



Discontinuous Galerkin Methods for a Class of Nonvariational Problems

Andreas Dedner¹ · Tristan Pryer²

Received: 22 September 2020 / Revised: 1 March 2021 / Accepted: 22 March 2021
© The Author(s) 2021

Abstract

We extend the finite element method introduced by Lakkis and Pryer (SIAM J. Sci. Comput. 33(2): 786–801, 2011) to approximate the solution of second-order elliptic problems in nonvariational form to incorporate the discontinuous Galerkin (DG) framework. This is done by viewing the “finite element Hessian” as an auxiliary variable in the formulation. Representing the finite element Hessian in a discontinuous setting yields a linear system of the same size and having the same sparsity pattern of the compact DG methods for variational elliptic problems. Furthermore, the system matrix is very easy to assemble; thus, this approach greatly reduces the computational complexity of the discretisation compared to the continuous approach. We conduct a stability and consistency analysis making use of the unified framework set out in Arnold et al. (SIAM J. Numer. Anal. 39(5): 1749–1779, 2001/2002). We also give an a posteriori analysis of the method in the case where the problem has a strong solution. The analysis applies to any consistent representation of the finite element Hessian, and thus is applicable to the previous works making use of continuous Galerkin approximations. Numerical evidence is presented showing that the method works well also in a more general setting.

Keywords Nonvariational problems · Discontinuous Galerkin · Error estimates · Adaptivity

Mathematics Subject Classification 65N15 · 65N30 · 65Y20 · 65N12

1 Introduction

Linear, second-order, nonvariational partial differential equations (PDEs) are those which are given in the form

✉ Andreas Dedner
A.S.Dedner@warwick.ac.uk

Tristan Pryer
tmp38@bath.ac.uk

¹ Mathematics Institute, University of Warwick, Coventry CV4 7AL, UK

² Department of Mathematical Sciences, University of Bath, Bath BA2 7AY, UK

$$-A : D^2u = f, \tag{1}$$

where $X:Y = \text{trace}(X^T Y)$ is the Frobenius inner product between matrices. If the matrix A is differentiable, then there is an equivalence between this problem and its variational sibling

$$-\text{div}(A \nabla u) + DA \nabla u = f, \tag{2}$$

where

$$DA = \left(\sum_{i=1}^d \partial_i a_{i,1}(x), \dots, \sum_{i=1}^d \partial_i a_{i,d}(x) \right), \tag{3}$$

and d is the dimension under consideration. Rewriting in this form is sometimes undesirable. For example, if the coefficient matrix A has near singular derivatives, the problem will become advection dominated and possibly unstable for conforming finite element methods. There is a wealth of material on the treatment of advection-dominated problems [c.f., 18, 19]. If A is not differentiable, then the problem has no variational structure. In this case, standard finite element methods cannot be applied.

In a previous work [32], a finite element method for the approximation of the nonvariational problem (1) was introduced. This involved the introduction of a *finite element Hessian* represented in the same finite element space as the solution (modulo boundary conditions). The applications of the discrete representation of a Hessian of a piecewise function are becoming broader, for example, it can be used to drive anisotropic adaptive algorithms [1, 40], as a notion of discrete convexity [2] and in the design of finite element methods for nonlinear fourth-order problems [37].

The algebraic formulation of the C^0 Galerkin approximation of the nonvariational problem requires the solution of large sparse $(d + 1)^2 N^2$ linear system [32, Lem 3.3], where N is the number of degrees of freedom. Equivalently, using a Schur complement argument, this can be reduced to an N^2 full linear system. The reason that this system is full is due to the global nature of the $L^2(\Omega)$ projection operator into a continuous finite element space. The motivation for extending the nonvariational finite element method into the discontinuous setting is the massive gain in computational efficiency over the continuous case. Indeed, due to the local representation of the projection operators in these discontinuous spaces, we are able to make massive computational savings, in that the system matrix becomes sparse and is the same size as that of a standard discontinuous Galerkin stiffness matrix for the Laplacian.

There has been a plethora of other finite element methods posed for linear nonvariational PDEs exist including [38] in which the authors pose a stable discontinuous Galerkin scheme for linear nonvariational problems using discrete analogs of the classical methods used to prove existence of strong solutions, see also [29] where similar techniques are used over curved domains. In [22–24] the author considers a least squares formulation, see also [31, 35] for related formulations. Many of these formulations share the same disadvantage in the conditioning of the system matrix—the condition number scales like a fourth-order problem. As already mentioned, one of the desirable features of the formulation presented here is that it scales in the same fashion as a discontinuous Galerkin method for a second-order variational problem, a property shared by the method in [25] under the Cordes condition.

We are particularly interested in nonvariational PDEs due to their relation to general fully nonlinear PDEs

$$\mathcal{F}(D^2u) = 0, \quad (4)$$

which are of significant current research. In the literature finite element methods have been presented to solve this general class of problem. For example in [8] the author presents a C^1 finite element method shows stability and consistency (hence convergence) of the scheme which requires a high degree of smoothness on the exact solution. In [20, 21] the authors give a method in which they approximate the general second-order fully nonlinear PDE by a sequence of fourth-order quasilinear PDEs. This is reminiscent of the vanishing viscosity method introduced for classically studying first-order fully nonlinear PDEs. Efficiency of any method used to approximate a problem such as this is key. Each of the methods is computationally costly due to their reliance on C^1 finite elements [8, 21] or mixed methods [20].

In [5], a generic framework was set up to prove convergence of numerical approximations to the viscosity solutions of degenerate elliptic fully nonlinear PDEs. This involved constructing monotone sequences of approximations which are typically applied to finite difference approximations of the nonlinear problem [c.f., 36]. The assumption of consistency made in the [5] framework is incompatible with finite element methods; however, an extremely important observation made in [28] is that the consistency condition may be weakened to incorporate the finite element case using a localisation argument (in the case of isotropic diffusion).

A posteriori analysis of linear nonvariational problems is less standard than for those of variational type. Typically, results are based on “closeness” conditions of Cordes type [12] which guarantees existence of a unique smooth, H^2 solution. Under these assumptions, it is relatively straightforward to derive a posteriori bounds in the natural norm for the problem, H^2 , or a mesh-dependent equivalent. For other methods, this has been done in [24]. It has even been shown that adaptive methods for these problems converge optimally in terms of the number of degrees of freedom [30].

We show similar a posteriori bounds that make it straightforward to incorporate the method into well-developed software package for finite element methods, shown here using DUNE [6, 7]. In this work, we are interested in the asymptotic behaviour of the discontinuous approximation and the computational gains using this method. In a subsequent work, we will study the computational gains using the discontinuous framework presented over the continuous one given in [32], as well as exploit the powerful parallelisation capabilities of the package.

The rest of the paper is set out as follows. In Sect. 2, we formally introduce the model problem and give a brief review of known classical facts about nonvariational PDEs. In Sect. 3, we examine the discretisation of the nonvariational method in the discontinuous Galerkin framework, making use of the unified framework set out in [4] to derive a very general formulation of the finite element Hessian represented as a discontinuous object. We present some examples and examine the natural question of what happens when we try to eliminate the finite element Hessian from the formulation. In Sect. 4, we conduct an a posteriori analysis of the method and show upper and lower bounds to the error in an H^2 -like mesh-dependent norm. Finally, in Sect. 5, we detail a summary of extensive numerical experiments aimed at examining convergence and robustness of the method presented.

2 Problem Formulation

In this section, we formulate the model problem, fix notation and give some basic assumptions. In addition, we review the existence and uniqueness of the nonvariational problems. Let $\Omega \subset \mathbb{R}^d$, $d = 2, 3$, be a connected domain with polygonal boundary. The Lebesgue spaces are defined as

$$L^2(\Omega) = \left\{ \phi : \int_{\Omega} |\phi(\mathbf{x})|^2 \, d\mathbf{x} < \infty \right\} \quad \text{and} \quad L^\infty(\Omega) = \left\{ \phi : \sup_{\mathbf{x} \in \Omega} |\phi(\mathbf{x})| < \infty \right\}, \tag{5}$$

and the Sobolev and Hilbert spaces

$$W^{k,p}(\Omega) = \{ \phi \in L^p(\Omega) : D^\alpha \phi \in L^p(\Omega), \quad \text{for } |\alpha| \leq k \} \quad \text{and} \quad H^k(\Omega) := W^{k,2}(\Omega). \tag{6}$$

These are equipped with the norms

$$\|\phi\|_{L^2(\Omega)}^2 = \int_{\Omega} |\phi|^2 \, d\mathbf{x}, \quad \|\phi\|_{L^\infty(\Omega)} = \sup_{\mathbf{x} \in \Omega} |\phi(\mathbf{x})|, \tag{7}$$

$$\|v\|_{W^{k,p}(\Omega)}^p = \sum_{|\alpha| \leq k} \|D^\alpha v\|_{L^p(\Omega)}^p \quad \text{and} \quad |v|_{W^{k,p}(\Omega)}^p = \sum_{|\alpha|=k} \|D^\alpha v\|_{L^p(\Omega)}^p, \tag{8}$$

where $\alpha = \{\alpha_1, \dots, \alpha_d\}$ is a multi-index, $|\alpha| = \sum_{i=1}^d \alpha_i$ and derivatives D^α are understood in a weak sense. We pay particular attention to the cases $k = 1, 2$ and

$$H_0^1(\Omega) := \text{closure of } C_0^\infty(\Omega) \text{ in } H^1(\Omega). \tag{9}$$

The model problem in strong form is: find $u \in H^2(\Omega) \cap H_0^1(\Omega)$ such that

$$\langle \mathcal{L}u, \phi \rangle = \langle f, \phi \rangle, \quad \forall \phi \in H_0^1(\Omega), \tag{10}$$

where the data $f \in L^2(\Omega)$ are prescribed and \mathcal{L} is a general linear, second-order, uniformly elliptic partial differential operator. Let $A \in L^\infty(\Omega)^{d \times d}$, we then define

$$\begin{aligned} \mathcal{L} : H^2(\Omega) \cap H_0^1(\Omega) &\rightarrow L^2(\Omega), \\ u &\mapsto \mathcal{L}u := -A : D^2u. \end{aligned} \tag{11}$$

We assume that A is uniformly positive definite, i.e., there exists a $\gamma > 0$ such that for all \mathbf{x}

$$\mathbf{y}^T A(\mathbf{x}) \mathbf{y} \geq \gamma |\mathbf{y}|^2, \quad \forall \mathbf{y} \in \mathbb{R}^d, \tag{12}$$

and we call γ the *ellipticity constant*.

Definition 1 (Strong solution) A *strong solution* of (1) is a function $u \in H^2(\Omega) \cap H_0^1(\Omega)$, that is a twice weakly differentiable function, which satisfies the problem almost everywhere.

Theorem 1 (Existence and regularity of a strong solution of (1) [12]) *Let $\Omega \subset \mathbb{R}^d$ be a convex polytope. Suppose $A \in L^\infty(\Omega)^{d \times d}$ is uniformly elliptic and $f \in L^2(\Omega)$. Suppose further that A satisfies the Cordes condition, that there exists $0 < \alpha \in L^\infty(\Omega)$ and $0 < \beta, \gamma \in \mathbb{R}$ with $\beta + \gamma < 1$ such that*

$$|\mathbf{I} : \mathbf{Y} - \alpha(x)\mathbf{A} : \mathbf{Y}| \leq \beta|\mathbf{Y}| + \gamma|\mathbf{I} : \mathbf{Y}|, \quad \forall \mathbf{Y} \in \text{Sym}^+(\mathbb{R}^{d \times d}). \tag{13}$$

Then, (1) admits a unique strong solution. There also exists a constant independent of u such that

$$\|u\|_{H^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}. \tag{14}$$

Remark 1 (Less regular solutions) Note that the theory of viscosity solutions has been developed for non-classical solutions of (10), if the problem data do not satisfy the regularity assumed above, see [27].

Assumption 1 (Inf-sup condition) From hereon in, we will assume that the problem satisfies an inf-sup condition, that is, with

$$a(w, v) := \int_{\Omega} -\mathbf{A} : \mathbf{D}^2 w v \, dx, \tag{15}$$

then, for all $w \in H^2(\Omega) \cap H_0^1(\Omega)$

$$\sup_{v \in L^2(\Omega)} \frac{a(w, v)}{\|v\|_{L^2(\Omega)}} \geq C\|\Delta w\|_{L^2(\Omega)}. \tag{16}$$

This is true under a variety of conditions. For example, those of Theorem 1. Note that for $d = 2$, this criterion is satisfied for all essentially bounded, symmetric positive definite \mathbf{A} .

3 Discretisation

Let \mathcal{T} be a shape regular triangulation of Ω , namely, \mathcal{T} is a finite family of sets such that

- 1) $K \in \mathcal{T}$ implies K is an open simplex (segment for $d = 1$, triangle for $d = 2$, tetrahedron for $d = 3$),
- 2) for any $K, J \in \mathcal{T}$ we have that $\overline{K} \cap \overline{J}$ is a full subsimplex (i.e., it is either \emptyset , a vertex, an edge, a face, or the whole of \overline{K} and \overline{J}) of either \overline{K} and \overline{J} and
- 3) $\bigcup_{K \in \mathcal{T}} \overline{K} = \overline{\Omega}$.

We use the convention where $h : \Omega \rightarrow \mathbb{R}$ denotes the *mesh size function* of \mathcal{T} , i.e.,

$$h(\mathbf{x}) := \max_{\overline{K} \ni \mathbf{x}} h_K, \tag{17}$$

where h_K is the diameter of K . We let \mathcal{E} be the skeleton (set of common interfaces) of the triangulation \mathcal{T} and say $e \in \mathcal{E}$ if e is on the interior of Ω and $e \in \partial\Omega$ if e lies on the boundary $\partial\Omega$.

Let $\mathbb{P}^k(\mathcal{T})$ denote the space of piecewise polynomials of degree k over the triangulation \mathcal{T} , i.e.,

$$\mathbb{P}^k(\mathcal{T}) = \{ \phi : \phi|_K \in \mathbb{P}^k(K) \}, \tag{18}$$

and introduce the *finite element spaces*

$$\mathbb{V} = \mathbb{V}(\mathcal{T}, k) := \mathbb{P}^k(\mathcal{T}), \tag{19}$$

$$\mathbb{V}_0 = \mathbb{V}_0(\mathcal{T}, k) := \{ \phi \in \mathbb{P}^k(\mathcal{T}) : \phi|_{\partial\Omega} = 0 \}, \tag{20}$$

to be the usual spaces of discontinuous piecewise polynomial functions.

Remark 2 (Generalised Hessian) Assume $v \in H^2(\Omega)$, let $\mathbf{n} : \partial\Omega \rightarrow \mathbb{R}^d$ be the outward pointing normal of Ω , then the Hessian D^2v of v , satisfies the following identity:

$$\int_{\Omega} D^2v \phi \, dx = - \int_{\Omega} \nabla v \otimes \nabla \phi \, dx + \int_{\partial\Omega} \nabla v \otimes \mathbf{n} \phi \, ds, \quad \forall \phi \in H^1(\Omega). \tag{21}$$

If $v \in H^1(\Omega)$, the right-hand side of (21) is still well defined in view of duality, in this case we set

$$\langle D^2v | \phi \rangle = - \int_{\Omega} \nabla v \otimes \nabla \phi \, dx + \int_{\partial\Omega} \nabla v \otimes \mathbf{n} \phi \, ds, \quad \forall \phi \in H^1(\Omega), \tag{22}$$

where the last term is understood as a duality pairing.

Definition 2 (Broken Sobolev spaces, trace spaces) We introduce the broken Sobolev space

$$H^k(\mathcal{T}) := \{ \phi : \phi|_K \in H^k(K), \text{ for each } K \in \mathcal{T} \}. \tag{22}$$

We also make use of functions defined in these broken spaces restricted to the skeleton of the triangulation. This requires an appropriate trace space

$$\mathcal{X}(\mathcal{E}) := \prod_{K \in \mathcal{T}} L^2(\partial K) = \prod_{K \in \mathcal{T}} H^{\frac{1}{2}}(K). \tag{24}$$

Definition 3 (Jumps, averages and tensor jumps) We define average, jump and tensor jump operators for arbitrary scalar functions $v \in \mathcal{X}(\mathcal{E})$, vectors $\mathbf{v} \in \mathcal{X}(\mathcal{E})^d$ and matrices $\mathbf{V} \in \mathcal{X}(\mathcal{E})^{d \times d}$ as

$$\{v\} = \frac{1}{2}(v|_{K_1} + v|_{K_2}), \quad \{\mathbf{v}\} = \frac{1}{2}(\mathbf{v}|_{K_1} + \mathbf{v}|_{K_2}), \tag{25}$$

$$[[v]] = v|_{K_1} \mathbf{n}_{K_1} + v|_{K_2} \mathbf{n}_{K_2}, \quad [[\mathbf{v}]] = (\mathbf{v}|_{K_1}) \cdot \mathbf{n}_{K_1} + (\mathbf{v}|_{K_2}) \cdot \mathbf{n}_{K_2}, \tag{26}$$

$$[[\mathbf{V}]] = \mathbf{V}|_{K_1} \mathbf{n}_{K_1} + \mathbf{V}|_{K_2} \mathbf{n}_{K_2}, \quad [[\mathbf{v}]]_{\otimes} = \mathbf{v}|_{K_1} \otimes \mathbf{n}_{K_1} + \mathbf{v}|_{K_2} \otimes \mathbf{n}_{K_2}. \tag{27}$$

Note that on the boundary of the domain $\partial\Omega$, the jump and average operators are defined as

$$\{v\}|_{\partial\Omega} := v, \quad \{\mathbf{v}\}|_{\partial\Omega} := \mathbf{v}, \tag{28}$$

$$[[v]]|_{\partial\Omega} := v \mathbf{n}, \quad [[\mathbf{v}]]|_{\partial\Omega} := \mathbf{v} \cdot \mathbf{n}, \tag{29}$$

$$\llbracket \mathbf{V} \rrbracket \Big|_{\partial\Omega} := \mathbf{V}\mathbf{n}, \quad \llbracket \mathbf{v} \rrbracket_{\otimes} \Big|_{\partial\Omega} := \mathbf{v} \otimes \mathbf{n}. \tag{30}$$

We will often use the following Proposition which we state in full for clarity but whose proof is merely using the identities in Definition 3.

Proposition 1 (Elementwise integration) *For a generic vector-valued function \mathbf{p} and scalar-valued function ϕ , we have*

$$\sum_{K \in \mathcal{T}} \int_K \operatorname{div}(\mathbf{p})\phi \, dx = \sum_{K \in \mathcal{T}} \left(- \int_K \mathbf{p} \cdot \nabla_h \phi \, dx + \int_{\partial K} \phi \mathbf{p} \cdot \mathbf{n}_K \, ds \right), \tag{31}$$

where $\nabla_h = (\mathbf{D}_h)^\top$ is the elementwise spatial gradient. Furthermore, If we have $\mathbf{p} \in \mathcal{T}(\mathcal{E} \cup \partial\Omega)^d$ and $\phi \in \mathcal{T}(\mathcal{E} \cup \partial\Omega)$, the following identity holds

$$\sum_{K \in \mathcal{T}} \int_{\partial K} \phi \mathbf{p} \cdot \mathbf{n}_K \, ds = \int_{\mathcal{E}} \llbracket \mathbf{p} \rrbracket \{ \phi \} \, ds + \int_{\partial\Omega} \llbracket \phi \rrbracket \cdot \{ \mathbf{p} \} \, ds = \int_{\partial\Omega} \llbracket \mathbf{p}\phi \rrbracket \, ds. \tag{32}$$

An equivalent tensor formulation of (31)–(32) is

$$\sum_{K \in \mathcal{T}} \int_K \mathbf{D}_h \mathbf{p} \phi \, dx = \sum_{K \in \mathcal{T}} \left(- \int_K \mathbf{p} \otimes \nabla_h \phi \, dx + \int_{\partial K} \phi \mathbf{p} \otimes \mathbf{n}_K \, ds \right), \tag{33}$$

where

$$\sum_{K \in \mathcal{T}} \int_{\partial K} \phi \mathbf{p} \otimes \mathbf{n}_K \, ds = \int_{\mathcal{E}} \llbracket \mathbf{p} \rrbracket_{\otimes} \{ \phi \} \, ds + \int_{\partial\Omega} \llbracket \phi \rrbracket \otimes \{ \mathbf{p} \} \, ds = \int_{\partial\Omega} \llbracket \mathbf{p}\phi \rrbracket_{\otimes} \, ds. \tag{34}$$

In addition for matrix-valued \mathbf{V} , we have that

$$\sum_{K \in \mathcal{T}} \int_K (\mathbf{D}_h \mathbf{p}) : \mathbf{V} \, dx = \sum_{K \in \mathcal{T}} \left(- \int_K \mathbf{p} : \mathbf{D}_h \mathbf{V} \, dx + \int_{\partial\Omega} (\mathbf{V}\mathbf{p}) \cdot \mathbf{n} \, ds \right) \tag{35}$$

and

$$\sum_{K \in \mathcal{T}} \int_{\partial\Omega} (\mathbf{V}\mathbf{p}) \cdot \mathbf{n} \, ds = \int_{\mathcal{E}} \llbracket \mathbf{V} \rrbracket \cdot \{ \mathbf{p} \} \, ds + \int_{\partial\Omega} \llbracket \mathbf{p} \rrbracket_{\otimes} : \{ \mathbf{V} \} \, ds = \int_{\partial\Omega} \llbracket \mathbf{V}\mathbf{p} \rrbracket \, ds. \tag{36}$$

3.1 Construction of an Appropriate Discrete Hessian

We now use the framework set out in [4] to construct a general notion of discrete Hessian. We first give a definition using a flux formulation.

Definition 4 (Generalised finite element Hessian: flux formulation) Let $u \in H^2(\mathcal{T})$, $\hat{U} : H^1(\mathcal{T}) \rightarrow \mathcal{T}(\mathcal{E} \cup \partial\Omega)$ be a linear form and $\hat{\mathbf{p}} : H^2(\mathcal{T}) \times H^1(\mathcal{T})^d \rightarrow \mathcal{T}(\mathcal{E} \cup \partial\Omega)^d$ a bilinear form representing approximations to u and ∇u over the skeleton of the triangulation. Then, we define the generalised finite element Hessian $\mathbf{H}[u]$ as the solution of

$$\int_K \mathbf{H}[u] \Phi \, dx = - \int_K \mathbf{p} \otimes \nabla_h \Phi \, dx + \int_{\partial K} \hat{\mathbf{p}}_K \otimes \mathbf{n} \Phi \, ds, \quad \forall \Phi \in H^1(\mathcal{T}), \quad (37)$$

$$\int_K \mathbf{p} \otimes \mathbf{q} \, dx = - \int_K u \, D_h \mathbf{q} \, dx + \int_{\partial K} \mathbf{q} \otimes \mathbf{n} \hat{U}_K \, ds \quad (38)$$

for all $\Phi \in \mathbb{V}$.

We now present the primal formulation for the generalised finite element Hessian.

Theorem 2 (Generalised finite element Hessian: primal form) *Let $u \in H^2(\mathcal{T})$ and let \hat{U} and $\hat{\mathbf{p}}$ be defined as in Definition 4. Then, the generalised finite element Hessian $H[u]$ is given for each $\Phi \in \mathbb{V}$ as*

$$\begin{aligned} \int_{\Omega} \mathbf{H}[u] \Phi \, dx &= - \int_{\Omega} \nabla_h u \otimes \nabla_h \Phi \, dx + \int_{\partial \cup \partial \Omega} \llbracket \Phi \rrbracket \otimes \{ \hat{\mathbf{p}} \} \, ds + \int_{\mathcal{E}} \{ \Phi \} \llbracket \hat{\mathbf{p}} \rrbracket_{\otimes} \, ds \\ &\quad - \int_{\mathcal{E}} \{ \hat{U} - u \} \llbracket \nabla_h \Phi \rrbracket_{\otimes} \, ds - \int_{\partial \cup \partial \Omega} \llbracket \hat{U} - u \rrbracket \otimes \{ \nabla_h \Phi \} \, ds. \end{aligned} \quad (39)$$

Proof Note that in view of Definition 3 for generic vector fields $\mathbf{q} \in \mathbb{W}$ and $v \in \mathbb{V}$, we have the following identity:

$$\sum_{K \in \mathcal{T}} \int_{\partial K} v \mathbf{q} \otimes \mathbf{n} \, ds = \int_{\partial \cup \partial \Omega} \llbracket v \rrbracket \otimes \{ \mathbf{q} \} \, ds + \int_{\mathcal{E}} \{ v \} \llbracket \mathbf{q} \rrbracket_{\otimes} \, ds. \quad (40)$$

Then summing (37) over $K \in \mathcal{T}$ and making use of the identity (40) we see

$$\begin{aligned} \int_{\Omega} \mathbf{H}[u] \Phi \, dx &= \sum_{K \in \mathcal{T}} \int_K \mathbf{H}[u] \Phi \, dx = \sum_{K \in \mathcal{T}} \left(- \int_K \mathbf{p} \otimes \nabla_h \Phi \, dx + \int_{\partial K} \hat{\mathbf{p}}_K \otimes \mathbf{n} \Phi \, ds \right) \\ &= - \int_{\Omega} \mathbf{p} \otimes \nabla_h \Phi \, dx + \int_{\partial \cup \partial \Omega} \llbracket \Phi \rrbracket \otimes \{ \hat{\mathbf{p}} \} \, ds + \int_{\mathcal{E}} \{ \Phi \} \llbracket \hat{\mathbf{p}} \rrbracket_{\otimes} \, ds. \end{aligned} \quad (41)$$

Using the same argument for (38)

$$\begin{aligned} \int_{\Omega} \mathbf{p} \otimes \mathbf{q} \, dx &= \sum_{K \in \mathcal{T}} \int_K \mathbf{p} \otimes \mathbf{q} \, dx = \sum_{K \in \mathcal{T}} \left(- \int_K u \, D_h \mathbf{q} \, dx + \int_{\partial K} \mathbf{q} \otimes \mathbf{n} \hat{U}_K \, ds \right) \\ &= - \int_{\Omega} u \, D_h \mathbf{q} \, dx + \int_{\partial \cup \partial \Omega} \llbracket \hat{U} \rrbracket \otimes \{ \mathbf{q} \} \, ds + \int_{\mathcal{E}} \{ \hat{U} \} \llbracket \mathbf{q} \rrbracket_{\otimes} \, ds. \end{aligned} \quad (42)$$

Note that, again making use of (40), we have for each $\mathbf{q} \in H^1(\mathcal{T})^d$ and $v \in H^1(\mathcal{T})$ that

$$\int_{\Omega} \mathbf{q} \otimes \nabla_h v \, dx = - \int_{\Omega} D_h \mathbf{q} v \, dx + \int_{\partial \cup \partial \Omega} \{ \mathbf{q} \} \otimes \llbracket v \rrbracket \, ds + \int_{\mathcal{E}} \llbracket \mathbf{q} \rrbracket_{\otimes} \{ v \} \, ds. \quad (43)$$

Taking $v = u$ in (43) and substituting into (38), we see

$$\int_{\Omega} \mathbf{p} \otimes \mathbf{q} \, dx = \int_{\Omega} \mathbf{q} \otimes \nabla_h u \, dx + \int_{\partial \cup \partial \Omega} \llbracket \hat{U} - u \rrbracket \otimes \{ \mathbf{q} \} \, ds + \int_{\mathcal{E}} \{ \hat{U} - u \} \llbracket \mathbf{q} \rrbracket_{\otimes} \, ds. \quad (44)$$

Now choosing $\mathbf{q} = \nabla_h \Phi$ and substituting (44) into (37), we arrive at the fully generalised finite element Hessian given by (39).

Remark 3 (Consistent representations of the gradient operator) If one were interested in consistent representations of other derivatives, for example the gradient operator, one would need to modify the proof of Theorem 2. Examples of consistent gradient representations can be found in [4]. See also [10, 11, 17]. Using this methodology, it should be possible to construct an entire hierarchy of derivatives.

Example 1 An example of a DG formulation for the approximation to the Hessian, D^2u , can be derived by taking the fluxes in the following way:

$$\widehat{U} = \begin{cases} \{u_h\} & \text{over } \mathcal{E}, \\ 0 & \text{on } \partial\Omega, \end{cases} \tag{45}$$

$$\widehat{\mathbf{p}} = \{\nabla_h u_h\} \text{ on } \mathcal{E} \cup \partial\Omega. \tag{46}$$

The result is a discrete representation of the Hessian $\mathbf{H}[u_h]$ as a unique element of $\mathbb{V}^{d \times d}$ such that

$$\int_{\Omega} \mathbf{H}[u_h] \Phi \, dx = - \int_{\Omega} \nabla_h u_h \otimes \nabla_h \Phi \, dx + \int_{\mathcal{E} \cup \partial\Omega} \llbracket u_h \rrbracket \otimes \{\nabla_h \Phi\} + \llbracket \Phi \rrbracket \otimes \{\nabla_h u_h\} \, ds.$$

We can also derive unsymmetric approximations with fluxes similar to those used for the NIPG method. This is given by taking $\theta = -1$ in the following, while the symmetric version is given by $\theta = 1$:

$$\begin{aligned} \int_{\Omega} \mathbf{H}[u_h] \Phi \, dx &= - \int_{\Omega} \nabla_h u_h \otimes \nabla_h \Phi \, dx \\ &\quad + \int_{\mathcal{E} \cup \partial\Omega} \theta \llbracket u_h \rrbracket \otimes \{\nabla_h \Phi\} + \llbracket \Phi \rrbracket \otimes \{\nabla_h u_h\} \, ds \\ &= \int_{\Omega} D_h^2 u_h \Phi \, dx - \int_{\mathcal{E}} \llbracket \nabla_h u_h \rrbracket \otimes \{\Phi\} \, ds \\ &\quad + \int_{\mathcal{E} \cup \partial\Omega} \theta \llbracket u_h \rrbracket \otimes \{\nabla_h \Phi\} \, ds, \quad \forall \Phi \in \mathbb{V}. \end{aligned} \tag{47}$$

3.2 The Discontinuous Nonvariational Finite Element Method

We are now in a position to state the numerical method for the approximation of (1). We look to find $u_h \in \mathbb{V}_0$ together with $\mathbf{H}[u_h] \in \mathbb{V}^{d \times d}$ such that

$$\mathcal{A}_h(u_h, \Psi) = l(\Psi), \quad \forall \Psi \in \mathbb{V}_0 \tag{48}$$

with

$$\mathcal{A}_h(u_h, \Psi) := \int_{\Omega} -\mathbf{A} : \mathbf{H}[u_h] \Psi \, dx + \int_{\mathcal{E} \cup \partial\Omega} \sigma h^{-1} \llbracket u_h \rrbracket \cdot \llbracket \Psi \rrbracket \, ds, \tag{49}$$

$$I(\Psi) := \int_{\Omega} f\Psi \, dx, \tag{50}$$

where the *penalisation parameter* $\sigma > 0$ is to be chosen sufficiently large.

Using the L^2 projection operator $P_k : L^2(\Omega) \rightarrow \mathbb{V}$ defined for $v \in L^2(\Omega)$ through

$$\int_{\Omega} P_k(v)\Psi \, dx = \int_{\Omega} v\Psi \, dx, \quad \forall \Psi \in \mathbb{V}, \tag{51}$$

it is possible to eliminate the finite element Hessian from the bilinear form for sufficiently smooth A .

Lemma 1 (Elimination of the finite element Hessian in a general setting) *If the fluxes are chosen as in Example 1, then*

$$\begin{aligned} \mathcal{A}_h(u_h, \Psi) &= \int_{\Omega} D_h(P_k(\Psi A)) \nabla_h u_h \, dx - \int_{\mathcal{S} \cup \partial\Omega} \theta \llbracket u_h \rrbracket \cdot \{D_h(P_k(\Psi A))\} \, ds \\ &\quad - \int_{\mathcal{S} \cup \partial\Omega} \llbracket P_k(\Psi A) \rrbracket \cdot \{\nabla_h u_h\} \, ds + \int_{\mathcal{S} \cup \partial\Omega} \sigma h^{-1} \llbracket u_h \rrbracket \cdot \llbracket \Psi \rrbracket \, ds. \end{aligned} \tag{52}$$

Proof This follows from the following identity:

$$\begin{aligned} \int_{\Omega} -A : \mathbf{H}[u_h] \Psi \, dx &= \int_{\Omega} -\mathbf{H}[u_h] : (\Psi A) \, dx = \int_{\Omega} -\mathbf{H}[u_h] : P_k(\Psi A) \, dx \\ &= \int_{\Omega} D_h(P_k(\Psi A)) \nabla_h u_h \, dx - \int_{\mathcal{S} \cup \partial\Omega} \theta \llbracket u_h \rrbracket \cdot \{D_h(P_k(\Psi A))\} \, ds \\ &\quad - \int_{\mathcal{S} \cup \partial\Omega} \llbracket P_k(\Psi A) \rrbracket \cdot \{\nabla_h u_h\} \, ds. \end{aligned} \tag{53}$$

Remark 4 The solution of the problem in this form is nontrivial due to the global $L^2(\Omega)$ projection appearing in the formulation. However, in the discontinuous setting, the global $L^2(\Omega)$ projection is in fact computable locally. We may actually exploit this fact to optimise our schemes efficiency. We will discuss this further in the sequel.

Example 2 (Laplacian formulation) Note that if in (1) we have that $A = I$, then we have that

$$f = -A : D^2 u = -\Delta u \tag{54}$$

and our bilinear form reduces to

$$\begin{aligned} \mathcal{A}_h(u_h, \Psi) &= \int_{\Omega} (\nabla_h \Psi) \cdot \nabla_h u_h \, dx - \int_{\mathcal{S} \cup \partial\Omega} \theta \llbracket u_h \rrbracket \cdot \{\nabla_h \Psi\} \, ds \\ &\quad - \int_{\mathcal{S} \cup \partial\Omega} (\llbracket \Psi \rrbracket \cdot \{\nabla_h u_h\} - \sigma h^{-1} \llbracket u_h \rrbracket \cdot \llbracket \Psi \rrbracket) \, ds, \end{aligned} \tag{55}$$

since $P_k(\Psi A) = \Psi I$.

The nonvariational finite element method, thus, coincides with the classical (symmetric) interior penalty method for the Laplacian [16].

Remark 5 (Relation to standard DG methods) It is not difficult to prove that choosing to numerical fluxes in the same way as presented in [4, Table 3.2] results in the same correlation to the DG methods summarised in the aforementioned paper for the case that \mathbf{A} is constant. For brevity, we will not prove this here.

Note that when \mathbf{A} is not constant, we have that the nonvariational finite element method does *not* coincide with its standard variational finite element counterpart. Indeed, the method is able to successfully cope with classes of advection dominated problems not falling under the variational framework without special treatment [32, §4.2] which is illustrated by the result of Lemma 1.

We conclude this section with a proof of consistency of the method and then show that Galerkin orthogonality holds.

Lemma 2 (Consistency) *Let $u \in H^2(\Omega)$ and assume that the numerical fluxes are chosen in a consistent fashion in the sense of [4, §3.1], that is,*

$$\hat{U} = u|_{\mathcal{E} \cup \partial\Omega}, \tag{56}$$

$$\hat{\mathbf{p}} = \nabla u|_{\mathcal{E} \cup \partial\Omega}. \tag{57}$$

Then for $\Phi \in \mathbb{V}$

$$\int_{\Omega} \mathbf{H}[u] \Phi \, dx = \int_{\Omega} D^2 u \Phi \, dx. \tag{58}$$

Therefore, we have that $\mathbf{H}[u] = P_k(D^2 u)$.

Proof Applying Proposition 1 to the first term in the definition of $\mathbf{H}[u]$ yields

$$\begin{aligned} \int_{\Omega} \mathbf{H}[u] \Phi \, dx &= \int_{\Omega} D^2 u \Phi \, dx + \int_{\mathcal{E}} \llbracket \hat{\mathbf{p}} - \nabla u \rrbracket_{\otimes} \{ \Phi \} \, ds + \int_{\mathcal{E} \cup \partial\Omega} \{ \hat{\mathbf{p}} - \nabla u \} \otimes \llbracket \Phi \rrbracket \, ds \\ &\quad - \int_{\mathcal{E}} \{ \hat{U} - u \} \llbracket \nabla_h \Phi \rrbracket_{\otimes} \, ds - \int_{\mathcal{E} \cup \partial\Omega} \llbracket \hat{U} - u \rrbracket \otimes \{ \nabla_h \Phi \} \, ds \\ &= \int_{\Omega} D^2 u \Phi \, dx, \quad \forall \Phi \in \mathbb{V}, \end{aligned} \tag{59}$$

which proves the results under the consistency conditions on the fluxes.

Lemma 3 (Galerkin orthogonality) *Let $u \in H^2(\Omega) \cap H_0^1(\Omega)$ be a strong solution to the problem (1) and let $u_h \in \mathbb{V}_0$ be its nonvariational finite element approximation. Assume that the numerical fluxes \hat{U} and $\hat{\mathbf{p}}$ are consistent, then we have the following orthogonality result:*

$$\mathcal{A}_h(u_h - u, \Psi) = J(\Psi), \quad \forall \Psi \in \mathbb{V}_0, \tag{60}$$

with the error functional given by

$$J(\Psi) = \int_{\Omega} (D^2 u - \mathbf{H}[u]) : (\mathbf{A}\Psi) \, dx. \tag{61}$$

Proof Using the consistency result and that $\llbracket u \rrbracket = 0$, we conclude

$$\begin{aligned} \mathcal{A}_h(u_h - u, \Psi) &= \mathcal{A}_h(u_h, \Psi) + \int_{\Omega} \mathbf{A} : \mathbf{H}[u] \Psi \, dx = l(\Psi) + \int_{\Omega} \mathbf{H}[u] : (\mathbf{A}\Psi) \, dx \\ &= - \int_{\Omega} \mathbf{A} : \mathbf{D}^2 u \Psi - \mathbf{H}[u] : (\mathbf{A}\Psi) \, dx = J(\Psi), \end{aligned}$$

concluding the proof.

Remark 6 If \mathbf{A} is piecewise constant, then since $\mathbf{H}[u] = \mathbf{P}_k(\mathbf{D}^2 u)$ we have $J(\Psi) = 0$ and we recover the usual Galerkin orthogonality $\mathcal{A}_h(u_h - u, \Psi) = 0$.

Definition 5 ($H^1(\mathcal{T})$ and $H^2(\mathcal{T})$ norms) We introduce the broken $H^1(\mathcal{T})$ and $H^2(\mathcal{T})$ norms as

$$\|u_h\|_{\text{DG},1}^2 := \|\nabla_h u_h\|_{L^2(\Omega)}^2 + h^{-1} \|\llbracket u_h \rrbracket\|_{L^2(\mathcal{E})}^2, \tag{62}$$

$$\|u_h\|_{\text{DG},2}^2 := \|\mathbf{D}_h^2 u_h\|_{L^2(\Omega)}^2 + h^{-1} \|\llbracket \nabla_h u_h \rrbracket\|_{L^2(\mathcal{E})}^2 + h^{-3} \|\llbracket u_h \rrbracket\|_{L^2(\mathcal{E})}^2. \tag{63}$$

These are equivalent to their continuous equivalent norms for functions in \mathbb{V} .

Proposition 2 (Projection approximation in \mathbb{V}) *Let $\mathbf{P}_k : L^2(\Omega) \rightarrow \mathbb{V}$ be the $L^2(\Omega)$ orthogonal projection operator defined by (51). Using standard approximation arguments, we have that*

$$\begin{cases} \|v - \mathbf{P}_k v\|_{\text{DG},1} \leq Ch^k |v|_{H^{k+1}(\Omega)}, \\ \|v - \mathbf{P}_k v\|_{L^2(\Omega)} \leq Ch^{k+1} |v|_{H^{k+1}(\Omega)}. \end{cases} \tag{64}$$

Lemma 4 (Stability of \mathbf{H} [37, Theorem 4.10]) *Let \mathbf{H} be defined as in Example 1. Then the DG Hessian is stable in the sense that*

$$\|\mathbf{D}_h^2 v_h - \mathbf{H}[v_h]\|_{L^2(\Omega)}^2 \leq C \left(\int_{\mathcal{E}} h^{-1} |\llbracket \nabla_h v_h \rrbracket|^2 + h^{-3} |\llbracket v_h \rrbracket|^2 \, ds \right). \tag{65}$$

Consequently, we have

$$\|\mathbf{H}[v_h]\|_{L^2(\Omega)}^2 \leq C \|v_h\|_{\text{DG},2}^2. \tag{66}$$

4 A Posteriori Analysis

The discrete bilinear form (49) only makes sense over $H^2(\mathcal{T}) \times H^2(\mathcal{T})$. To allow for an a posteriori bound we require an extension to ensure the appropriate stability arguments can be applied. We require the bilinear form to be extended to $H^2(\mathcal{T}) \times L^2(\Omega)$. To do this for $(u, v) \in H^2(\mathcal{T}) \times L^2(\Omega)$, we define

$$\mathcal{A}_h(u, v) := \int_{\Omega} \mathbf{A} : \mathbf{H}(u)v + \int_{\mathcal{E}} \sigma h_e^{-1} \llbracket u \rrbracket \cdot \llbracket \mathbf{P}_\kappa v \rrbracket. \tag{67}$$

Remark 7 Notice that the modified bilinear form (67) coincides with (49) over $\mathbb{V} \times \mathbb{V}$ and that it satisfies

$$\mathcal{A}_h(u, v) \leq C \|u\|_{\text{DG},2} \|v\|_{L^2(\Omega)} \text{ for } (u, v) \in (H^2(\mathcal{T}) \times L^2(\Omega)). \tag{68}$$

Assumption 2 (Existence of an H^2 reconstruction) We will assume there exists an operator $\mathcal{R} : \mathbb{V} \rightarrow H^2(\Omega)$ such that

$$\|\mathcal{R}(u_h) - u_h\|_{\text{DG},2}^2 \leq C \left(\|h^{-1/2} \llbracket \nabla u_h \rrbracket \|_{L^2(\mathcal{E})}^2 + \|h^{-3/2} \llbracket u_h \rrbracket \|_{L^2(\mathcal{E})}^2 \right). \tag{69}$$

Such an example is given, for $d = 2$, in [26, Lem 3.1] and consists of averaging techniques onto macro-elements on the Hsieh-Clough-Tocher space or in [9], for $d = 3$, using virtual element spaces.

Proposition 3 (Abstract upper bound) *Let $u \in H^2(\Omega)$ solve (1), $u_h \in \mathbb{V}$ be the finite element approximation given by (48) and $\mathcal{R}(u_h) \in H^2(\Omega)$ be a post-processor satisfying (69). Then,*

$$\begin{aligned} & \|u - \mathcal{R}(u_h)\|_{H^2(\Omega)} \\ & \leq \frac{1}{C} \sup_{v \in L^2(\Omega)} \frac{\left(l(v) - \mathcal{A}_h(u_h, v) + \mathcal{A}_h(u_h - \mathcal{R}(u_h), v) + \mathcal{A}_h(\mathcal{R}(u_h), v) - \mathcal{A}(\mathcal{R}(u_h), v) \right)}{\|v\|_{L^2(\Omega)}}. \end{aligned} \tag{70}$$

Proof Making use of the Miranda-Talenti inequality [34, 39] and the inf-sup condition (16) we see that

$$C \|u - \mathcal{R}(u_h)\|_{H^2(\Omega)} \leq \tilde{C} \|\Delta(u - \mathcal{R}(u_h))\|_{L^2(\Omega)} \leq \sup_{v \in L^2(\Omega)} \frac{\mathcal{A}(u - \mathcal{R}(u_h), v)}{\|v\|_{L^2(\Omega)}}. \tag{71}$$

Now, adding and subtracting appropriate terms, we have

$$\begin{aligned} \mathcal{A}(u - \mathcal{R}(u_h), v) &= l(v - v_h) - \mathcal{A}_h(u_h, v - v_h) \\ &\quad + \mathcal{A}_h(u_h - \mathcal{R}(u_h), v) + \mathcal{A}_h(\mathcal{R}(u_h), v) - \mathcal{A}(\mathcal{R}(u_h), v), \end{aligned} \tag{72}$$

and the result follows.

Theorem 3 (A posteriori bound) *Let $u \in H^2(\Omega)$ solve (1) and $u_h \in \mathbb{V}$ be the finite element approximation given by (48). Then*

$$\|u - u_h\|_{\text{DG},2} \leq C \left(\sum_{K \in \mathcal{T}} \eta_K^2 + \sum_{e \in \mathcal{E}} \eta_e^2 \right)^{1/2}, \tag{73}$$

where

$$\eta_K^2 := \|f + \mathbf{A} : \mathbf{H}(u_h)\|_{L^2(K)}^2, \tag{74}$$

$$\eta_e^2 := h_e^{-1} \|[\![\nabla u_h]\!] \|_{L^2(e)}^2 + h_e^{-3} \|[\![u_h]\!] \|_{L^2(e)}^2. \tag{75}$$

Proof Beginning from Proposition 3 we note that the bound can be split into a residual component \mathcal{I}_1 , a nonconformity component \mathcal{I}_2 and an inconsistency component \mathcal{I}_3 as follows:

$$\begin{aligned} \|u - \mathcal{R}(u_h)\|_{H^2(\Omega)} &\leq \frac{1}{C} \left(\sup_{v \in L^2(\Omega), \|v\|_{L^2(\Omega)} \leq 1} (l(v) - \mathcal{A}_h(u_h, v)) \right. \\ &\quad + \sup_{v \in L^2(\Omega), \|v\|_{L^2(\Omega)} \leq 1} \mathcal{A}_h(u_h - \mathcal{R}(u_h), v) \\ &\quad \left. + \sup_{v \in L^2(\Omega), \|v\|_{L^2(\Omega)} \leq 1} \mathcal{A}_h(\mathcal{R}(u_h), v) - \mathcal{A}(\mathcal{R}(u_h), v) \right) \\ &=: \mathcal{I}_1 + \mathcal{I}_2 + \mathcal{I}_3. \end{aligned} \tag{76}$$

Now, we proceed to bound these term by term. In view of Cauchy-Schwartz, we have

$$\begin{aligned} \mathcal{I}_1 &\leq \sup_{v \in L^2(\Omega), \|v\|_{L^2(\Omega)} \leq 1} \int_{\Omega} (f - \mathbf{A} : \mathbf{H}(u_h)) v \\ &\leq \left(\sum_{K \in \mathcal{T}} \|f - \mathbf{A} : \mathbf{H}(u_h)\|_{L^2(K)}^2 \right)^{1/2}. \end{aligned} \tag{77}$$

For the nonconformity term, we note that

$$\mathcal{I}_2 \leq C \| \mathcal{R}(u_h) - u_h \|_{DG,2} \leq C \left(\sum_{e \in K} h_e^{-1} \|[\![\nabla_h v_h]\!] \|_{L^2(e)}^2 + h_e^{-3} \|[\![v_h]\!] \|_{L^2(e)}^2 \right)^{1/2}, \tag{78}$$

by the properties of the reconstruction given in (69).

For the inconsistency term, we have that

$$\begin{aligned} \mathcal{I}_3 &= \sup_{v \in L^2(\Omega), \|v\|_{L^2(\Omega)} \leq 1} \frac{1}{C} \int_{\Omega} \mathbf{A} : (\mathbf{H}(\mathcal{R}(u_h)) - \mathbf{D}^2 \mathcal{R}(u_h)) v \\ &\leq C \| \mathbf{A} \|_{L^\infty(\Omega)} \| \mathbf{H}(\mathcal{R}(u_h)) - \mathbf{D}^2 \mathcal{R}(u_h) \|_{L^2(\Omega)} \\ &\leq C \| \mathbf{A} \|_{L^\infty(\Omega)} \left(\| \mathbf{H}(\mathcal{R}(u_h)) - \mathbf{H}(u_h) \|_{L^2(\Omega)} + \| \mathbf{H}(u_h) - \mathbf{D}_h^2 u_h \|_{L^2(\Omega)} \right. \\ &\quad \left. + \| \mathbf{D}_h^2 u_h - \mathbf{D}^2 \mathcal{R}(u_h) \|_{L^2(\Omega)} \right). \end{aligned} \tag{79}$$

Making use of Lemma 4 and the properties of the reconstruction (69), we see that

$$\begin{aligned} \|\mathbf{H}(\mathcal{R}(u_h)) - \mathbf{H}(u_h)\|_{L^2(\Omega)} &\leq C \|\mathcal{R}(u_h) - u_h\|_{\text{DG},2} \\ &\leq C \left(\sum_{e \in \mathcal{E}} h_e^{-1} \|\llbracket \nabla_h v_h \rrbracket\|_{L^2(e)}^2 + h_e^{-3} \|\llbracket v_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2}. \end{aligned} \tag{80}$$

Again from Lemma 4 we have that

$$\|\mathbf{H}(u_h) - \mathbf{D}_h^2 u_h\|_{L^2(\Omega)} \leq C \left(\sum_{e \in \mathcal{E}} h_e^{-1} \|\llbracket \nabla_h v_h \rrbracket\|_{L^2(e)}^2 + h_e^{-3} \|\llbracket v_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2}, \tag{81}$$

and finally (69) gives that

$$\|\mathbf{D}_h^2 u_h - \mathbf{D}^2 \mathcal{R}(u_h)\|_{L^2(\Omega)} \leq C \left(\sum_{e \in \mathcal{E}} h_e^{-1} \|\llbracket \nabla_h v_h \rrbracket\|_{L^2(e)}^2 + h_e^{-3} \|\llbracket v_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2}. \tag{82}$$

Hence, we have that

$$\mathcal{J}_3 \leq C \left(\sum_{e \in \mathcal{E}} h_e^{-1} \|\llbracket \nabla_h v_h \rrbracket\|_{L^2(e)}^2 + h_e^{-3} \|\llbracket v_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2}. \tag{83}$$

Collecting (77), (78) and (83) yields the desired result.

Proposition 4 (A posteriori lower bound) *Using the notation of Theorem 3, through standard a posteriori techniques, we have a lower bound of the form*

$$\eta_K + \sum_{e \in \partial K} \eta_e \leq C \left(\|\mathbf{D}^2 u - \mathbf{H}[u_h]\|_{L^2(\hat{K})} + \sum_{K \in \hat{\mathcal{K}}} \eta_{I,K} \right), \tag{84}$$

where

$$\eta_{I,K}^2 := \|(P_{k-2} \mathbf{A} - \mathbf{A}) : \mathbf{H}(u_h)\|_{L^2(K)}^2 + \|f - P_{k-2} f\|_{L^2(K)}^2. \tag{85}$$

5 Numerical Experiments

In this section, we detail numerical experiments carried out in the finite element package DUNE-FEM [15] which is based on the DUNE software framework [6, 7]. The code makes use of the newly developed Python frontend [15] and the unified form language [3] is used to provide the problem data. The code will be made freely available within the DUNE-FEM-tutorial in a future release [13].

We present some benchmark problems designed such that the exact solution is known. In each of the experiments the domain $\Omega = [0, 1]^2$ and we consider the coefficient matrix to be

$$\mathbf{A}(\mathbf{x}) = \begin{bmatrix} 1 & b(\mathbf{x}) \\ b(\mathbf{x}) & a(\mathbf{x}) \end{bmatrix} \tag{86}$$

varying $a(\mathbf{x})$ and $b(\mathbf{x})$. We study three different choices using $\mathbf{x} = (x_1, x_2)$ as follows.

- 1) (Coercive) In this test, we take the components of \mathbf{A} such that the differential operator can be written in variational form and is coercive, fitting into a standard analytical framework:

$$a(\mathbf{x}) = 1 - \ln\left((x_1 - 1/2)^2 + 10^{-4}\right), \tag{87}$$

$$b(\mathbf{x}) = 0. \tag{88}$$

- 2) (Continuous not H^2) In this test, we take \mathbf{A} such that it is comparable to [32, §4.4]

$$a(\mathbf{x}) = 2, \tag{89}$$

$$b(\mathbf{x}) = (x_1^2 x_2^2)^{1/3}. \tag{90}$$

- 3) (Discontinuous not $W^{1,\infty}$) In our third test, a is discontinuous and not in $W^{1,\infty}$. We take $b \equiv 0$ and choose

$$a(\mathbf{x}) = \alpha(|\mathbf{x}|_\infty) \tag{91}$$

with

$$\alpha(s) = \frac{1}{100} + 1000 \begin{cases} \sqrt{\frac{1}{4} - s}, & \text{if } s < \frac{1}{4}, \\ \cos \pi s, & \text{otherwise.} \end{cases} \tag{92}$$

Note that the initial grid is chosen so that it aligns with the discontinuity.

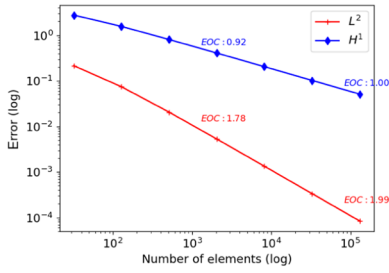
We study the behaviour of our method for polynomial degrees $k = 1, 2, 3$, choosing the forcing such that the exact solution is given by the following two choices.

- 1) (Smooth solution): $u(\mathbf{x}) = \sin 2\pi x_1 \sin 2\pi x_2$.
 2) ($H^2 \setminus H^3$ solution): $u(\mathbf{x}) = \begin{cases} \frac{1}{4} \left(\cos 8\pi \left| \mathbf{x} - \frac{1}{2} \right|^2 + 1 \right), & \text{if } \left| \mathbf{x} - \frac{1}{2} \right|^2 \leq \frac{1}{8}, \\ 0, & \text{otherwise.} \end{cases}$

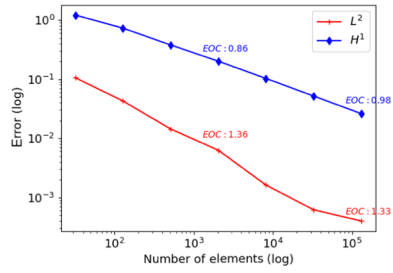
The penalty parameter is chosen as $\sigma = \lambda_{A,\max} 5k(k+1)$ where $\lambda_{A,\max}$ is the maximum eigenvalue of \mathbf{A} .

Remark 8 (Compatibility of spaces for u_h and $\mathbf{H}(u_h)$) Note that it is not actually required that the approximations u_h and $\mathbf{H}(u_h)$ are represented in the same finite element space. Computationally we observe similar results when $\mathbf{H}(u_h)$ is one degree lower than u_h and suboptimal convergence rates when $\mathbf{H}(u_h)$ is two orders lower than u_h . Other choices do not appear to be stable.

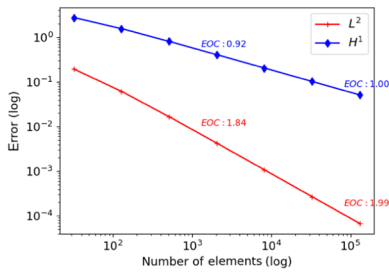
Benchmarking for these tests is shown in Figs. 1, 2 and 3. For linear polynomials, the $H^2(\mathcal{T})$ -error does not converge since the piecewise Hessian of u_h vanishes. So we only show the errors in the $L^2(\Omega)$ and $H^1(\mathcal{T})$ norms. For the smooth solution (left column), we clearly see that the method converges optimally in both norms. While for the less smooth solution convergence is slowed as would be expected when approximating an $H^2(\Omega)$ solution. For polynomial orders 2 and 3 convergence is optimal for all norms studied here when approximating a smooth solution and between 1 and 1.5 for



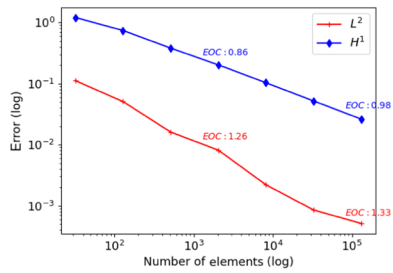
(a) $u \in C^\infty(\Omega)$ and A is coercive



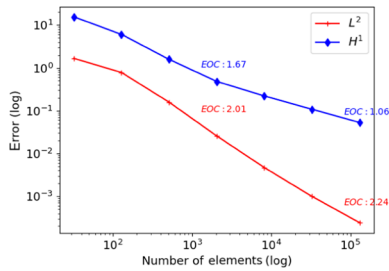
(b) $u \in H^2(\Omega)/H^3(\Omega)$ and A is coercive



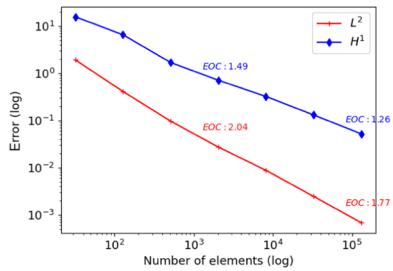
(c) $u \in C^\infty(\Omega)$ and A is continuous but not $H^2(\Omega)$



(d) $u \in H^2(\Omega)/H^3(\Omega)$ and A is continuous but not $H^2(\Omega)$



(e) $u \in C^\infty(\Omega)$ and A is discontinuous



(f) $u \in H^2(\Omega)/H^3(\Omega)$ and A is discontinuous

Fig. 1 Convergence rates for the error measured in $H^1(\Omega)$ and $L^2(\Omega)$ for $k = 1$ testing the case when u is either a prescribed smooth solution or $u \in H^2(\Omega)/H^3(\Omega)$. We also test three different choices for diffusion coefficient A that are coercive, continuous but not H^2 and discontinuous. The rates are optimal in $H^1(\Omega)$ in all cases and in $L^2(\Omega)$ in most cases. Note we did not present the $H^2(\Omega)$ rates nor the estimate since neither will converge for $k = 1$

the $H^2(\Omega)$ solution. The convergence rate hardly depends on the smoothness of A . The only case where some clear dependency is visible is in the $L^2(\Omega)$ errors for quadratic polynomials with the discontinuous A . In this case, the order of the $L^2(\Omega)$ error seems to equal the convergence in the $H^1(\mathcal{T})$ norm, i.e., is not optimal (see Fig. 2).

We also study the condition number of the system matrix generated when assembling (49). The results are shown in Fig. 4. In contrast to many other methods, which rewrite the

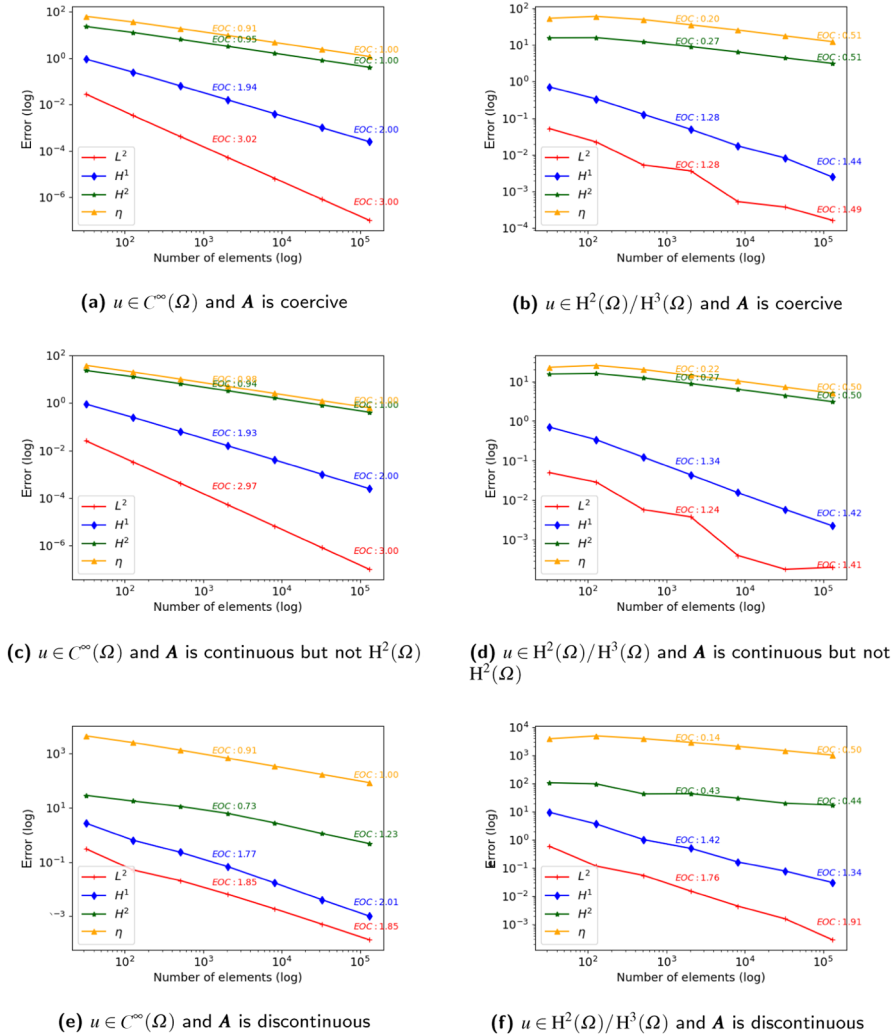
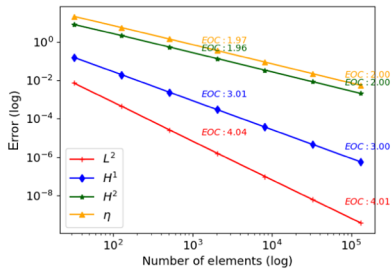


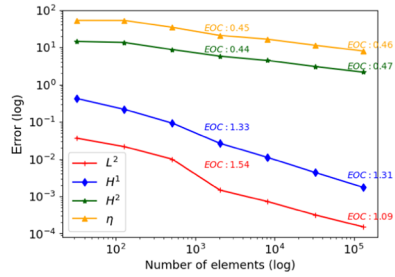
Fig. 2 Convergence rates for the error measured in $H^2(\Omega)$, $H^1(\Omega)$ and $L^2(\Omega)$ for $k = 2$ and the estimator given in Theorem 3. We test the case when u is either a prescribed smooth solution or $u \in H^2(\Omega)/H^3(\Omega)$. We also test three different choices for diffusion coefficient A that are coercive, continuous but not H^2 and discontinuous. The rates are optimal in $H^2(\Omega)$ and $H^1(\Omega)$ in all cases and in $L^2(\Omega)$ in most cases. When the solution is not smooth the rates are slowed to an expected rate in-line with the regularity of u . In all cases, the estimate is efficient and robust

nonvariational model as a fourth-order problem, the condition number depends on the grid spacing in the same way as it does for second-order variational problem, i.e., it is $O(h^{-2})$.

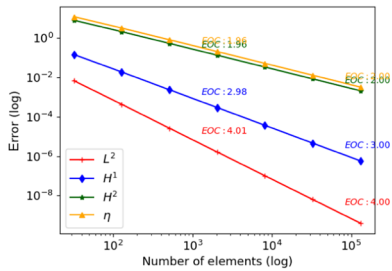
Finally, we study the behaviour of the residual error indicator and an adaptive scheme for the approximation of the solutions to these problems. The indicator is included in Figs. 2, 3 where its reliability is clearly visible. A comparison of these globally refined simulations with an adaptive simulation using an equal distribution strategy for locally refining the grid are given in Figs. 5 and 6. As to be expected, no advantage can be gained when approximating a smooth solution.



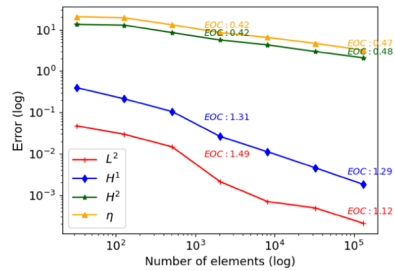
(a) $u \in C^\infty(\Omega)$ and A is coercive



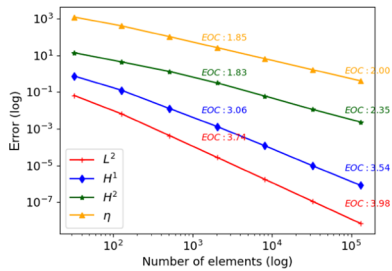
(b) $u \in H^2(\Omega)/H^3(\Omega)$ and A is coercive



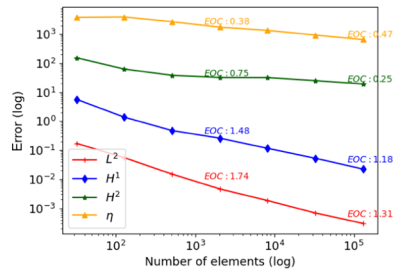
(c) $u \in C^\infty(\Omega)$ and A is continuous but not $H^2(\Omega)$



(d) $u \in H^2(\Omega)/H^3(\Omega)$ and A is continuous but not $H^2(\Omega)$



(e) $u \in C^\infty(\Omega)$ and A is discontinuous

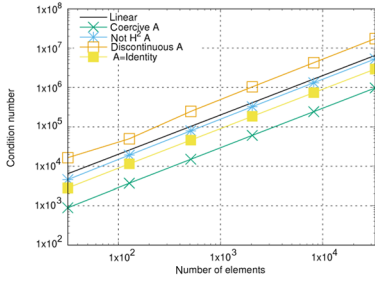


(f) $u \in H^2(\Omega)/H^3(\Omega)$ and A is discontinuous

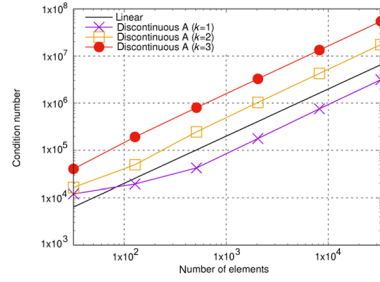
Fig. 3 Convergence rates for the error measured in $H^2(\Omega)$, $H^1(\Omega)$ and $L^2(\Omega)$ for $k = 3$ and the estimator given in Theorem 3. We test the case when u is either a prescribed smooth solution or $u \in H^2(\Omega)/H^3(\Omega)$. We also test three different choices for diffusion coefficient A that are coercive, continuous but not H^2 and discontinuous. The rates are optimal in $H^2(\Omega)$ and $H^1(\Omega)$ in all cases and in $L^2(\Omega)$ in most cases. When the solution is not smooth the rates are slowed to an expected rate in-line with the regularity of u . In all cases, the estimate is efficient and robust.

For the $H^2(\Omega)$ solution, the adaptive simulation approximately doubles the convergence rate of the scheme.

We conclude with a simulation for which we do not have an exact solution. We use the discontinuous A and choose a constant forcing $f \equiv 1000$. As before boundary conditions are equal to zero. Results are shown in Fig. 6 where we show the convergence of the residual indicator under global and local refinement. A visualisation of the function a , the resulting discrete solution, and the adaptive grid are given in Fig. 7.

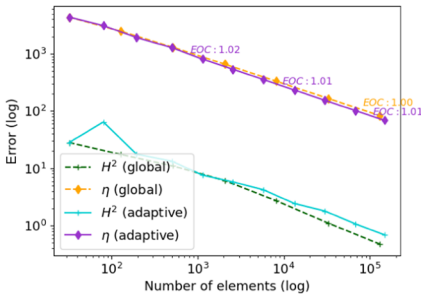


(a) Condition numbers for each of the operators (5.1). For comparison purposes we included the case of the Laplace problem ($A = I$) where our method corresponds to the IP-DG method

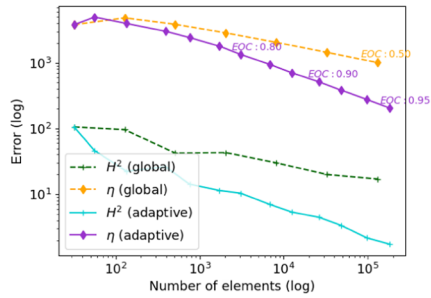


(b) Condition numbers for discontinuous coefficient A (91) with various polynomial degrees $k = 1, 2, 3$

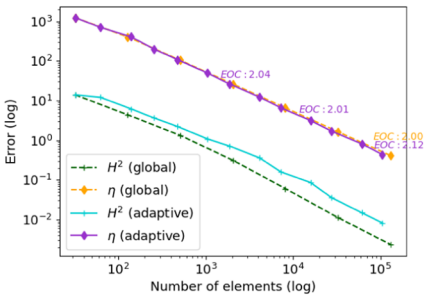
Fig. 4 Condition number estimates of the system matrix given by (49). Notice that the condition number grows in the same asymptotic fashion as the IP-DG method for the Laplacian ($A = Identity$). Also, the complexity does not change asymptotically as k is increased



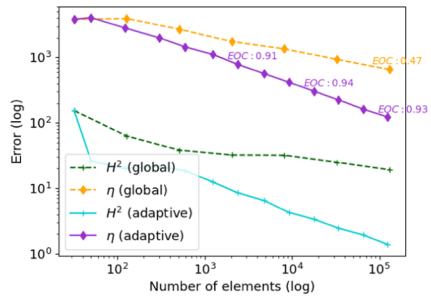
(a) $u \in C^\infty(\Omega)$ and $k = 2$



(b) $u \in H^2(\Omega)/H^3(\Omega)$ and $k = 2$

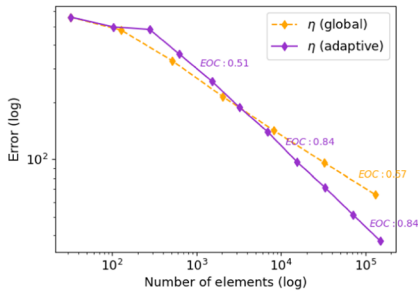


(c) $u \in C^\infty(\Omega)$ and $k = 3$

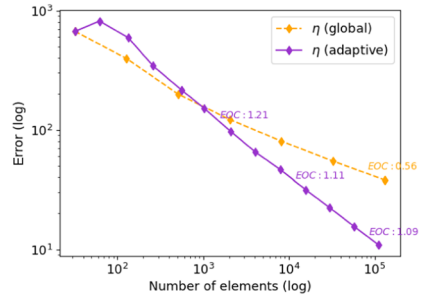


(d) $u \in H^2(\Omega)/H^3(\Omega)$ and $k = 3$

Fig. 5 Comparison of simulations carried out on adaptive and globally refined grids driven by the estimate from Theorem 3. We test the case when u is either a prescribed smooth solution or $u \in H^2(\Omega)/H^3(\Omega)$ and when A is discontinuous. For the smooth solution, we find the uniform and adaptive schemes behave similarly, for the solution that is not $H^3(\Omega)$, the adaptive scheme clearly outperforms the uniform one, although there is no gain from using $k = 3$

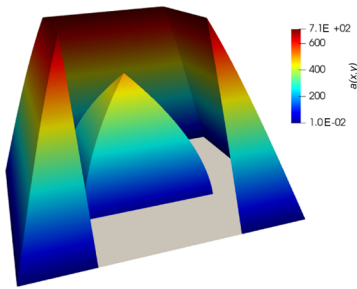


(a) Indicator with constant forcing and $k = 2$

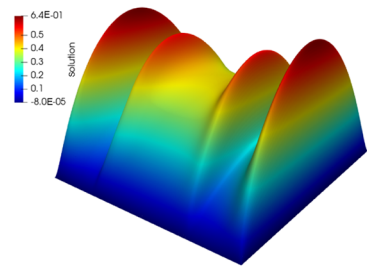


(b) Indicator with constant forcing and $k = 3$

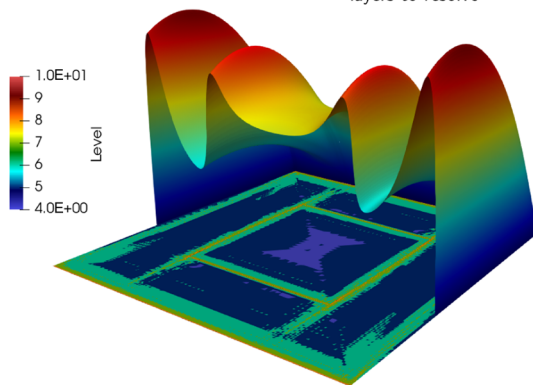
Fig. 6 Convergence of estimator under adaptive and globally refined grids. Here, we consider the case when A is discontinuous, shown in Fig. 7a, and the forcing is chosen $f = 100$. No exact solution is known for this case. There is a clear indication that the adaptive scheme outperforms the uniform one



(a) A visualisation of $a(x)$ where a slice has been cut out to show the discontinuity. The figure is scaled by 10^{-3}



(b) The adaptive approximation to the solution. Note there are boundary and interior layers to resolve



(c) The adaptive approximation has been sliced and overlaid on the underlying adaptive mesh. Colour indicates refinement level from the initial macro triangulation

Fig. 7 Visualisation of the problem coefficient, the solution and the underlying adaptive mesh refinement level. Notice the algorithm refines where the problem data are discontinuous and well as near boundary layers caused by the anisotropy of the diffusion tensor

6 Conclusions and Outlook

In this work, we have extended the framework from [32] for linear nonvariational problems to incorporate discontinuous approximations. We have derived a posteriori bounds for this problem and shown they are useful to drive adaptive algorithms. We would like to point out that the approach presented here can be directly applied to the case where the approximation u_h is chosen in a continuous finite element space but the finite element Hessian is defined in the discontinuous fashion described here.

In the numerical experiments, we note the method is well posed and converges optimally even for A that are discontinuous. The method is well suited to solve nonlinear problems [33] and this will be the topic of ongoing research.

Compliance with Ethical Standards

Conflict of Interest On behalf of all authors, the corresponding author states that there is no conflict of interest.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Agouzal, A., Vassilevski, Y.: On a discrete Hessian recovery for P_1 finite elements. *J. Numer. Math.* **10**(1), 1–12 (2002)
2. Aguilera, N.E., Morin, P.: On convex functions and the finite element method. *SIAM J. Numer. Anal.* **47**(4), 3139–3157 (2009)
3. Alnæs, M.S., Logg, A., Ølgaard, K.B., Rognes, M.E., Wells, G.N.: Unified form language: a domain-specific language for weak formulations of partial differential equations. CoRR (2012). arXiv:1211.4047
4. Arnold, D.N., Brezzi, F., Cockburn, B., Marini, L.D.: Unified analysis of discontinuous Galerkin methods for elliptic problems. *SIAM J. Numer. Anal.* **39**(5), 1749–1779 (2002)
5. Barles, G., Souganidis, P.E.: Convergence of approximation schemes for fully nonlinear second order equations. *Asymptot. Anal.* **4**(3), 271–283 (1991)
6. Bastian, P., Blatt, M., Dedner, A., Engwer, C., Klöforn, R., Ohlberger, M., Sander, O.: A generic grid interface for parallel and adaptive scientific computing. I. Abstract framework. *Computing* **82**(2/3), 103–119 (2008)
7. Bastian, P., Blatt, M., Dedner, A., Engwer, C., Klöforn, R., Kornhuber, R., Ohlberger, M., Sander, O.: A generic grid interface for parallel and adaptive scientific computing. II. Implementation and tests in DUNE. *Computing* **82**(2/3), 121–138 (2008)
8. Böhmer, K.: On finite element methods for fully nonlinear elliptic equations of second order. *SIAM J. Numer. Anal.* **46**(3), 1212–1249 (2008)
9. Brenner, S.C., Sung, L.-Y.: Virtual enriching operators. *Calcolo* **56**(4), 1–25 (2019)
10. Buffa, A., Ortner, C.: Compact embeddings of broken Sobolev spaces and applications. *IMA J. Numer. Anal.* **29**(4), 827–855 (2009)
11. Burman, E., Ern, A.: Discontinuous Galerkin approximation with discrete variational principle for the nonlinear Laplacian. *C. R. Math. Acad. Sci. Paris* **346**(17/18), 1013–1016 (2008)
12. Cordes, H.O.: Über die erste randwertaufgabe bei quasilinearen differentialgleichungen zweiter ordnung in mehr als zwei variablen. *Math. Ann.* **131**(3), 278–312 (1956)

13. Dedner, A., Klöforn, R., Nolte, M.: Python bindings for the DUNE-FEM module. Zenodo (2020). <https://doi.org/10.5281/zenodo.3706993>
14. Dedner, A., Klöforn, R., Nolte, M., Ohlberger, M.: A generic interface for parallel and adaptive scientific computing: abstraction principles and the DUNE-FEM module. *Computing* **90**, 165–196 (2010)
15. Dedner, A., Nolte, M.: The Dune-Python Module. *CoRR* (2018). [arXiv:1807.05252](https://arxiv.org/abs/1807.05252)
16. Douglas, J. Jr., Dupont, T.: Interior penalty procedures for elliptic and parabolic Galerkin methods. In: *Computing Methods in Applied Sciences (Second Internat. Sympos., Versailles, 1975)*. Lecture Notes in Phys., vol. 58, pp. 207–216. Springer, Berlin (1976)
17. Di Pietro, D.A., Ern, A.: Discrete functional analysis tools for discontinuous Galerkin methods with application to the incompressible Navier-Stokes equations. *Math. Comput.* **79**(271), 1303–1330 (2010)
18. Elman, H.C., Silvester, D.J., Wathen, A.J.: Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics. In: *Numerical Mathematics and Scientific Computation*. Oxford University Press, New York (2005)
19. Ern, A., Guermond, J.-L.: Theory and practice of finite elements. In: Antman, S.S., Marsden, J.E., Sirovich, L. (eds) *Applied Mathematical Sciences*, vol. 159. Springer, New York (2004)
20. Feng, X., Neilan, M.: Mixed finite element methods for the fully nonlinear Monge-Ampère equation based on the vanishing moment method. *SIAM J. Numer. Anal.* **47**(2), 1226–1250 (2009)
21. Feng, X., Neilan, M.: Vanishing moment method and moment solutions for fully nonlinear second order partial differential equations. *J. Sci. Comput.* **38**(1), 74–98 (2009)
22. Feng, X., Hennings, L., Neilan, M.: Finite element methods for second order linear elliptic partial differential equations in non-divergence form. *Math. Comput.* **86**(307), 2025–2051 (2017)
23. Feng, X., Neilan, M., Schnake, S.: Interior penalty discontinuous Galerkin methods for second order linear non-divergence form elliptic PDES. *J. Sci. Comput.* **74**(3), 1651–1676 (2018)
24. Gallistl, D.: Variational formulation and numerical analysis of linear elliptic equations in nondivergence form with Cordes coefficients. *SIAM J. Numer. Anal.* **55**(2), 737–757 (2017)
25. Gallistl, D.: Numerical approximation of planar oblique derivative problems in nondivergence form. *Math. Comput.* **88**(317), 1091–1119 (2019)
26. Georgoulis, E.H., Houston, P., Virtanen, J.: An a posteriori error indicator for discontinuous Galerkin approximations of fourth-order elliptic problems. *IMA J. Numer. Anal.* **31**(1), 281–298 (2011)
27. Gilbarg, D., Trudinger, N.S.: *Elliptic Partial Differential Equations of Second Order*, 2nd edn. Springer, Berlin (1983)
28. Jensen, M., Smears, I.: On the convergence of finite element methods for Hamilton-Jacobi-Bellman equations. *Technical report*, 01 (2011)
29. Kawecki, E.L.: A DGFEM for nondivergence form elliptic equations with Cordes coefficients on curved domains. *Numer. Methods Partial Differ. Equations* **35**(5), 1717–1744 (2019)
30. Kawecki, E.L., Smears, I.: Convergence of adaptive discontinuous Galerkin and c^0 -interior penalty finite element methods for Hamilton-Jacobi-Bellman and Isaacs equations (2020). [arXiv:2006.07215](https://arxiv.org/abs/2006.07215)
31. Lakkis, O., Mousavi, A.: A least-squares Galerkin approach to gradient and Hessian recovery for non-divergence-form elliptic equations (2019). [arXiv:1909.00491](https://arxiv.org/abs/1909.00491)
32. Lakkis, O., Pryer, T.: A finite element method for second order nonvariational elliptic problems. *SIAM J. Sci. Comput.* **33**(2), 786–801 (2011)
33. Lakkis, O., Pryer, T.: A nonvariational finite element method for fully nonlinear elliptic problems. *Submitted–Tech report* (2012). [arXiv:1103.2970](https://arxiv.org/abs/1103.2970)
34. Miranda, C.: Sulle equazioni ellittiche del secondo ordine di tipo non variazionale, a coefficienti discontinui. *Ann. Mat.* **63**(1), 353–386 (1963)
35. Mu, L., Ye, X.: A simple finite element method for non-divergence form elliptic equations. *Int. J. Numer. Anal. Model.* **14**(2), 306–311 (2017)
36. Oberman, A.M.: Convergent difference schemes for degenerate elliptic and parabolic equations: Hamilton-Jacobi equations and free boundary problems. *SIAM J. Numer. Anal.* **44**(2), 879–895 (2006) (electronic)
37. Pryer, T.: Discontinuous Galerkin methods for the p -biharmonic equation from a discrete variational perspective. *Electron. Trans. Numer. Anal.* **41**, 328–349 (2014)
38. Smears, I., Süli, E.: Discontinuous Galerkin finite element approximation of nondivergence form elliptic equations with Cordès coefficients. *SIAM J. Numer. Anal.* **51**(4), 2088–2106 (2013)
39. Talenti, G.: Sopra una classe di equazioni ellittiche a coefficienti misurabili. *Ann. Mat.* **69**(1), 285–304 (1965)
40. Vallet, M.-G., Manole, C.-M., Dompierre, J., Dufour, S., Guibault, F.: Numerical comparison of some Hessian recovery techniques. *Int. J. Numer. Methods Eng.* **72**(8), 987–1007 (2007)