, as, e.g., in Fig. 2, despite the diffusion in the dual the estimator does not smear.

Dörfler marking, [10], would suggest to refine the elements in the middle of Fig. 2, namely in the interval  $I_{\text{ref}} = [-1, 1]$ , such that the sum of the absolute values of the estimators in the set which is going to be refined,  $\mathcal{E}_{\text{ref}}$ , is larger than a specific percentage of the sum of the absolute values of the estimators in the whole set  $\mathcal{E}$ ,

$$\sum_{E \in \mathcal{E}_{\text{ref}}} |\eta_E^{hk} + \eta_E^{*hk}| \geq (1 - \Theta) \sum_{E \in \mathcal{E}} |\eta_E^{hk} + \eta_E^{*hk}|,$$

for some  $\Theta \in (0, 1)$ . In the test case at hand, refinement in  $I_{\text{ref}}$  is achieved with  $1 - \Theta \approx 1 - 10^{-6}$ , showing that most of the estimated error is in  $I_{\text{ref}}$ .

The summation over each element of the signed local spatial error estimators on the uniform grid brings the global estimator  $\eta_{\text{uni}}^{hk}$ , while the sum of the estimators over the locally refined, grid brings  $\eta_{\text{ref}}^{hk}$ .

TABLE 1. Dependence of the global error estimators and the error in the goal functional on the grid size. Uniform grid size  $h$  is marked in bold.  $\varepsilon = 0.1$

$h$	$\eta_{\text{uni}}^{hk}$	$\eta_{\text{ref}}^{hk}$	$ J(u_0) - J(u_{0,\text{uni}}^{hk}) $	$ J(u_0) - J(u_{0,\text{ref}}^{hk}) $
<b>0.5</b> /0.25	0.0674	0.0552	0.0832	0.0466
<b>0.25</b> /0.125	0.0437	0.0324	0.0466	0.0262
<b>0.125</b> /0.0625	0.0258	0.0174	0.0262	0.0140
<b>0.0625</b> /0.03125	0.0140	0.0080	0.0140	0.0062

Table 1 shows that the global error estimator on a uniform grid is greater than the estimator on a mesh which is locally refined once by bisection according to the error indicators. The bold  $h$  indicates the uniform grid size, the normal style  $h$  is the size of the refined elements. Also the error in the goal functional evaluated on a locally refined grid is smaller than on the uniform grid. On each element, the quadrature rule is the same, such that the approximation of the integral is better in the refined elements. But the numerically evaluated error estimator does not satisfy the error identity. This could be caused by several reasons, for instance by a quadrature error or by a non-adjoint consistent implementation [16]. However, as the difference between estimator and error decrease with decreasing element size  $h$ , we do not investigate this issue further.

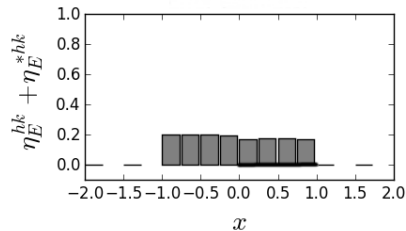


FIGURE 3. Absolute value of local error estimators on a locally refined grid with  $h = 0.5, 0.25$ , and  $\varepsilon = 0$

If the dual equation is not modified and the computations are done nevertheless by evaluating only

$$(20) \quad \eta_{E0} := \rho_E^{hk} (u_0^{hk}, z_0 - z_0^{hk}) + \rho_E^{*hk} (z_0^{hk}, u_0 - u_0^{hk}),$$

compare equation (18), the local error estimators are even more evenly distributed on the area which is expected to be refined, as shown in Fig. 3. For the computation it was naively assumed, that  $\partial_x z_0 = \partial_t z_0 = 0$  in  $[-2, 2]$ , since this is true almost everywhere – and in particular in the chosen quadrature points.

Concluding, the modification of the dual equation does not harm the local error indication, and even the approach without modification – ignoring the unboundedness in the analytic case – results in reasonable local error indication. So far, there seems to be no advantage in the modification, but this is different for the global error estimation:

The quality of the global error estimators is measured by the effectivity index, see, e.g., [3, 4], which is the ratio of the estimator to the true error. Here it is

$$\text{eff} = \frac{J(u_0) - J(u_0^{hk})}{\eta^{hk}}.$$

Fig. 4 shows the behavior of the effectivity index with respect to the spatial grid size. The index for the global error estimator without viscosity, e.g., equation (20), is increasing at first. If it ever converges to one, it is much later as in case of the modified error estimator.

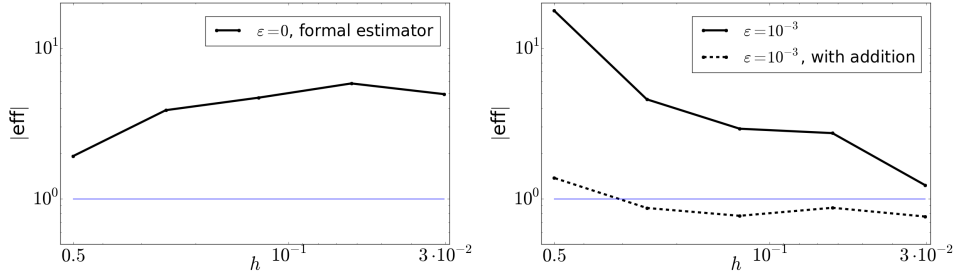


FIGURE 4. Effectivity of the global estimator without artificial viscosity in the dual equation (left) and with and without the additional residual,  $\rho^*(z_\varepsilon, u_0 - u_0^{hk})$  with viscosity  $\varepsilon = 0.001$  in the dual equation (right).

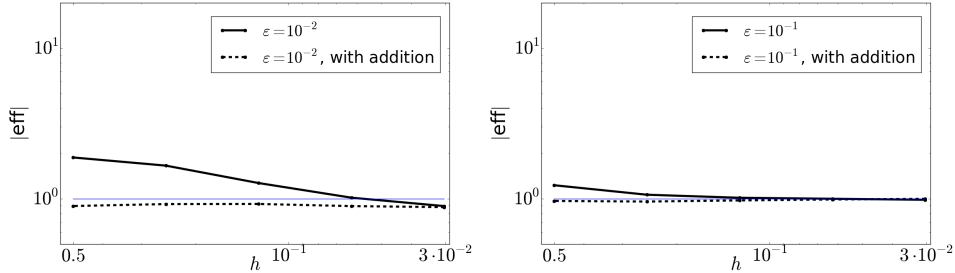


FIGURE 5. Effectivity of the global estimator with and without the additional residual,  $\rho^*(z_\varepsilon, u_0 - u_0^{hk})$ , for  $\varepsilon = 0.01$  (left) and  $\varepsilon = 0.1$  (right).

The right hand side of Fig. 4 shows that the error estimator including the additional residual gains a better effectivity on coarse grids than the estimator without the additional term, e.g., equation (19) with and without the last residual. Fig. 5 depicts this relation also for different values of  $\varepsilon$ . For any tested  $\varepsilon \in [0.0001, 0.1]$  the

effectivity of the estimator including the addition was closer to one, but obviously depending on the diffusion coefficient. Thus, the relation of the error and the error estimator to the diffusion parameter  $\varepsilon$  is of interest. Table 2 shows the different global spatial error estimators for a decreasing diffusion coefficient. For stability reasons, the time step size was chosen to be  $k = 0.0001$  and the spatial grid size was fixed at  $h = 0.0625$ . Since neither the primal problem nor the goal functional are modified, the error in the goal functional is constant for a fixed grid size.

TABLE 2. Dependence of the global spatial error estimators and on the dual diffusion coefficient  $\varepsilon$ , with  $J(u_0) - J(u_0^{hk}) = 0.0140$

$\varepsilon$	$ \eta_{\text{uni}}^{hk} $	eff
0.0	0.0028	4.9515
0.0001	0.0076	1.8389
0.001	0.0051	2.7269
0.01	0.0138	1.0154
0.1	0.0140	0.9986

While the error in the goal functional is not influenced by the modification in the dual equation, the error estimator and thus the effectivity is. Notice, that for an exact evaluation of the residuals the effectivity is always one - since an error identity is evaluated. However, with a fixed integration accuracy smaller values of  $\varepsilon$  increase the quadrature error and consequently effectivity deteriorates. Once the mesh is sufficiently refined the quadrature - fixed per element - gains accuracy and thus the effectivity converges to one. The same effect has to be expected when numerically recovering the unknown primal and dual solutions for the weights, as the accuracy of the discrete primal and dual solutions are fixed on a given mesh and can only be increased by refinement.

A ratio of the advection to the diffusion is given by the Peclet number, see, e.g., [25]. Here,  $P_h$  shall be the approximation of the Peclet number for a constant advection velocity of one, depending on the mesh size as  $P_h = \frac{h}{\varepsilon}$ .

TABLE 3. Effectivity and approximated Peclet number for  $\varepsilon = 0.0001$  (left),  $\varepsilon = 0.01$  (middle), and  $\varepsilon = 0.1$  (right), with  $k = 10^{-4}$  constant and 40 quadrature points for a composite trapezoidal rule.

$h$	$\varepsilon = 0.0001$		$\varepsilon = 0.01$		$\varepsilon = 0.1$	
	$P_h$	eff	$P_h$	eff	$P_h$	eff
0.5	5000	-1.995	50	1.924	5	1.235
0.25	2500	-3.815	25	1.678	2.5	1.065
0.125	1250	7.290	12.5	1.277	1.25	1.016
0.0625	625	1.885	6.26	1.018	0.626	1.000

Table 3 shows that the effectivity of the global error estimator is getting better, the more the diffusion is of influence in the discretized scheme. Thus, it is suggested that, if the diffusion is resolved sufficiently, the modified dual weighted residual error estimator gives an effective approximation of the global error in the goal functional.

Concluding, these experiments suggest that in this setting the modified dual weighted residual error estimator for a spatial refinement is a reasonable indicator for grid refinement with respect to some goal functional and moreover the modified global error estimator is in this case of discontinuities a better approximation of the actual global error than the classical approach.

**6.2. 2D Example.** In this section, we consider problem (1) with  $\Omega = (-2, 2) \times (-2, 2)$  and a goal functional causing discontinuities in the dual solution, as in the previous section. But here we will also address the topic of patch wise interpolation for approximation of the discretization error.

6.2.1. *2D Setting.* We assume the initial condition

$$(21) \quad u_0(x, y, 0) = u_{\text{ini}}(x, y) = \begin{cases} 1, & -1 \leq x, y \leq 0, \\ 0, & \text{else.} \end{cases}$$

For this initial condition and  $f(\cdot)$  chosen as the identity, the solution of the weak advection equation is a translation along the characteristics, namely,

$$(22) \quad u_0(x, y, t) = u_{\text{ini}}(x - t, y - t) = \begin{cases} 1, & -1 + t \leq x, y \leq t, \\ 0, & \text{else.} \end{cases}$$

Since this problem does not include viscosity, i.e.,  $\varepsilon = 0$ , the solution is, again, denoted with subscript 0.

Choosing the goal functional analogously to the 1D test case as

$$(23) \quad J(u_0) = \int_{\mathbb{R}^2} u_0(x, y, T) z_T(x, y) \, dx \, dy,$$

with the weight  $z_T$  indicating an area of interest

$$(24) \quad z_T(x, y) := \begin{cases} 1, & 0 \leq x, y \leq 1, \\ 0, & \text{else,} \end{cases}$$

gives again a dual problem of the above advection equation. The analytic solution  $z_0$ , which coincides with  $u_0$  as in the 1D example, can be obtained by the application of a Green's function. But for the sake of generality, we assume the analytic solution as unknown.

6.2.2. *Discretization schemes and approximation of weights.* As in the 1D test case, the evaluation of  $a_0(\cdot, \cdot)$  at  $(u_0, z_0 - z_0^{hk})$  is not well defined. Hence, the residuals are evaluated with the more regular dual solution of the advection-diffusion equation,  $z_\varepsilon$  and  $z_\varepsilon^{hk}$ , respectively. Again, we use a second order DG scheme for space discretization, but now on the domain  $\bar{\Omega} = [-2, 2] \times [-2, 2]$  and an explicit Euler time stepping. For the approximation of the spatial discretization errors  $u_0 - u_0^h$  and  $z_\varepsilon - z_\varepsilon^h$  we use a patch-wise linear interpolation of the discrete solutions, see [5], resulting in the approximation of the weights by

$$(25) \quad w_u = u_0 - u_0^{hk} \approx I_{hk}^{(1)} u_0^{hk} - u_0^{hk}$$

and

$$(26) \quad w_z = z_\varepsilon - z_\varepsilon^{hk} \approx I_{hk}^{(1)} z_\varepsilon^{hk} - z_\varepsilon^{hk},$$

with  $I_{hk}^{(1)} := I_h^{(1)} I_k^{(1)}$ .

For the spatial interpolation on an element, we used a linear interpolation between the function values of three neighboring nodes, e.g.,  $P_1, P_4, P_5$ , to obtain the function value in  $Q$ , see Fig. 6. We obtained function values in  $P$  and  $R$  analogously by linear interpolation between  $P_3, P_5, P_6$  and  $P_2, P_4, P_6$ . The values in  $P, Q$ , and  $R$  define a plane, which allows linear extrapolations back to the nodes  $P_1$  to  $P_6$ .

For our explicit Euler method in time, we use a linear interpolation to approximate the discretization error. [15] did this for the implicit Euler method. Interpolation such that the effectivity of the error estimator is converging towards one is up to now only known for the Euler method, but not for Runge-Kutta methods. The



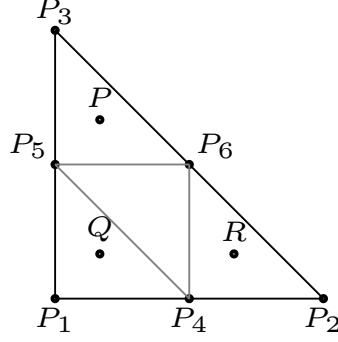


FIGURE 6. Nodes of quadratic basis functions,  $P_i, i = 1, \dots, 6$ , and interpolation points  $P, Q$ , and  $R$ .

explicit Euler and the corresponding linear interpolation are used in the application section of this paper.

For the explicit Euler method, the linear interpolation for equidistant time steps  $k = t_i - t_{i-1}, i = 1, \dots, M$  is

$$I_k^{(1)} u_0^{hk}(t) = \frac{t_i - t}{k} u_{0,i}^{hk} + \frac{t - t_{i-1}}{k} u_{0,i+1}^{hk},$$

for  $t \in [t_{i-1}, t_i]$ .

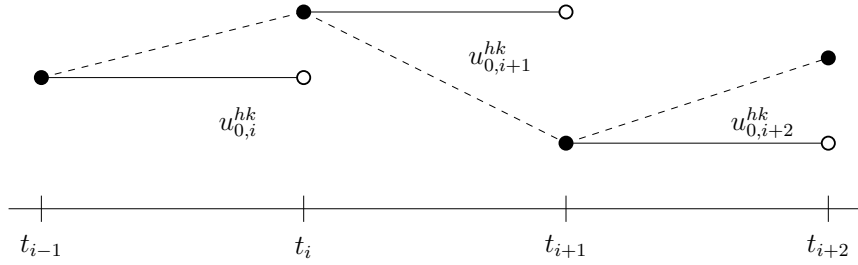


FIGURE 7. Linear interpolation (dashed line) of the piecewise constant explicit Euler solution.

With this interpolation operator, the derivative for  $t \in (t_{i-1}, t_i)$  is

$$\partial_t I_k^{(1)} u_0^{hk}(t) = \frac{u_{0,i+1}^{hk} - u_{0,i}^{hk}}{k},$$

which is of use for the discrete evaluation of the residuals.

The discrete dual solution  $z_\varepsilon^{hk}$  is interpolated in space and time in the same way and the error estimator (18) can be evaluated.

**6.2.3. Numerical experiments in 2D.** In this section, we study the dependence of the absolute value of the additional residual on the spatial grid size in 2D, as well as the behavior of the local error estimators with and without the additional dual residual. We find again the global error estimator including the additional residual to gain a better effectivity index as the formal global estimator.

In the following, we used a spatial DG discretization on  $\bar{\Omega} = [-2, 2] \times [-2, 2]$  with basis and test function polynomials of order 2. We used a composite box rule on each element for numerical integration. The value of the goal functional for the discrete solution,  $J(u_0^{hk})$ , is also determined exactly by a second order spatial

quadrature rule. The goal value of the analytic solution,  $J(u_0)$ , is one, as in the previous example.

In this setting, we computed the global value of the additional residual for different spatial resolutions, a fixed time step size of  $k = 10^{-5}$  and a fixed diffusion coefficient  $\varepsilon = 0.1$ . Fig. 8 shows the absolute global value of the additional residual, namely  $\rho^* \left( I_{hk}^{(1)} z_\varepsilon^{hk}; I_{hk}^{(1)} u_0^{hk} - u_0^{hk} \right)$ . In contrast to the 1D test case, the ad-

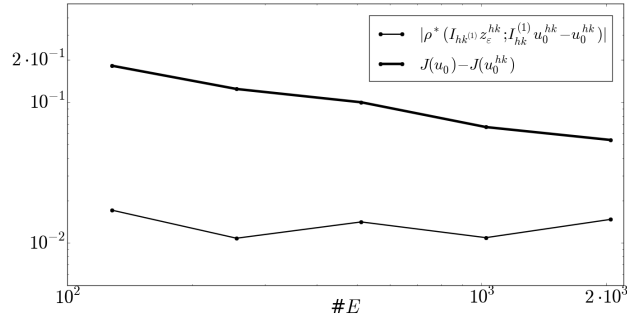


FIGURE 8. Absolute value of the additional residual,  $k = 10^{-5}$ ,  $\varepsilon = 0.1$ .

ditional term seems not to converge to zero. However, with the used parameters ( $\#E = 128, 256, 512, 1024, 2048$  and  $\varepsilon = 0.1$ ,  $k = 10^{-5}$ ) the additional residual is approximately ten times smaller than the error in the quantity of interest.

The formal elementwise error estimator, analogous to (20), neglects the diffusion and the additional dual residual. Both estimators, modified and formal, provide similar error indicators for grid refinement, as in 1D. The effectivity index for the example at hand is shown in Table 4. For a fixed time step size of  $k = 10^{-5}$  and a dual viscosity  $\varepsilon = 0.1$ , the effectivity index of the modified error estimator improves, though not monotone.

$\#E$	$\eta_{\text{uni}}^{hk}$	$ J(u_0) - J(u_0^{hk}) $	eff
128	0.0615	0.1818	2.9584
256	0.0247	0.1242	5.0346
512	0.0467	0.0999	2.1386
1024	0.0283	0.0666	2.3554
2048	0.0477	0.0539	1.130

TABLE 4. Dependence of the modified global error estimator and the error in the goal functional on the grid size (2D)

These results have to be considered in relation to the performance of the formal error estimator. Fig. 9 shows dependency of the effectivity index of the formal error estimator on the total number of elements of the uniform grid. It also depicts the effectivity index of the modified error estimator with and without the additional residual. The efficiency of the formal error estimator is worse than the efficiency of the modified one. This is even true if the residuals are evaluated with the solution of the dual advection-diffusion equation, but without the additional residual  $\rho^* \left( I_{hk}^{(1)} z_\varepsilon^{hk}; I_{hk}^{(1)} u_0^{hk} - u_0^{hk} \right)$ . With the additional term though, the efficiency is even better, see also Table 5.

Although further investigation with different parameters is needed, we conclude that our modified DWR error estimator proposed in this paper is a more efficient estimator as the formal one, also in the general case with approximated weights.

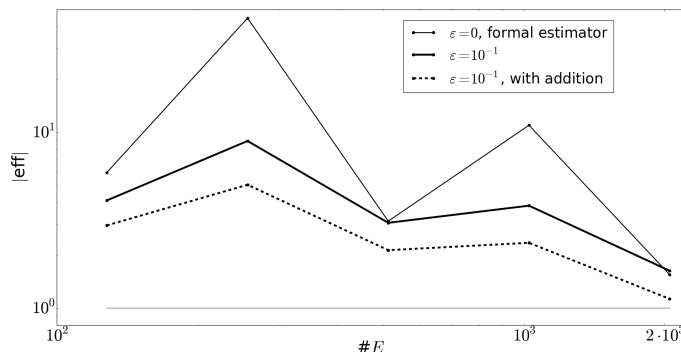


FIGURE 9. Effectivity of the formal global estimator and the modified estimator with and without the additional residual, with viscosity  $\varepsilon = 0.1$  in the dual equation.

TABLE 5. Effectivity index of the modified and of the formal error estimator

$\#E$	eff  (modified)	eff  (formal)
128	2.9584	5.9125
256	5.0346	44.624
512	2.1386	3.126
1024	2.3554	10.977
2048	1.130	1.552

#### ACKNOWLEDGMENTS

The first author acknowledges partial support by *Forschungs- und Wissenschaftsstiftung Hamburg*. The first and third author acknowledges the partial support by the DFG SFB-TRR154 in project A08.

We also like to thank Dr. Stefan Vater from the University of Hamburg/ClisAP for the code he shared, which gave a basis for the 1D dG advection code. Last but not least, the authors would like to thank Stefan Ulbrich for some insightful discussions on the subtleties of adjoints for hyperbolic problems.

#### REFERENCES

- [1] R. A. ADAMS AND J. J. F. FOURNIER, *Sobolev Spaces*, vol. 140 of Pure and Applied Mathematics (Amsterdam), Elsevier/Academic Press, Amsterdam, second ed., 2003.
- [2] M. AINSWORTH AND J. T. ODEN, *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics (New York), Wiley-Interscience [John Wiley & Sons], New York, 2000, doi:10.1002/9781118032824.
- [3] I. BABUŠKA, R. DURÁN, AND R. RODRÍGUEZ, *Analysis of the efficiency of an a posteriori error estimator for linear triangular finite elements*, SIAM J. Numer. Anal., 29 (1992), pp. 947–964, doi:10.1137/0729058.
- [4] I. BABUŠKA AND W. C. RHEINOLDT, *A-posteriori error estimates for the finite element method*, Internat. J. Numer. Methods Engrg., 12 (1978), pp. 1597–1615, doi:10.1002/nme.1620121010.
- [5] W. BANGERTH AND R. RANNACHER, *Adaptive Finite Element Methods for Differential Equations*, Lectures in Mathematics ETH Zürich, Birkhäuser, Basel, Boston, Berlin, 1. ed., 2003.
- [6] C. BARDOS, A. Y. LEROUX, AND J. C. NEDELEC, *First order quasilinear equations with boundary conditions*, Comm. Partial Differential Equations, 4 (1979), pp. 1017–1034, doi:10.1080/03605307908820117.
- [7] R. BECKER AND R. RANNACHER, *An optimal control approach to a posteriori error estimation*, Acta Numer., 10 (2001), pp. 1–102, doi:10.1017/S0962492901000010.

- [8] F. BOUCHUT AND F. JAMES, *One-dimensional transport equations with discontinuous coefficients*, *Nonlinear Anal.*, 32 (1998), pp. 891–933, doi:10.1016/S0362-546X(97)00536-1.
- [9] B. COCKBURN, G. E. KARNIADAKIS, AND C.-W. SHU, *The development of discontinuous Galerkin methods*, in *Discontinuous Galerkin methods* (Newport, RI, 1999), vol. 11 of *Lect. Notes Comput. Sci. Eng.*, Springer, Berlin, 2000, pp. 3–50, doi:10.1007/978-3-642-59721-3\_1.
- [10] W. DÖRFLER, *A convergent adaptive algorithm for Poisson’s equation*, *SIAM J. Numer. Anal.*, 33 (1996), pp. 1106–1124, doi:10.1137/0733054.
- [11] K. ERIKSSON, D. ESTEP, P. HANSBO, AND C. JOHNSON, *Introduction to adaptive methods for differential equations*, *Acta Numer.*, 4 (1995), pp. 105–158, doi:10.1017/S0962492900002531.
- [12] M. B. GILES, N. PIERCE, AND E. SÜLI, *Progress in adjoint error correction for integral functionals*, *Comput. Visual Sci.*, 6 (2004), pp. 113–121, doi:10.1007/BF02663040.
- [13] M. B. GILES AND S. ULBRICH, *Convergence of linearized and adjoint approximations for discontinuous solutions of conservation laws. part 1: Linearized approximations and linearized output functionals*, *SIAM J. Numer. Anal.*, 48 (2010), pp. 882–904, doi:10.1137/080727464.
- [14] M. B. GILES AND S. ULBRICH, *Convergence of linearized and adjoint approximations for discontinuous solutions of conservation laws. part 2: Adjoint approximations and extensions*, *SIAM J. Numer. Anal.*, 48 (2010), pp. 905–921, doi:10.1137/09078078X.
- [15] C. GOLL, R. RANNACHER, AND W. WOLLNER, *The damped Crank-Nicolson time-marching scheme for the adaptive solution of the Black-Scholes equation*, *J. Comput. Finance*, 18 (2015), pp. 1–37, doi:10.21314/JCF.2015.301.
- [16] R. HARTMANN, *Adjoint consistency analysis of discontinuous galerkin discretizations*, *SIAM J. Numer. Anal.*, 45 (2007), pp. 2671–2696, doi:10.1137/060665117.
- [17] R. HARTMANN AND P. HOUSTON, *Adaptive discontinuous galerkin finite element methods for nonlinear hyperbolic conservation laws*, *SIAM J. Sci. Comput.*, 24 (2002), pp. 979–1004, doi:10.1137/S1064827501389084.
- [18] C. JOHNSON AND A. SZEPESSY, *Adaptive finite element methods for conservation laws based on a posteriori error estimates*, *Comm. Pure Appl. Math.*, 48 (1995), pp. 199–234, doi:10.1002/cpa.3160480302.
- [19] S. MARTIN, *First order quasilinear equations with boundary conditions in the  $L^\infty$  framework*, *J. Differential Equations*, 236 (2007), pp. 375–406, doi:10.1016/j.jde.2007.02.007.
- [20] D. MEIDNER, R. RANNACHER, AND J. VIHAREV, *Goal-oriented error control of the iterative solution of finite element equations*, *J. Numer. Math.*, 17 (2009), pp. 143–172, doi:10.1515/JNUM.2009.009.
- [21] C. MEYER, A. RADEMACHER, AND W. WOLLNER, *Adaptive optimal control of the obstacle problem*, *SIAM J. Sci. Comput.*, 37 (2015), pp. A918–A945, doi:10.1137/140975863.
- [22] I. NEITZEL AND B. VEXLER, *A priori error estimates for space-time finite element discretization of semilinear parabolic optimal control problems*, *Numer. Math.*, 120 (2012), pp. 345–386, doi:10.1007/s00211-011-0409-9.
- [23] R. H. NOCHETTO AND A. VEESER, *Primer of adaptive finite element methods*, in *Multiscale and Adaptivity: Modeling, Numerics and Applications*, G. Naldi and G. Russo, eds., vol. 2040 of *Lecture Notes in Mathematics*, Springer Verlag, 2012, pp. 125–226, doi:10.1007/978-3-642-24079-9\_3.
- [24] M. OHLBERGER, *A review of a posteriori error control and adaptivity for approximations of non-linear conservation laws*, *Int. J. Numer. Meth. Fluids*, 59 (2009), pp. 333–354, doi:10.1002/flid.1686.
- [25] S. V. PATANKAR, *Numerical heat transfer and fluid flow*, *Series in Computational Methods in Mechanics and Thermal Sciences*, CRC press, 1980.
- [26] N. A. PIERCE AND M. B. GILES, *Adjoint recovery of superconvergent functionals from PDE approximations*, *SIAM Rev.*, 42 (2000), pp. 247–264, doi:10.1137/S0036144598349423.
- [27] P. W. POWER, M. D. PIGGOTT, F. FANG, G. J. GORMAN, C. C. PAIN, D. P. MARSHALL, A. J. H. GODDARD, AND I. M. NAVON, *Adjoint goal-based error norms for adaptive mesh ocean modelling*, *Ocean modelling*, 15 (2006), pp. 3–38, doi:10.1016/j.ocemod.2006.05.001.
- [28] R. RANNACHER, A. WESTENBERGER, AND W. WOLLNER, *Adaptive finite element solution of eigenvalue problems: Balancing of discretization and iteration error*, *J. Numer. Math.*, 18 (2010), pp. 303–327, doi:10.1515/JNUM.2010.015.
- [29] F. RAUSER, P. KORN, AND J. MAROTZKE, *Predicting goal error evolution from near-initial-information: A learning algorithm*, *J. Comput. Physics*, 230 (2011), pp. 7284–7299, doi:10.1016/j.jcp.2011.05.029.
- [30] W. H. REED AND T. R. HILL, *Triangular mesh methods for the neutron transport equation*, *Los Alamos Report LA-UR-73-479*, (1973).

- [31] T. RICHTER AND T. WICK, *Variational localizations of the dual weighted residual estimator*, J. Comput. Appl. Math., 279 (2015), pp. 192–208, doi:10.1016/j.cam.2014.11.008.
- [32] B. RIVIÈRE, *Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations*, vol. 35 of Frontiers in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2008, doi:10.1137/1.9780898717440. Theory and implementation.
- [33] M. SCHMICH AND B. VEXLER, *Adaptivity with dynamic meshes for space-time finite element discretizations of parabolic equations*, SIAM J. Sci. Comput., 30 (2008), pp. 369–393, doi:10.1137/060670468.
- [34] J. SCHÜTZ, G. MAY, AND S. NOELLE, *Analytical and numerical investigation of the influence of artificial viscosity in discontinuous Galerkin methods on an adjoint-based error estimator*, in Computational fluid dynamics 2010, Springer, 2011, pp. 203–209, doi:10.1007/978-3-642-17884-9\_24.
- [35] E. SÜLI AND P. HOUSTON, *Adaptive finite element approximation of hyperbolic problems*, in Error Estimation and Adaptive Discretization Methods in Computational Fluid Dynamics, vol. 25 of Lecture Notes in Computational Science and Engineering, Springer Berlin Heidelberg, 2003, pp. 269–344, doi:10.1007/978-3-662-05189-4\_6.
- [36] S. ULBRICH, *Optimal Control of Nonlinear Hyperbolic Conservation Laws with Source Terms*, habilitation, Technische Universität München, 2001.
- [37] S. ULBRICH, *A sensitivity and adjoint calculus for discontinuous solutions of hyperbolic conservation laws with source terms*, SIAM J. Control Optim., 41 (2002), pp. 740–797, doi:10.1137/S0363012900370764.
- [38] S. ULBRICH, *Adjoint-based derivative computations for the optimal control of discontinuous solutions of hyperbolic conservation laws*, Systems Control Lett., 48 (2003), pp. 313–328, doi:10.1016/S0167-6911(02)00275-X.
- [39] D. A. VENDITTI AND D. L. DARMOFAL, *Adjoint error estimation and grid adaptation for functional outputs: Application to quasi-one-dimensional flow*, J. Comput. Physics, 164 (2000), pp. 204–227, doi:10.1006/jcph.2000.6600.
- [40] R. VERFÜRTH, *A Posteriori Error Estimation Techniques for Finite Element Methods*, Oxford University Press, 2013, doi:10.1093/acprof:oso/9780199679423.001.0001.
- [41] M. F. WHEELER, *An elliptic collocation-finite element method with interior penalties*, SIAM J. Numer. Anal., 15 (1978), pp. 152–161, doi:10.1137/0715010.
- [42] W. WOLLNER, *A posteriori error estimates for a finite element discretization of interior point methods for an elliptic optimization problem with state constraints*, Comput. Optim. Appl., 47 (2010), pp. 133–159, doi:10.1007/s10589-008-9209-2.

(S. Beckers) DEPT OF MATHEMATICS, TECHNISCHE UNIVERSITÄT DARMSTADT, DOLIVOSTR. 15, 64293 DARMSTADT, GERMANY  
*E-mail address:* [beckers@mathematik.tu-darmstadt.de](mailto:beckers@mathematik.tu-darmstadt.de)

(J. Behrens) DEPT OF MATHEMATICS, UNIVERSITÄT HAMBURG, BUNDESSTRASSE 55, 20146 HAMBURG, GERMANY  
*E-mail address:* [joernbehrens@uni-hamburg.de](mailto:joernbehrens@uni-hamburg.de)

(W. Wollner) DEPT OF MATHEMATICS, TECHNISCHE UNIVERSITÄT DARMSTADT, DOLIVOSTR. 15, 64293 DARMSTADT, GERMANY  
*E-mail address:* [wollner@mathematik.tu-darmstadt.de](mailto:wollner@mathematik.tu-darmstadt.de)