

# Robust DPG Methods for Transient Convection-Diffusion

Truman Ellis, Jesse Chan, Leszek Demkowicz

## Abstract

We introduce two robust DPG methods for transient convection-diffusion problems. Once a variational formulation is selected, the choice of test norm critically influences the quality of a particular DPG method. It is desirable that a test norm produce convergence of the solution in a norm equivalent to  $L^2$  while producing optimal test functions that can be accurately computed and maintaining good conditioning of the optimal test function solve on highly adaptive meshes. Two such *robust* norms are introduced and proven to guarantee close to  $L^2$  convergence of the primary solution variable. Numerical experiments demonstrate robust convergence of the two methods.

## 1 Introduction

The discontinuous Petrov-Galerkin finite element method presents an attractive new framework for developing robust numerical methods for computational mechanics. DPG contains the promise of being an automated scientific computing technology – it provides stability for any variational formulation, optimal convergence rates in a user-defined norm, virtually no pre-asymptotic stability issues on coarse meshes, and a measure of the error residual which can be used to robustly drive adaptivity. The method also delivers Hermitian positive definite stiffness matrices for any problem, weak enforcement of boundary conditions, and several other attractive properties [11, 12, 9]. For the most recent review of DPG, see [13]. The process of developing robust DPG methods for steady convection-diffusion was explored in [15, 4]. *In the sense, the main challenge is to come up with a correct test norm. The residual is measured in the dual test norm, and the DPG method minimizes the residual. The residual can be interpreted as a special energy norm. In other words, the DPG method delivers an orthogonal projection in the energy norm. The task is especially challenging for singular perturbation problems. Given a trial norm, we strive to determine a quasi-optimal test norm such that the corresponding energy norm is robustly equivalent to the trial norm of choice. An additional difficulty comes from the fact that the optimal test functions should be easily approximated with a simple enrichment strategy. For convection dominated diffusion, this means that the test functions should not develop boundary layers. The task of determining the quasi optimal test norm (we call it a *robust test norm* leads then to a stability analysis for the adjoint equation which is the subject of this paper. For a more general discussion on the subject, see [13]. We start with an abstract derivation of the DPG framework then define the concept of a robust test norm and specialize to transient convection-diffusion. Two new robust norms are derived and numerical verifications of the theory are presented.*

*It is worth mentioning connections to other modern stabilized finite element methods. DPG can be thought of as a generalization of least-squares finite element methods, and in fact simplifies to this case when the  $L^2$  topology is chosen for the test space. Connections to multi-scale methods have been studied in [8] and [6]. In the particular case of ultra-weak variational formulation and scaled adjoint graph norm used for the test norm, the DPG method delivers optimal test functions of Barrett and Morton (see [7]). For each trial basis function, these test functions are defined as solutions of the global adjoint equation with the basis function as the forcing term. It can be shown that the Petrov-Galerkin method with the resulting test space delivers  $L^2$ -projection. DPG approximates this test space (though not each individual optimal test function) using an enriched (i.e. multiscale) mortar least squares formulation. In a similar sense, DPG can also be interpreted as a Variational Multiscale Method which approximates the fine-scale contribution by enforcing the orthogonality of the fine scales under a specific inner product. The connection to HDG is less clear; DPG shares much with the more classical mortar method, but its connection to HDG has primarily been explored in [22]. Similarities include the use of trace unknowns which live on the mesh skeleton and*

the use of static condensation to eliminate the interior degrees of freedom, but the means of stabilization for the two methods are very different.

## 1.1 Space-Time Finite Elements

Most finite element simulations of transient phenomena use a semi-discrete formulation: the PDE is first discretized in space using finite elements and then the leftover system of ordinary differential equations in time is usually solved by a finite difference method. But it is possible to treat time as just another dimension to be discretized with finite elements. Some of the earliest proponents of this approach were Kaczkowski[20], Argyris and Scharpf[2], Fried[17], and Oden[25]. These techniques were built on the underlying concept of Hamilton's principle.

Van der Vegt and van der Ven[30] have advanced a space-time discontinuous Galerkin method for 3D inviscid compressible moving boundary problems. Klaij *et al.* [21] then extended the method to compressible Navier-Stokes while Rhebergen *et al.* [27] developed the method for incompressible Navier-Stokes. Rhebergen and Cockburn[26] also developed a space-time HDG method for incompressible Navier-Stokes. Tezduyar and Behr[29] develop a deforming-spatial-domain/space-time procedure coupled with Galerkin/least-squares to handle incompressible Navier-Stokes flows with moving boundaries and later Aliabadi and Tezduyar[1] apply the procedure to compressible flows. Hughes and Stewart[19] develop a general space-time multiscale framework for deriving stabilized methods for transient phenomena.

It is possible to use the semi-discrete approach to solving transient problems with DPG, but it doesn't appear to be a natural fit with the adaptive nature of DPG. The Courant-Friedrichs-Lewy (CFL) condition is not binding with implicit time integration schemes, but it can be a guiding principle for temporal accuracy. If we are interested in temporally accurate solutions, we are limited by the fact that our smallest mesh elements (which may be orders of magnitude smaller than the largest elements) are constrained to proceed at a much smaller time step than the mesh as a whole. We can either restrict the whole mesh to the smallest time step, or we can attempt some sort of local time stepping. A space-time DPG formulation presents an attractive choice as we will be able to preserve our natural adaptivity from the steady problems while extending it in time. Thus we achieve an adaptive solution technique for transient problems in a unified framework. This paper expands previous work developing robust DPG methods for steady convection-diffusion to the space-time form.

## 2 Overview of DPG

### 2.1 A Generalized Minimum Residual Method

We begin with a well posed variational problem: find  $u \in U$  such that

$$b(u, v) = l(v) \quad \forall v \in V$$

where  $b(u, v)$  is a bilinear (sesquilinear) form on  $U \times V$  and  $l \in V'$ . Introducing operator  $B : U \rightarrow V'$  ( $V'$  is the dual space to  $V$ ) defined by  $b(u, v) = \langle Bu, v \rangle_{V' \times V}$ , we can reformulate the equation in operator form:

$$Bu = l \in V'.$$

We wish to find the element  $u_h$  of a finite dimensional subspace which minimizes the residual  $Bu - l$  in  $V'$ :

$$u_h = \arg \min_{w_h \in U_h} \frac{1}{2} \|Bu - l\|_{V'}^2.$$

This mathematical framework is very natural, but it is not yet practical as the  $V'$  norm is not especially amenable to computations. With the assumption that we are working with Hilbert spaces, we can use the Riesz representation theorem to find a complementary object in  $V$  rather than  $V'$ . Let  $R_V : V \ni v \rightarrow (v, \cdot) \in V'$  be the Riesz map, which is an isometry. Then the inverse Riesz map lets us represent our residual in  $V$ :

$$u_h = \arg \min_{w_h \in U_h} \frac{1}{2} \|R_V^{-1}(Bu - l)\|_V^2.$$

Since this is a convex minimization problem, the solution is given by the critical points where the Gâteaux derivative is zero in all directions  $\delta u \in U_h$ :

$$(R_V^{-1}(Bu_h - l), R_V^{-1}B\delta u)_V = 0, \quad \forall \delta u \in U.$$

By definition of the Riesz map this is equivalent to the duality pairing

$$\langle Bu_h - l, R_V^{-1}B\delta u_h \rangle = 0 \quad \forall \delta u_h \in U_h.$$

We can define an optimal test function  $v_{\delta u_h} := R_V^{-1}B\delta u_h$  for each trial function  $\delta u_h$ . This allows us to revert back to our original bilinear form with a finite dimensional set of trial and test functions:

$$b(u_h, v_{\delta u_h}) = l(v_{\delta u_h}).$$

Note that  $v_{\delta u_h} \in V$  comes from the auxiliary problem

$$(v_{\delta u_h}, \delta v)_V = \langle R_V v_{\delta u_h}, \delta v \rangle = \langle B\delta u_h, \delta v \rangle = b(\delta u_h, \delta v) \quad \forall \delta v \in V.$$

We might refer to this as an *optimal Petrov-Galerkin* method. We arrive at the same method by realizing the supremum in the inf-sup condition (see [12]), motivating the *optimal* nomenclature. As a minimum residual method, optimal Petrov-Galerkin methods produce Hermitian, positive-definite stiffness matrices since

$$b(u_h, v_{\delta u_h}) = (v_{u_h}, v_{\delta u_h})_V = \overline{(v_{\delta u_h}, v_{u_h})} = \overline{b(\delta u_h, v_{u_h})}.$$

The energy norm of the error is directly related to the residual:

$$\|u_h - u\|_E = \|B(u_h - u)\|_{V'} = \|Bu_h - l\|_{V'} = \|R_V^{-1}(Bu_h - l)\|_V,$$

where we designate  $R_V^{-1}(Bu_h - l)$  the *error representation function*. This has proven to be a very robust *a-posteriori* error estimator for driving adaptivity [14].

## 2.2 Transient Convection-Diffusion

### 2.2.1 Problem Description

In order to better illustrate choice of the  $U$  and  $V$  spaces, we introduce the transient convection-diffusion problem. Consider spatial domain  $\Omega$  and corresponding space-time domain  $Q = \Omega \times [0, T]$  with boundary  $\Gamma = \Gamma_- \cup \Gamma_+ \cup \Gamma_0 \cup \Gamma_T$  where  $\Gamma_-$  is the inflow boundary ( $\beta \cdot \mathbf{n}_x < 0$ , where  $\beta$  is the convection vector and  $\mathbf{n}_x$  is the outward spatial normal),  $\Gamma_+$  is the outflow boundary ( $\beta \cdot \mathbf{n}_x \geq 0$ ),  $\Gamma_0$  is the initial time boundary, and  $\Gamma_T$  is the final time boundary. Let  $\Gamma_h := \bigcup \partial K$  denote the entire mesh skeleton, where  $\partial K$  denotes the boundary of element  $K$ .  $\Gamma_{h_x}$  denotes any parts of the skeleton with a nonzero spatial normal and  $\Gamma_{h_t}$  have a nonzero temporal normal.

The transient convection-diffusion equation is

$$\frac{\partial u}{\partial t} + \nabla \cdot (\beta u) - \epsilon \Delta u = f,$$

where  $u$  is the quantity of interest, often interpreted to be a concentration of some quantity,  $\epsilon$  is the diffusion coefficient, and  $f$  is the source term.

We apply flux boundary conditions on the inflow and trace boundary conditions on the outflow

$$\begin{aligned} \text{tr}(\beta \cdot u - \epsilon \nabla u) \cdot \mathbf{n}_x &= t_- \quad \text{on } \Gamma_- \\ \text{tr}(u) &= u_+ \quad \text{on } \Gamma_+ \\ \text{tr}(u) &= u_0 \quad \text{on } \Gamma_0. \end{aligned}$$

We note that Dirichlet boundary conditions also induce Dirichlet boundary conditions for the adjoint problem. Since the direction of convection is reversed for the adjoint convection-diffusion problem, this results in boundary layer adjoint solutions, which must be controlled using special weighted norms [28, 15]. However, since the convection-diffusion operator is not self-adjoint, the Cauchy inflow boundary condition induces a Neumann boundary condition for the adjoint problem. As a result, the adjoint solution does not contain boundary layers, simplifying the construction of a robust DPG method.

### 2.2.2 Relevant Sobolev Spaces

We begin by defining operators  $\nabla_{xt}u := \begin{pmatrix} \nabla u \\ \frac{\partial u}{\partial t} \end{pmatrix}$  and  $\nabla_{xt} \cdot \mathbf{u} := \nabla \cdot \mathbf{u}_x + \frac{\partial u_t}{\partial t}$ , where  $\mathbf{u} = (\mathbf{u}_x, u_t)$ . We will need the following Sobolev spaces defined on our space-time domain.

$$\begin{aligned} H^1(Q) &= \{u \in L^2(Q) : \nabla u \in \mathbf{L}^2(Q)\} \\ H_{xt}^1(Q) &= \{u \in L^2(Q) : \nabla_{xt}u \in \mathbf{L}^2(Q)\} \\ \mathbf{H}(\text{div}, Q) &= \{\boldsymbol{\sigma} \in \mathbf{L}^2(Q) : \nabla \cdot \boldsymbol{\sigma} \in L^2(Q)\} \\ \mathbf{H}(\text{div}_{xt}, Q) &= \{\boldsymbol{\sigma} \in \mathbf{L}^2(Q) : \nabla_{xt} \cdot \boldsymbol{\sigma} \in L^2(Q)\} \end{aligned}$$

We will also need the corresponding broken Sobolev spaces.

$$\begin{aligned} H^1(Q_h) &= \{u \in L^2(Q) : u|_K \in H^1(K), K \in Q_h\} &= \prod_{K \in Q_h} H^1(K) \\ H_{xt}^1(Q_h) &= \{u \in L^2(Q) : u|_K \in H_{xt}^1(K), K \in Q_h\} &= \prod_{K \in Q_h} H_{xt}^1(K) \\ \mathbf{H}(\text{div}, Q_h) &= \{\boldsymbol{\sigma} \in \mathbf{L}^2(Q) : u|_K \in \mathbf{H}(\text{div}, K), K \in Q_h\} &= \prod_{K \in Q_h} \mathbf{H}(\text{div}, K) \\ \mathbf{H}(\text{div}_{xt}, Q_h) &= \{\boldsymbol{\sigma} \in \mathbf{L}^2(Q) : u|_K \in \mathbf{H}(\text{div}_{xt}, K), K \in Q_h\} &= \prod_{K \in Q_h} \mathbf{H}(\text{div}_{xt}, K) \end{aligned}$$

Consider the following trace operators:

$$\begin{aligned} \text{tr}_{\text{grad}}^K u &= u|_{\partial K_x} & u &\in H^1(K) \\ \text{tr}_{\text{div}_{xt}}^K \boldsymbol{\sigma} &= \boldsymbol{\sigma}|_{\partial K_{xt}} \cdot \mathbf{n}_{K_{xt}} & \boldsymbol{\sigma} &\in \mathbf{H}(\text{div}_{xt}, K) \end{aligned}$$

where  $\partial K_x$  refers to spatial faces of element  $K$ ,  $\partial K_{xt}$  to the full space-time boundary, and  $\mathbf{n}_{K_{xt}}$  is the unit outward normal on  $\partial K_{xt}$ . The operators  $\text{tr}_{\text{grad}}$  and  $\text{tr}_{\text{div}_{xt}}$  perform the same operation element by element to produce the linear maps

$$\begin{aligned} \text{tr}_{\text{grad}} : H^1(Q_h) &\rightarrow \prod_{K \in Q_h} H^{1/2}(\partial K_x) \\ \text{tr}_{\text{div}_{xt}} : \mathbf{H}(\text{div}_{xt}, Q_h) &\rightarrow \prod_{K \in Q_h} H^{-1/2}(\partial K_{xt}) \end{aligned}$$

Finally, we define spaces of interface functions. In order that our functions be single valued, we use the following definitions.

$$\begin{aligned} H^{1/2}(\Gamma_{h_x}) &= \text{tr}_{\text{grad}} H^1(Q), \\ H_{xt}^{-1/2}(\Gamma_h) &= \text{tr}_{\text{div}_{xt}} \mathbf{H}(\text{div}_{xt}, Q). \end{aligned}$$

For more details on broken and trace Sobolev spaces, see [3].

### 2.2.3 Variational Formulations

There are many possible manipulations that could be performed before arriving at a variational formulation. We begin by reformulating the problem in terms of the first order system:

$$\begin{aligned} \frac{1}{\epsilon} \boldsymbol{\sigma} - \nabla u &= 0 \\ \nabla_{xt} \cdot \begin{pmatrix} \beta u - \boldsymbol{\sigma} \\ u \end{pmatrix} &= f. \end{aligned} \tag{1}$$

Multiplying (1) by test functions  $\boldsymbol{\tau} \in \mathbf{L}^2(Q)$  and  $v \in L^2(Q)$ , we obtain the following “trivial” variational formulation equivalent to the strong form:

$$\begin{aligned}
u &\in H_{xt}^1(Q) & u &= u_+ \quad \text{on } \Gamma_+ \\
& & u &= u_0 \quad \text{on } \Gamma_0 \\
\boldsymbol{\sigma} &\in \mathbf{H}(\text{div}, Q) & (\beta u - \epsilon \nabla u) \cdot \mathbf{n} &= t_- \quad \text{on } \Gamma_- \\
\left( \frac{1}{\epsilon} \boldsymbol{\sigma}, \boldsymbol{\tau} \right) - (\nabla u, \boldsymbol{\tau}) & & &= 0 \quad \forall \boldsymbol{\tau} \in \mathbf{L}^2(Q) \\
\left( \nabla_{xt} \cdot \begin{pmatrix} \beta u - \boldsymbol{\sigma} \\ u \end{pmatrix}, v \right) & & &= f \quad \forall v \in L^2(Q).
\end{aligned} \tag{2}$$

We can now choose either to relax (integrate by parts and build in the boundary conditions) or strongly enforce each equation. The steady state case and resulting options are explored and analyzed in further detail in [10] and are termed the trivial formulation (don’t relax anything), the classical formulation (relax the second equation), the mixed formulation (relax the first equation), and the ultra-weak formulation (relax both equations). The stability constants for the four formulations are related, but the functional settings and norms of convergence change. Early DPG work emphasized the ultra-weak formulation since in many ways it was the easiest to analyze, though recently the classical formulation has been under very active consideration. In the interests of simpler analysis, we focus on the ultra-weak formulation in this paper.

$$\begin{aligned}
u &\in L^2(Q), \boldsymbol{\sigma} \in \mathbf{L}^2(Q) \\
\left( \frac{1}{\epsilon} \boldsymbol{\sigma}, \boldsymbol{\tau} \right) + (u, \nabla \cdot \boldsymbol{\tau}) &= 0 \quad \forall \boldsymbol{\tau} \in \mathbf{H}(\text{div}, Q) : \boldsymbol{\tau} \cdot \mathbf{n}_x = 0 \text{ on } \Gamma_- \\
- \left( \begin{pmatrix} \beta u - \boldsymbol{\sigma} \\ u \end{pmatrix}, \nabla_{xt} v \right) &= f \quad \forall v \in H_{xt}^1(Q) : v = 0 \text{ on } \Gamma_+ \cup \Gamma_0,
\end{aligned} \tag{3}$$

We can remove the conditions on the test functions by introducing trace unknowns

$$\begin{aligned}
\hat{u} &= \text{tr}(u) & & \text{on } \partial Q_x \\
\hat{t} &= \text{tr} \left( \begin{pmatrix} \beta u - \boldsymbol{\sigma} \\ u \end{pmatrix} \cdot \mathbf{n}_{xt} \right) & & \text{on } \partial Q_{xt}.
\end{aligned}$$

Our new ultra-weak formulation with conforming test functions is

$$\begin{aligned}
u &\in L^2(Q), \boldsymbol{\sigma} \in \mathbf{L}^2(Q) \\
\hat{u} &\in H^{1/2}(\partial Q_x), & \hat{u} &= u_+ \text{ on } \Gamma_+ \\
\hat{t} &\in H_{xt}^{-1/2}(\partial Q), & \hat{t} &= t_- \text{ on } \Gamma_-, \quad \hat{t} = -u_0 \text{ on } \Gamma_0 \\
\left( \frac{1}{\epsilon} \boldsymbol{\sigma}, \boldsymbol{\tau} \right) + (u, \nabla \cdot \boldsymbol{\tau}) - \langle \hat{u}, \boldsymbol{\tau} \cdot \mathbf{n}_x \rangle &= 0 \quad \forall \boldsymbol{\tau} \in \mathbf{H}(\text{div}, Q) \\
- \left( \begin{pmatrix} \beta u - \boldsymbol{\sigma} \\ u \end{pmatrix}, \nabla_{xt} v \right) + \langle \hat{t}, v \rangle &= f \quad \forall v \in H_{xt}^1(Q).
\end{aligned} \tag{4}$$

#### 2.2.4 Broken Test Functions

One of the key insights that led to the development of the DPG framework was the process of breaking test functions, that is testing with functions from larger broken Sobolev spaces, replacing  $H_{xt}^1(Q)$  with  $H_{xt}^1(Q_h)$  and  $\mathbf{H}(\text{div}, Q)$  with  $\mathbf{H}(\text{div}, Q_h)$ . Discretizing such spaces is much simpler than standard spaces which require enforcement of global continuity conditions. The cost of introducing broken spaces is that we have to extend our interface unknowns  $\hat{u}$  and  $\hat{t}$  to live on the mesh skeleton. Our ultra-weak formulation with broken test

functions looks like

$$\begin{aligned}
u &\in L^2(Q), \sigma \in \mathbf{L}^2(Q) \\
\hat{u} &\in H^{1/2}(\Gamma_{h_x}), & \hat{u} &= u_+ \text{ on } \Gamma_+ \\
\hat{t} &\in H_{xt}^{-1/2}(\Gamma_h), & \hat{t} &= t_- \text{ on } \Gamma_-, \quad \hat{t} = -u_0 \text{ on } \Gamma_0 \\
\left(\frac{1}{\epsilon}\sigma, \tau\right) + (u, \nabla \cdot \tau) - \langle \hat{u}, \tau \cdot \mathbf{n}_x \rangle &= 0 \quad \forall \tau \in \mathbf{H}(\text{div}, Q_h) \\
-\left(\left(\begin{array}{c} \beta u - \sigma \\ u \end{array}\right), \nabla_{xt} v\right) + \langle \hat{t}, v \rangle &= f \quad \forall v \in H_{xt}^1(Q_h).
\end{aligned} \tag{5}$$

The main consequence of breaking test functions is that it reduces the cost of solving for optimal test functions from a global solve to an embarrassingly parallel solve element-by-element. Now that we've derived a suitable variational formulation, we are left with the task of selecting a test norm with which to compute our optimal test functions.

### 3 Robust Test Norms

The final unresolved choice is what norm to apply to the  $V$  space. This is one of the most important factors in designing a robust DPG method as the corresponding Riesz operator needs to be inverted to solve for the optimal test functions. If the norm produces unresolved boundary layers in the auxiliary problem, then many of the attractive features of DPG may fall apart. This is the primary emphasis of this paper. The problem of constructing stable test norms for steady convection-diffusion was addressed in [15, 4]. In this paper, we extend that work to transient convection-diffusion in space-time.

We define a robust test norm such that the  $L^2$  norm of the solution is bounded by the energy norm of the solution with a constant independent of  $\epsilon$ . We can rewrite any ultra-weak formulation with broken test functions as the following bilinear form with group variables:

$$b((u, \hat{u}), v) = (u, A^* v)_{L^2} + \langle \hat{u}, \llbracket v \rrbracket \rangle_{\Gamma_h}$$

where  $A^*$  represents the adjoint. In the case of convection-diffusion,  $u := \{u, \sigma\}$ ,  $\hat{u} := \{\hat{u}, \hat{t}\}$ ,  $v := \{v, \tau\}$ .

Note that for conforming  $v^*$  satisfying  $A^* v^* = u$

$$\begin{aligned}
\|u\|_{L^2}^2 &= b(u, v^*) = \frac{b(u, v^*)}{\|v^*\|_V} \|v^*\|_V \\
&\leq \sup_{v^* \neq 0} \frac{|b(u, v^*)|}{\|v^*\|} \|v^*\| = \|u\|_E \|v^*\|_V.
\end{aligned}$$

This defines a necessary condition for robustness, namely that

$$\|v^*\|_V \lesssim \|u\|_{L^2}. \tag{6}$$

If this condition is satisfied, then we get our final result:

$$\|u\|_{L^2} \lesssim \|u\|_E.$$

So far, we've assumed that our finite set of optimal test functions are assembled from an infinite dimensional space. In practice, we have found it to be sufficient to use an "enriched" space of higher polynomial dimension than the trial space [18]. This adds an additional requirement when assembling a robust test norm, namely that our optimal test functions should be adequately representable within this enriched space. We illustrate this point by considering three norms which satisfy the above conditions for 1D steady convection-diffusion. The graph norm is  $\left(\|A^* v\|_{L^2}^2 + \|v\|_{L^2}^2\right)^{\frac{1}{2}}$ :

$$\|(v, \tau)\|^2 = \|\nabla \cdot \tau - \beta \cdot \nabla v\|^2 + \left\| \frac{1}{\epsilon} \tau + \nabla v \right\|^2 + \|v\|^2 + \|\tau\|^2.$$

*Remark.* In the DPG technology, the test norm must be *localizable*, i.e.,

$$\|v\|_V^2 = \sum_K \|v\|_{V(K)}^2$$

where  $\|v\|_{V(K)}$  denotes a test norm (and not just a seminorm) for the element test space. In practice this means the addition of properly scaled  $L^2$ -terms. Without those terms, we could not invert the Riesz operator on the element level. Addition of the  $L^2$  terms does not necessarily contradict the robustness of the norm, see the discussion in [13] on bounded below operators. An alternate strategy has been explored in [16] where we enforce element conservation property by securing the presence of a constant function in the element test space. The residual is then minimized only over the orthogonal complement to the constants which eliminates the need for adding the  $L^2$ -term to the test norm.

The robust norm was derived in [4]:

$$\|(v, \boldsymbol{\tau})\|^2 = \|\boldsymbol{\beta} \cdot \nabla v\|^2 + \epsilon \|\nabla v\|^2 + \min\left(\frac{\epsilon}{h^2}, 1\right) \|v\|^2 + \|\nabla \cdot \boldsymbol{\tau}\|^2 + \min\left(\frac{1}{h^2}, \frac{1}{\epsilon}\right) \|\boldsymbol{\tau}\|^2.$$

The case for the coupled robust norm was made in [5]:

$$\|(v, \boldsymbol{\tau})\|^2 = \|\boldsymbol{\beta} \cdot \nabla v\|^2 + \epsilon \|\nabla v\|^2 + \min\left(\frac{\epsilon}{h^2}, 1\right) \|v\|^2 + \|\nabla \cdot \boldsymbol{\tau} - \boldsymbol{\beta} \cdot \nabla v\|^2 + \min\left(\frac{1}{h^2}, \frac{1}{\epsilon}\right) \|\boldsymbol{\tau}\|^2.$$

The argument for the coupled norm was that in certain cases we noticed pollution of  $u$  from errors in  $\boldsymbol{\sigma}$ , for example at singularities in  $\boldsymbol{\sigma}$ ,  $u$  also exhibited degraded quality with the robust norm. The coupled robust norm seemed to relax this behavior, i.e. errors in  $u$  appear more independent of errors in  $\boldsymbol{\sigma}$ .

The bilinear form and test norm define a mapping from input trial functions to an optimal test function:

$$T = R_V^{-1} B : U \rightarrow V.$$

Below, we plot the optimal test functions produced given  $\epsilon = 10^{-2}$ , a representative trial function  $u = x - \frac{1}{2}$ , and either the graph norm, the robust norm, or the coupled robust norm. Note that the optimal test functions will be different for any other trial function. In the left column, we see the fully resolved *ideal* optimal test function that DPG theory relies on. On the right, we see the approximated optimal test function using a enriched cubic test space.

Mathematically, the graph norm satisfies the necessary condition to be a robust norm, but the ideal optimal test functions contain strong boundary layers which can not be realistically approximated with the provided enriched space. If the approximated optimal test functions can not come sufficiently close to the ideal, then the whole DPG theory falls apart. See [18] for more discussion. This provides an additional condition on a test norm before we can truly call it robust: the ideal test functions must be adequately representable within the provided enriched space. This ultimately comes down to an analysis of the relative magnitudes of individual terms within the test norm, usually attempting to bound reactive or convective terms by diffusive terms. The coupled robust norm satisfies condition (6) and also produces relatively smooth optimal test functions that can be sufficiently approximated with a cubic polynomial space. Niemi *et al.* attempted to approximate boundary layers in optimal shape functions with Shishkin meshes [23, 24].

### 3.1 Application to Transient Convection-Diffusion

Now we present the analysis leading to two robust norms for transient convection-diffusion. Consider the problem with homogeneous boundary conditions

$$\begin{aligned} \frac{1}{\epsilon} \boldsymbol{\sigma} - \nabla u &= 0 \\ \frac{\partial u}{\partial t} + \boldsymbol{\beta} \cdot \nabla u - \nabla \cdot \boldsymbol{\sigma} &= f \\ \beta_n u - \epsilon \frac{\partial u}{\partial n} &= 0 \text{ on } \Gamma_- \\ u &= 0 \text{ on } \Gamma_+ \\ u &= u_0 \text{ on } \Gamma_0. \end{aligned}$$

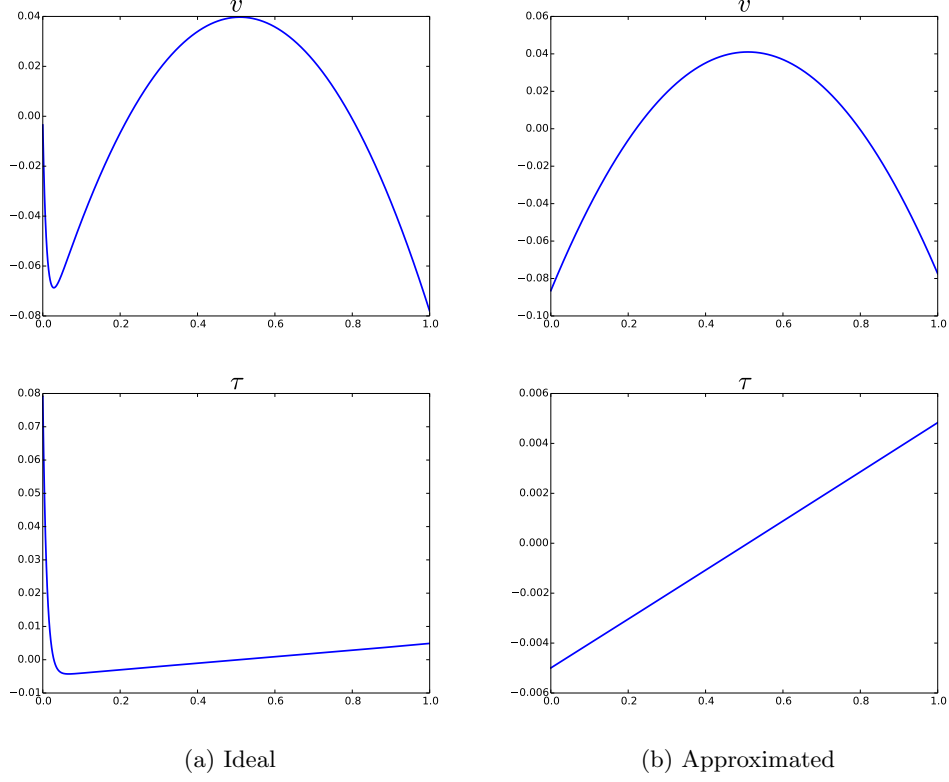


Figure 1: Graph norm optimal test functions for  $u = x - \frac{1}{2}$

Let  $\tilde{\beta} := \begin{pmatrix} \beta \\ 1 \end{pmatrix}$ , then we can rewrite this as

$$\begin{aligned}
 \frac{1}{\epsilon} \sigma - \nabla u &= 0 \\
 \tilde{\beta} \cdot \nabla_{xt} u - \nabla \cdot \sigma &= f \\
 \beta_n u - \epsilon \frac{\partial u}{\partial n} &= 0 \text{ on } \Gamma_- \\
 u &= 0 \text{ on } \Gamma_+ \\
 u &= u_0 \text{ on } \Gamma_0.
 \end{aligned}$$

The adjoint operator  $A^*$  is given by

$$A^*(v, \tau) = \left( \frac{1}{\epsilon} \tau + \nabla v, -\tilde{\beta} \cdot \nabla_{xt} v + \nabla \cdot \tau \right).$$

We decompose now the continuous adjoint problem

$$A^*(v, \tau) = (f, g)$$



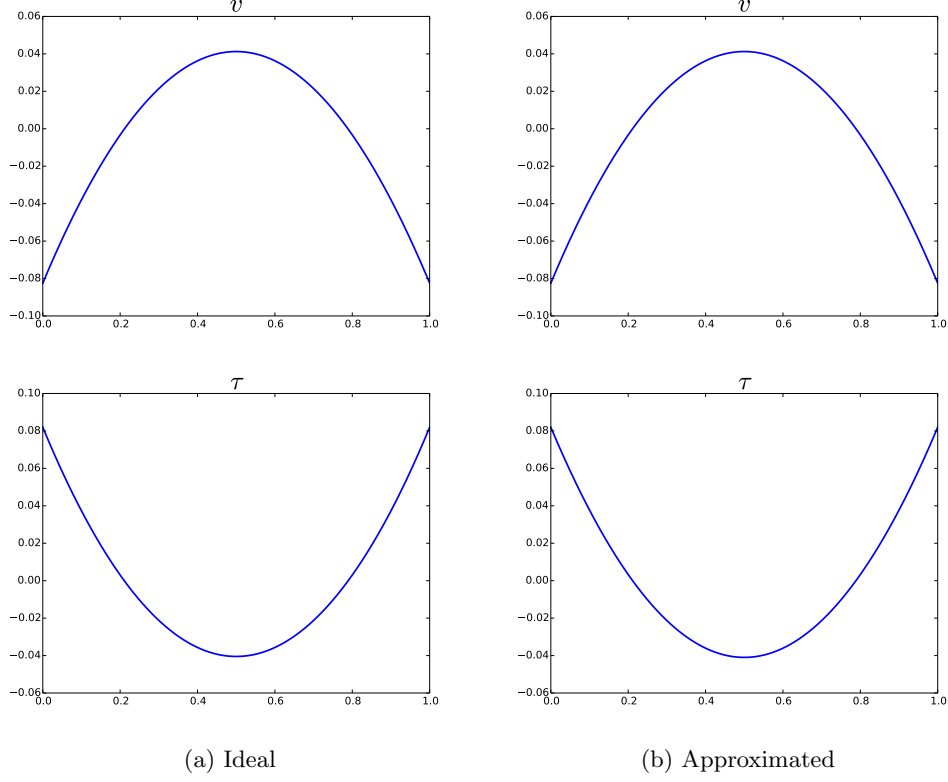


Figure 2: Robust norm optimal test functions for  $u = x - \frac{1}{2}$

into two cases a continuous part with forcing term  $g$

$$\begin{aligned}
 \frac{1}{\epsilon} \tau_1 + \nabla v_1 &= 0 \\
 -\tilde{\beta} \cdot \nabla_{xt} v_1 + \nabla \cdot \tau_1 &= g \\
 \tau_1 \cdot \mathbf{n}_x &= 0 \text{ on } \Gamma_- \\
 v_1 &= 0 \text{ on } \Gamma_+ \\
 v_1 &= 0 \text{ on } \Gamma_T,
 \end{aligned}$$

and a continuous part with forcing  $f$

$$\begin{aligned}
 \frac{1}{\epsilon} \tau_2 + \nabla v_2 &= f \\
 -\tilde{\beta} \cdot \nabla_{xt} v_2 + \nabla \cdot \tau_2 &= 0 \\
 \tau_2 \cdot \mathbf{n}_x &= 0 \text{ on } \Gamma_- \\
 v_2 &= 0 \text{ on } \Gamma_+ \\
 v_2 &= 0 \text{ on } \Gamma_T.
 \end{aligned}$$

(The boundary conditions can be derived by taking the ultra-weak formulation and choosing boundary conditions such that the temporal flux and spatial flux terms  $\langle \hat{u}, \llbracket \tau_n \rrbracket \rangle_{\Gamma_{out}}$  and  $\langle \hat{t}_n, \llbracket v \rrbracket \rangle_{\Gamma_{in}}$  are zero.)

We can then derive that the test norms

$$\|(v, \tau)\|_{V,K}^2 := \left\| \tilde{\beta} \cdot \nabla_{xt} v \right\|_K^2 + \epsilon \|\nabla v\|_K^2 + \|v\|_K^2 + \|\nabla \cdot \tau\|_K^2 + \frac{1}{\epsilon} \|\tau\|_K^2, \quad (7)$$

and

$$\|(v, \tau)\|_{V,K}^2 := \left\| \tilde{\beta} \cdot \nabla_{xt} v \right\|_K^2 + \epsilon \|\nabla v\|_K^2 + \|v\|_K^2 + \left\| \nabla \cdot \tau - \tilde{\beta} \cdot \nabla_{xt} v \right\|_K^2 + \frac{1}{\epsilon} \|\tau\|_K^2, \quad (8)$$

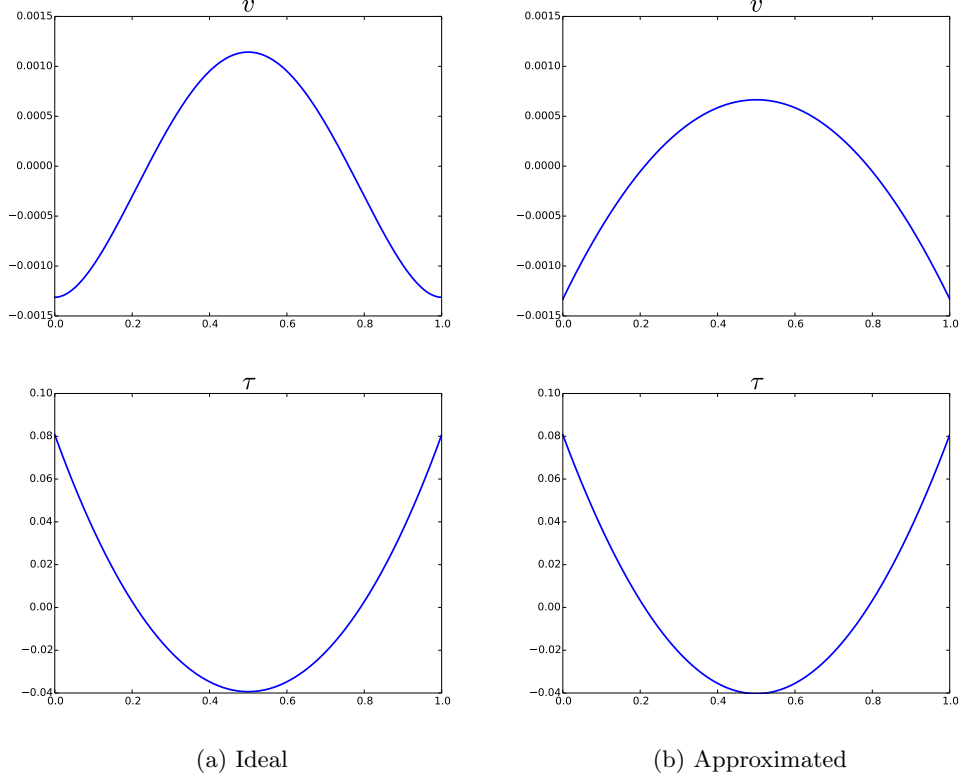


Figure 3: Coupled robust norm optimal test functions for  $u = x - \frac{1}{2}$

respectively designated the *robust* test norm and the *coupled robust* test norm, provide the necessary bound  $\|v^*\|_V \lesssim \|u\|_{L^2(Q)}$ .

In the following lemmas we establish the following bounds:

- Bound on  $\|(v_1, \tau_1)\|_V$ . Lemma 3.2 gives  $\|\tilde{\beta} \cdot \nabla_{xt} v_1\| \leq \|g\|$ . Since  $\nabla \cdot \tau_1 = g + \tilde{\beta} \cdot \nabla_{xt} v_1$ ,

$$\|\nabla \cdot \tau_1\| \leq \|g\| + \|\tilde{\beta} \cdot \nabla_{xt} v_1\| \leq 2\|g\|.$$

Or, the fact that  $\nabla \cdot \tau - \tilde{\beta} \cdot \nabla_{xt} v_1 = g$  clearly gives

$$\|\nabla \cdot \tau - \tilde{\beta} \cdot \nabla_{xt} v_1\| = \|g\|.$$

Lemma 3.1 gives  $\|v_1\|^2 + \epsilon \|\nabla v_1\|^2 \leq \|g\|^2$ . Since  $\epsilon^{1/2} \nabla v_1 = -\epsilon^{-1/2} \tau_1$ ,

$$\frac{1}{\epsilon} \|\tau_1\|^2 \leq \|g\|^2.$$

Thus, all  $(v_1, \tau_1)$  terms in (7) and (8) are accounted for, guaranteeing at least robust control of  $u$ .

- Bound on  $\|(v_2, \tau_2)\|_V$ . The fact that  $\nabla \cdot \tau - \tilde{\beta} \cdot \nabla_{xt} v = 0$  clearly gives

$$\|\nabla \cdot \tau - \tilde{\beta} \cdot \nabla_{xt} v_2\| = 0 \leq \|f\|.$$

Lemma 3.1 gives  $\|v_2\|^2 + \epsilon \|\nabla v_2\|^2 \leq \epsilon \|f\|^2$ . Since  $\epsilon^{1/2} \nabla v_2 = f - \epsilon^{-1/2} \tau_2$ ,

$$\frac{1}{\epsilon} \|\tau_2\|^2 \leq (1 + \epsilon) \|f\|^2.$$

We have not been able to develop bounds on  $\|\tilde{\beta} \cdot \nabla_{xt} v_2\|$  and  $\|\nabla \cdot \tau\|$  which means that we can not guarantee robust control of  $\sigma$  with with provided test norms.

We proceed now with the technical estimates.

**Lemma 3.1.** *For the duration of this lemma, let  $v := v_1 + v_2$ . Assuming the advection field  $\beta$  is incompressible, i.e.  $\nabla \cdot \beta = 0$ ,*

$$\|v\|^2 + \epsilon \|\nabla v\|^2 \leq \|g\|^2 + \epsilon \|\mathbf{f}\|^2.$$

*Proof.* Define  $w = e^t v$  and note that  $\frac{\partial w}{\partial t} = \left(\frac{\partial v}{\partial t} + v\right) e^t$  while all spatial derivatives go through. Multiplying the adjoint by  $w$  and integrating over  $Q$  gives

$$-\int_Q \tilde{\beta} \cdot \nabla_{xt} v w - \epsilon \Delta v w = \int_Q g w - \epsilon \int_Q \nabla \cdot \mathbf{f} w$$

or

$$-\int_Q e^t v \tilde{\beta} \cdot \nabla_{xt} v - \epsilon \int_Q e^t v \Delta v = \int_Q e^t g v - \epsilon \int_Q e^t v \nabla \cdot \mathbf{f}$$

Integrating by parts:

$$\begin{aligned} \int_Q \nabla_{xt} \cdot (e^t \tilde{\beta} v) v - \int_{\Gamma} e^t \tilde{\beta} \cdot \mathbf{n} v^2 + \epsilon \int_Q e^t \nabla v \cdot \nabla v - \epsilon \int_{\Gamma_x} e^t v \cdot \nabla v \cdot \mathbf{n}_x \\ = \int_Q e^t g v + \epsilon \int_Q e^t \nabla v \cdot \mathbf{f} - \epsilon \int_{\Gamma_x} e^t v \mathbf{f} \cdot \mathbf{n}_x \end{aligned}$$

Note that  $\nabla_{xt} \cdot e^t v \tilde{\beta} = e^t (\tilde{\beta} \cdot \nabla_{xt} v + v)$  if  $\nabla \cdot \beta = 0$ . Moving some terms to the right hand side, we get

$$\begin{aligned} \int_Q e^t v^2 + \int_Q \epsilon e^t \nabla v \cdot \nabla v \\ = \int_Q e^t g v + \epsilon \int_Q e^t \nabla v \cdot \mathbf{f} - \epsilon \int_{\Gamma_x} e^t v \mathbf{f} \cdot \mathbf{n}_x \\ - \int_Q e^t \tilde{\beta} \cdot \nabla_{xt} v v + \int_{\Gamma} e^t \tilde{\beta} \cdot \mathbf{n} v^2 + \epsilon \int_{\Gamma_x} e^t v \cdot \nabla v \cdot \mathbf{n}_x \end{aligned}$$

Note that  $1 \leq \|e^t\|_\infty = e^T$ . Then

$$\begin{aligned}
& \|v\|^2 + \epsilon \|\nabla v\|^2 \\
& \leq e^T \left( \int_Q gv + \epsilon \int_Q \nabla v \cdot \mathbf{f} - \epsilon \int_{\Gamma_-} v \underbrace{\mathbf{f} \cdot \mathbf{n}_x}_{=\tau_n + \frac{0}{\partial \mathbf{n}_x}} - \epsilon \int_{\Gamma_+} \underbrace{v}_{=0} \mathbf{f} \cdot \mathbf{n}_x \right. \\
& \quad \left. - \int_Q \tilde{\beta} \cdot \nabla_{xt} v v + \int_\Gamma \tilde{\beta} \cdot \mathbf{n} v^2 + \epsilon \int_{\Gamma_-} v \cdot \nabla v \cdot \mathbf{n}_x + \epsilon \int_{\Gamma_+} \underbrace{v}_{=0} \frac{\partial v}{\partial \mathbf{n}_x} \right) \\
& \quad \text{Note: boundary conditions give } \tau_n = 0 \text{ on } \Gamma_- \text{ and } v = 0 \text{ on } \Gamma_+ \\
& = e^T \left( \int_Q gv + \epsilon \int_Q \nabla v \cdot \mathbf{f} - \epsilon \int_{\Gamma_-} v \frac{\partial v}{\partial \mathbf{n}_x} + \epsilon \int_{\Gamma_x} v \frac{\partial v}{\partial \mathbf{n}_x} - \frac{1}{2} \int_Q \tilde{\beta} \cdot \nabla_{xt} v^2 + \int_\Gamma \tilde{\beta} \cdot \mathbf{n} v^2 \right) \\
& \quad \text{Note: } \Gamma_x = \Gamma_- \cup \Gamma_+ \text{ and } v = 0 \text{ on } \Gamma_- \\
& = e^T \left( \int_Q gv + \epsilon \int_Q \nabla v \cdot \mathbf{f} + \frac{1}{2} \int_Q \cancel{\nabla_{xt} \cdot \tilde{\beta} v^2}^0 - \frac{1}{2} \int_\Gamma \tilde{\beta} \cdot \mathbf{n} v^2 + \int_\Gamma \tilde{\beta} \cdot \mathbf{n} v^2 \right) \\
& \quad \text{Note: Integration by parts of } -\frac{1}{2} \int_Q \tilde{\beta} \cdot \nabla_{xt} v^2 \text{ and } \nabla \cdot \beta = 0 \\
& = e^T \left( \int_Q gv + \epsilon \int_Q \nabla v \cdot \mathbf{f} + \frac{1}{2} \left( \int_{\Gamma_0} \underbrace{-v^2}_{\leq 0} + \int_{\Gamma_T} \cancel{v^2}^0 + \int_{\Gamma_-} \underbrace{\beta \cdot \mathbf{n}_x v^2}_{\leq 0} + \int_{\Gamma_+} \beta \cdot \mathbf{n}_x \cancel{v^2}^0 \right) \right) \\
& \quad \text{Note: Split boundary term into components, } v = 0 \text{ on } \Gamma_+ \text{ and } \Gamma_T \\
& \leq e^T \left( \int_Q gv + \epsilon \int_Q \nabla v \cdot \mathbf{f} \right) \\
& \leq e^T \left( \frac{\|g\|^2}{2} + \epsilon \frac{\|\mathbf{f}\|^2}{2} + \frac{\|v\|^2}{2} + \epsilon \frac{\|\nabla v\|^2}{2} \right)
\end{aligned}$$

Note: Young's inequality

□

**Lemma 3.2.** If  $\|\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}\|_{L^\infty} \leq C_\beta$  we can bound

$$\|\tilde{\beta} \cdot \nabla_{xt} v_1\| \lesssim \|g\|.$$

*Proof.* Multiply  $-\tilde{\beta} \cdot \nabla_{xt} v_1 = g - \nabla \cdot \tau_1$  by  $-\tilde{\beta} \cdot \nabla_{xt} v_1$  and integrate over  $Q$  to get

$$\|\tilde{\beta} \cdot \nabla_{xt} v_1\|^2 = - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \int_Q \tilde{\beta} \cdot \nabla_{xt} v_1 \nabla \cdot \tau_1. \quad (9)$$

Note that

$$\frac{1}{\epsilon} \int_Q \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla \cdot \boldsymbol{\tau}_1 = - \int_Q \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla \cdot \nabla v_1$$

$$\text{Note: } \boldsymbol{\tau}_1 = \epsilon \nabla v_1$$

$$= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \int_Q \nabla (\tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1) \cdot \nabla v_1$$

$$\text{Note: Integration by parts}$$

$$\begin{aligned} &= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \int_Q (\nabla \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1) \cdot \nabla v_1 + \int_Q \tilde{\boldsymbol{\beta}} \cdot \nabla \nabla_{xt} v_1 \cdot \nabla v_1 \\ &= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \int_Q (\nabla \boldsymbol{\beta} \cdot \nabla v_1) \cdot \nabla v_1 + \frac{1}{2} \int_Q \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} (\nabla v_1 \cdot \nabla v_1) \end{aligned}$$

$$\text{Note: } \nabla \nabla_{xt} v_1 \cdot \nabla v_1 = \nabla_{xt} \nabla v_1 \cdot \nabla v_1 = \frac{1}{2} \nabla_{xt} (\nabla v_1 \cdot \nabla v_1)$$

$$= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \int_Q (\nabla \boldsymbol{\beta} \cdot \nabla v_1) \cdot \nabla v_1 + \frac{1}{2} \int_{\Gamma} \tilde{\boldsymbol{\beta}} \cdot \mathbf{n} (\nabla v_1 \cdot \nabla v_1) - \frac{1}{2} \int_Q \nabla_{xt} \cdot \tilde{\boldsymbol{\beta}} (\nabla v_1 \cdot \nabla v_1)$$

$$\text{Note: Integration by parts}$$

$$= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \int_Q (\nabla \boldsymbol{\beta} \cdot \nabla v_1) \cdot \nabla v_1 + \frac{1}{2} \int_{\Gamma} \tilde{\boldsymbol{\beta}} \cdot \mathbf{n} (\nabla v_1 \cdot \nabla v_1) - \frac{1}{2} \int_Q \nabla \cdot \boldsymbol{\beta} (\nabla v_1 \cdot \nabla v_1)$$

$$\text{Note: } \nabla_{xt} \cdot \tilde{\boldsymbol{\beta}} = \nabla \cdot \boldsymbol{\beta}$$

$$= - \int_{\Gamma_x} \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \frac{1}{2} \int_{\Gamma} \tilde{\boldsymbol{\beta}} \cdot \mathbf{n} (\nabla v_1 \cdot \nabla v_1) + \int_Q \nabla v_1 (\nabla \boldsymbol{\beta} - \frac{1}{2} \nabla \cdot \boldsymbol{\beta} \mathbf{I}) \nabla v_1$$

$$\text{Note: } (\nabla \boldsymbol{\beta} \cdot \nabla v_1) \cdot \nabla v_1 - \frac{1}{2} \nabla \cdot \boldsymbol{\beta} (\nabla v_1 \cdot \nabla v_1) = \nabla v_1 (\nabla \boldsymbol{\beta} - \frac{1}{2} \nabla \cdot \boldsymbol{\beta} \mathbf{I}) \nabla v_1$$

Plugging this into (9), we get

$$\begin{aligned}
\left\| \tilde{\beta} \cdot \nabla_{xt} v_1 \right\|^2 &= - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\quad - \epsilon \int_{\Gamma_x} \tilde{\beta} \cdot \nabla_{xt} v_1 \nabla v_1 \cdot \mathbf{n}_x + \frac{\epsilon}{2} \int_{\Gamma} \tilde{\beta} \cdot \mathbf{n} (\nabla v_1 \cdot \nabla v_1) \\
&= - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\quad - \epsilon \int_{\Gamma_-} \tilde{\beta} \cdot \nabla_{xt} v_1 \underbrace{\nabla v_1 \cdot \mathbf{n}_x}_{=0} - \epsilon \int_{\Gamma_+} \left( \underbrace{\frac{\partial v_1}{\partial t}}_{=0} + \beta \cdot \nabla v_1 \right) \nabla v_1 \cdot \mathbf{n}_x
\end{aligned}$$

Note:  $\nabla v_1 \cdot \mathbf{n}_x = \tau_{1n} = 0$  on  $\Gamma_-$ ,  $v_1 = 0$  on  $\Gamma_+$

$$\begin{aligned}
&+ \frac{\epsilon}{2} \int_{\Gamma_-} \underbrace{\beta \cdot \mathbf{n}_x}_{<0} (\nabla v_1 \cdot \nabla v_1) + \frac{\epsilon}{2} \int_{\Gamma_+} \beta \cdot \mathbf{n}_x (\nabla v_1 \cdot \nabla v_1) \\
&+ \frac{\epsilon}{2} \int_{\Gamma_0} \underbrace{n_t}_{<0} (\nabla v_1 \cdot \nabla v_1) + \frac{\epsilon}{2} \int_{\Gamma_T} n_t \underbrace{(\nabla v_1 \cdot \nabla v_1)}_{=0}
\end{aligned}$$

Note:  $v_1 = 0$  on  $\Gamma_T$

$$\begin{aligned}
&\leq - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\quad + \epsilon \int_{\Gamma_+} \left( - \frac{\partial v_1}{\partial \mathbf{n}_x} \beta + \frac{1}{2} \beta \cdot \mathbf{n}_x \nabla v_1 \right) \cdot \nabla v_1
\end{aligned}$$

Note: Dropped negative terms from RHS

$$\begin{aligned}
&= - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\quad + \epsilon \int_{\Gamma_+} \left( - \frac{\partial v_1}{\partial \mathbf{n}_x} \beta + \frac{1}{2} \beta \cdot \mathbf{n}_x \frac{\partial v_1}{\partial \mathbf{n}_x} \mathbf{n}_x \right) \cdot \frac{\partial v_1}{\partial \mathbf{n}_x} \mathbf{n}_x
\end{aligned}$$

Note:  $\nabla v_1 \cdot \nabla v_1 = \nabla v_1 \cdot \nabla v_1 \mathbf{n}_x \cdot \mathbf{n}_x = (\nabla v_1 \cdot \mathbf{n}_x \mathbf{n}_x) \cdot (\nabla v_1 \cdot \mathbf{n}_x \mathbf{n}_x)$

$$\begin{aligned}
&= - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\quad - \underbrace{\frac{\epsilon}{2} \int_{\Gamma_+} \left( \frac{\partial v_1}{\partial \mathbf{n}_x} \right)^2 \beta \cdot \mathbf{n}_x}_{<0}
\end{aligned}$$

$$\begin{aligned}
&\leq - \int_Q g \tilde{\beta} \cdot \nabla_{xt} v_1 + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1 \\
&\leq \frac{\|g\|^2}{2} + \frac{\left\| \tilde{\beta} \cdot \nabla_{xt} v_1 \right\|^2}{2} + \epsilon \int_Q \nabla v_1 (\nabla \beta - \frac{1}{2} \nabla \cdot \beta \mathbf{I}) \nabla v_1
\end{aligned}$$

Note: Young's inequality

$$\leq \frac{\|g\|^2}{2} + \frac{\left\| \tilde{\beta} \cdot \nabla_{xt} v_1 \right\|^2}{2} + \epsilon C_\beta \|\nabla v_1\|^2$$

Note: Assumption on  $\beta$

$$\leq \left( \frac{1}{2} + C_\beta \right) \|g\|^2 + \frac{\left\| \tilde{\beta} \cdot \nabla_{xt} v_1 \right\|^2}{2}$$

□

In conclusion, with either robust test norm, we can claim the following stability result,

$$\begin{aligned}\|u - u_h\| &\lesssim \|(u, \boldsymbol{\sigma}, \hat{u}, \hat{t}) - (u_h, \boldsymbol{\sigma}_h, \hat{u}_h, \hat{t}_h)\|_E \\ &= \inf_{(u_h, \boldsymbol{\sigma}_h, \hat{u}_h, \hat{t}_h)} \|(u, \boldsymbol{\sigma}, \hat{u}, \hat{t}) - (u_h, \boldsymbol{\sigma}_h, \hat{u}_h, \hat{t}_h)\|_E.\end{aligned}$$

Notice that, contrary to the steady-state case, we have not been able to secure a robust  $L^2$  bound for the stress. The best approximation error in the energy norm can be estimated locally, i.e. element-wise, see [15, 4]. This leads to an ultimate, final  $h$  estimate but not necessarily with robust constants. The loss of robustness in the best approximation error estimate is the consequence of rescaling the  $L^2$ -terms to avoid boundary layers in the optimal test functions. However, similarly to the steady-state case, with refinements, the mesh-dependent  $L^2$ -terms converge to the optimal ones so we hope to regain robustness in the limit. We do not attempt to analyze the best approximation error in this contribution and restrict ourselves to numerical experiments only.

## 4 Numerical Tests

The norms given in (7) and (8) are robust, but the reaction (0<sup>th</sup> order) terms still dominate the diffusion terms which produces boundary layers in optimal test functions and prohibits their resolution with a simple enrichment strategy. We can mitigate this by introducing mesh-dependent norms:

$$\|(v, \boldsymbol{\tau})\|_{V,K}^2 := \|\tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v\|_K^2 + \epsilon \|\nabla v\|_K^2 + \min\left(\frac{\epsilon}{h^2}, 1\right) \|v\|_K^2 + \|\nabla \cdot \boldsymbol{\tau}\|_K^2 + \min\left(\frac{1}{\epsilon}, \frac{1}{h^2}\right) \|\boldsymbol{\tau}\|_K^2, \quad (10)$$

and

$$\|(v, \boldsymbol{\tau})\|_{V,K}^2 := \|\tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v\|_K^2 + \epsilon \|\nabla v\|_K^2 + \min\left(\frac{\epsilon}{h^2}, 1\right) \|v\|_K^2 + \|\nabla \cdot \boldsymbol{\tau} - \tilde{\boldsymbol{\beta}} \cdot \nabla_{xt} v\|_K^2 + \min\left(\frac{1}{\epsilon}, \frac{1}{h^2}\right) \|\boldsymbol{\tau}\|_K^2. \quad (11)$$

Note that any version of (7) and (8) with smaller coefficients also satisfies the criteria for robustness. The mesh dependent coefficients were chosen in an attempt to balance the relative size of “reaction” terms like  $\|v\|$  which scale like  $h^d$  with “diffusive” terms like  $\epsilon \|\nabla v\|$  which scale like  $h^{d-2}$ . This is also the mechanism by which we avoid creating sharp boundary layers in our optimal test functions – by correctly balancing reactive and diffusive terms. In the following numerical experiments, we compute with these mesh dependent norms.

We verify robust convergence of our transient coupled robust norm on an analytical solution (shown in Figure 4) that decays to a steady state Eriksson-Johnson problem:

$$u = \exp(-lt) [\exp(\lambda_1 x) - \exp(\lambda_2 x)] + \cos(\pi y) \frac{\exp(s_1 x) - \exp(r_1 x)}{\exp(-s_1) - \exp(-r_1)},$$

where  $l = 4$ ,  $\lambda_{1,2} = \frac{-1 \pm \sqrt{1-4\epsilon l}}{-2\epsilon}$ ,  $r_1 = \frac{1 + \sqrt{1+4\pi^2 \epsilon^2}}{2\epsilon}$ , and  $s_1 = \frac{1 - \sqrt{1+4\pi^2 \epsilon^2}}{2\epsilon}$ . The problem domain is  $[-1, 0] \times [-0.5, 0.5]$  and  $\boldsymbol{\beta} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ . We show robustness for  $\epsilon = 10^{-2}, 10^{-4}, 10^{-6}, 10^{-8}$  for linear, quadratic, and

quartic polynomial trial functions. Flux boundary conditions were applied based on the exact solution at  $x = -1$  and  $t = 0$  while trace boundary conditions were set at  $y = -0.5$ ,  $y = 0.5$ , and  $x = 0$ . An adaptive solve was undertaken using a greedy refinement strategy in which any elements with at least 20% of the energy error of highest energy error element were refined at each step. See [14] for details on adaptivity within the DPG context.

In the plot legends,  $L^2$  indicates  $\left(\|u - u_{\text{exact}}\|_L^2 + \|\boldsymbol{\sigma} - \boldsymbol{\sigma}_{\text{exact}}\|_{L^2}\right)^{\frac{1}{2}}$  while  $V^*$  indicates the energy error reported by the method. Despite a lack of guaranteed control  $\boldsymbol{\sigma}$  by norms (10) and (11),  $\|\boldsymbol{\sigma} - \boldsymbol{\sigma}_{\text{exact}}\|_{L^2}$  is included in the  $L^2$  error computation and does appear to be under control in the problems considered here. When plotted in isolation, the  $L^2$  error in  $\boldsymbol{\sigma}$  was usually orders of magnitude smaller than  $\|u - u_{\text{exact}}\|_{L^2}$ .

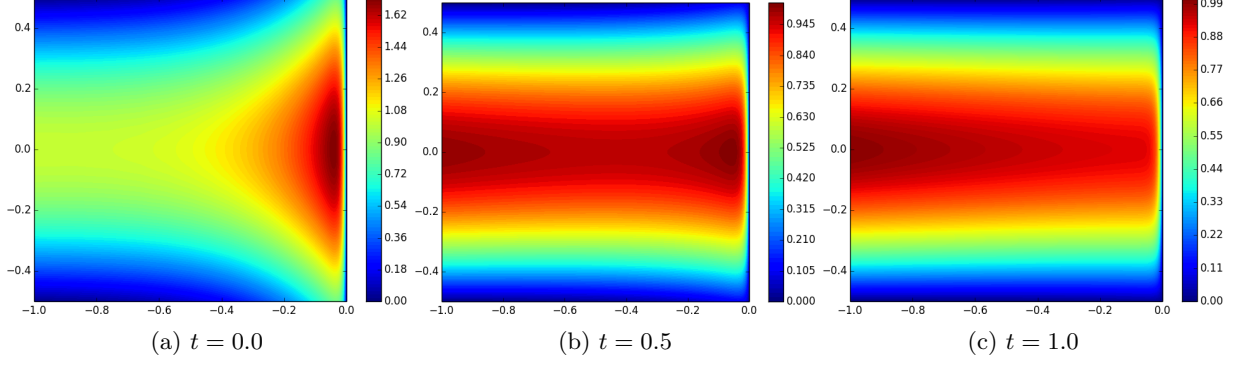


Figure 4: Transient Eriksson-Johnson solution

We provide surface plots of temporal slices of the solution at  $t = 0.2$  for the two norms with  $\epsilon = 10^{-2}$ , and  $p = 2$  after 4 adaptive refinements. The results conform to our previous experience with steady convection-diffusion where the coupled robust norm tends to produce smoother results in regions with sharp gradients.

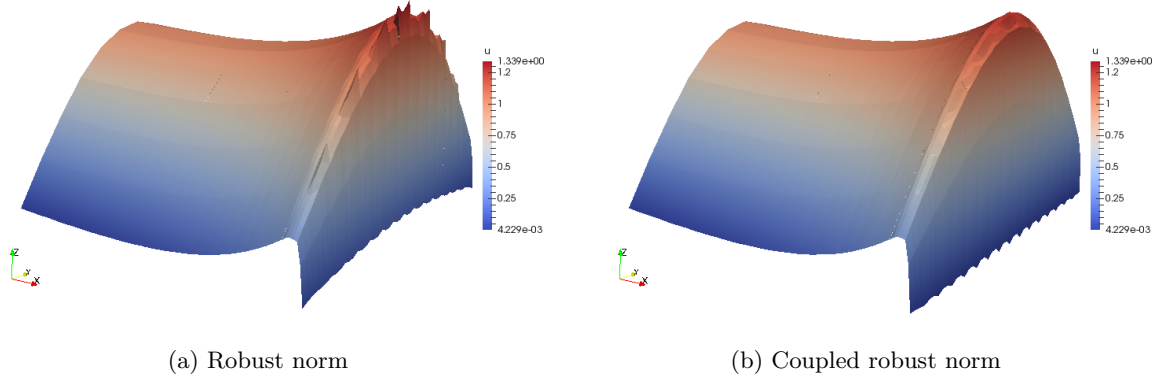


Figure 5:  $u$  at  $t = 0.2$  for  $\epsilon = 10^{-2}$  and  $p = 2$  after 4 adaptive refinements

## 5 Conclusions

As expected, convergence of the energy error appears to be a reliable predictor of convergence of the  $L^2$  error. This relation is especially tight for moderate values of  $\epsilon$ . We've developed two robust test norms for transient convection-diffusion, though neither one guarantees robust control over  $\sigma$  as we had with their steady analogs.

## References

- [1] S.K. Aliabadi and T.E. Tezduyar. Space-time finite element computation of compressible flows involving moving boundaries and interfaces. *Comput. Methods in Appl. Mech. Eng.*, 107(12):209 – 223, 1993.
- [2] J.H. Argyris and D.W. Scharpf. Finite elements in time and space. *Nucl. Eng. Des.*, 10(4):456 – 464, 1969.



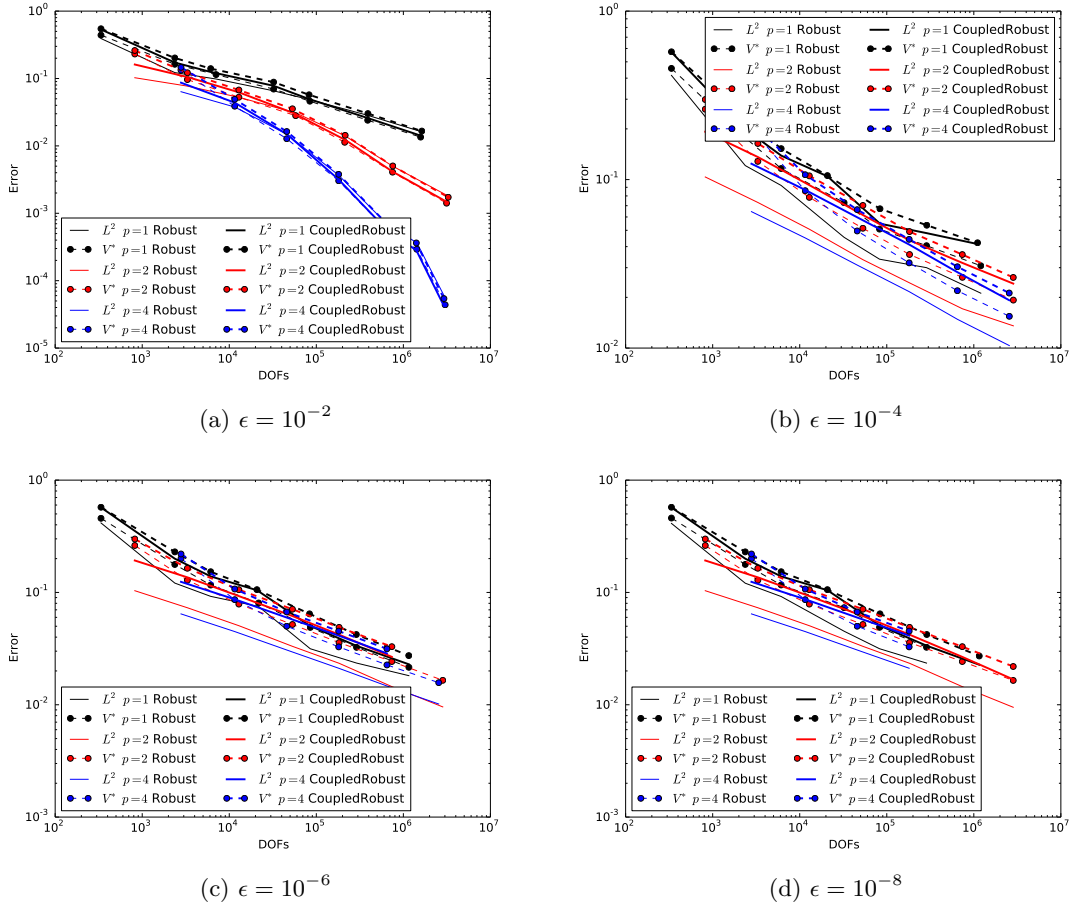


Figure 6: Convergence to analytical solution

- [3] C. Carstensen and L.F. Demkowicz J. Gopalakrishnan. Breaking spaces and forms for the DPG method and applications including Maxwell equations. Technical Report 15-18, ICES, 2015.
- [4] J. Chan, N. Heuer, T. Bui-Thanh, and L. Demkowicz. A robust DPG method for convection-dominated diffusion problems II: Adjoint boundary conditions and mesh-dependent test norms. *Comp. Math. Appl.*, 67(4):771 – 795, 2014. High-order Finite Element Approximation for Partial Differential Equations.
- [5] J.L. Chan. *A DPG Method for Convection-Diffusion Problems*. PhD thesis, University of Texas at Austin, 2013.
- [6] J.L. Chan, J.A. Evans, and W. Qiu. A dual Petrov-Galerkin finite element method for the convection-diffusion equation. *Comput. Math. Appl.*, 68(11):1513 – 1529, 2014. Minimum Residual and Least Squares Finite Element Methods.
- [7] J.L. Chan, J. Gopalakrishnan, and L.F. Demkowicz. Global properties of DPG test spaces for convection-diffusion problems. Technical Report 13-05, ICES, February 2013.
- [8] A. Cohen, W. Dahmen, and G. Welper. Adaptivity and variational stabilization for convection-diffusion equations. *ESAIM: Math. Model. Numer. Anal.*, 46(5):1247–1273, 3 2012.
- [9] W. Dahmen, C. Huang, C. Schwab, and G. Welper. Adaptive Petrov-Galerkin methods for first order transport equations. *SIAM J. Numer. Anal.*, 50(5):2420–2445, 2012.

- [10] L.F. Demkowicz. Various variational formulations and closed range theorem. Technical Report 15-03, ICES, January 2015.
- [11] L.F. Demkowicz and J. Gopalakrishnan. A class of discontinuous Petrov-Galerkin methods. Part II: Optimal test functions. *Numer. Meth. Part. D. E.*, 2010.
- [12] L.F. Demkowicz and J. Gopalakrishnan. *Recent Developments in Discontinuous Galerkin Finite Element Methods for Partial Differential Equations* (eds. X. Feng, O. Karakashian, Y. Xing), volume 157, chapter An Overview of the DPG Method, pages 149–180. IMA Volumes in Mathematics and its Applications, 2014.
- [13] L.F. Demkowicz and J. Gopalakrishnan. Discontinuous Petrov-Galerkin (DPG) method. Technical Report 15-20, ICES, December 2015.
- [14] L.F. Demkowicz, J. Gopalakrishnan, and A.H. Niemi. A class of discontinuous Petrov-Galerkin methods. Part III: Adaptivity. *Appl. Numer. Math.*, 62(4):396–427, April 2012.
- [15] L.F. Demkowicz and N. Heuer. Robust DPG method for convection-dominated diffusion problems. *SIAM J. Numer. Anal.*, 51(5):1514–2537, 2013.
- [16] T.E. Ellis, L.F. Demkowicz, and J.L. Chan. Locally conservative discontinuous Petrov-Galerkin finite elements for fluid problems. *Comp. Math. Appl.*, 68(11):1530 – 1549, 2014.
- [17] I. Fried. Finite-element analysis of time-dependent phenomena. *AIAA J.*, 7(6):1170–1173, 1969.
- [18] J. Gopalakrishnan and W. Qiu. An analysis of the practical DPG method. *Math. Comp.*, 83(286):537–552, March 2014.
- [19] T.J.R. Hughes and J.R. Stewart. A space-time formulation for multiscale phenomena. *J. Comput. Appl. Math.*, 74(12):217 – 229, 1996.
- [20] Z. Kaczkowski. The method of finite space-time elements in dynamics of structures. *J. Tech. Phys.*, 16(1):69 – 84, 1975.
- [21] C.M. Klaij, J.J.W. van der Vegt, and H. van der Ven. Space-time discontinuous Galerkin method for the compressible Navier-Stokes equations. *J. Comp. Phys.*, 217(2):589 – 611, 2006.
- [22] D. Moro, N.C. Nguyen, and J. Peraire. A hybridized discontinuous Petrov-Galerkin scheme for compressible flows. Master’s thesis, Massachusetts Institute of Technology, 2011.
- [23] A.H. Niemi, N.O. Collier, and V.M. Calo. Automatically stable discontinuous Petrov-Galerkin methods for stationary transport problems: Quasi-optimal test space norm. *Comput. Math. Appl.*, 66(10):2096–2113, December 2013.
- [24] A.H. Niemi, N.O. Collier, and V.M. Calo. Discontinuous Petrov-Galerkin method based on the optimal test space norm for steady transport problems in one space dimension. *Journal of Computational Science*, 4(3):157–163, 2013.
- [25] J.T. Oden. A general theory of finite elements. II. Applications. *Int. J. Numer. Meth. Eng.*, 1(3):247–259, 1969.
- [26] S. Rhebergen and B. Cockburn. A space-time hybridizable discontinuous Galerkin method for incompressible flows on deforming domains. *J. Comp. Phys.*, 231(11):4185 – 4204, 2012.
- [27] S. Rhebergen, B. Cockburn, and J.J.W. Van Der Vegt. A space-time discontinuous Galerkin method for the incompressible Navier-Stokes equations. *J. Comput. Phys.*, 233:339–358, January 2013.
- [28] H. G. Roos, M. Stynes, and L. Tobiska. *Robust Numerical Methods for Singularly Perturbed Differential Equations*, volume 24 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2nd edition, 2008.

- [29] T.E. Tezduyar, M. Behr, and J. Liou. A new strategy for finite element computations involving moving boundaries and interfaces – The deforming-spatial-domain/space-time procedure: I. The concept and the preliminary numerical tests. *Comput. Methods in Appl. Mech. Eng.*, 94(3):339 – 351, 1992.
- [30] J.J.W. van der Vegt and H. van der Ven. Space-time discontinuous Galerkin finite element method with dynamic grid motion for inviscid compressible flows: I. General formulation. *J. Comp. Phys.*, 182(2):546 – 585, 2002.