



HAL
open science

Improved Error estimates of hybridizable interior penalty methods using a variable penalty for highly anisotropic diffusion problems

Grégory Etangsale, Marwan Fahs, Vincent Fontaine, Nalitiana Rajaonison

► To cite this version:

Grégory Etangsale, Marwan Fahs, Vincent Fontaine, Nalitiana Rajaonison. Improved Error estimates of hybridizable interior penalty methods using a variable penalty for highly anisotropic diffusion problems. 2020. hal-02893064v1

HAL Id: hal-02893064

<https://hal.science/hal-02893064v1>

Preprint submitted on 8 Jul 2020 (v1), last revised 20 Jun 2021 (v2)

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

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

Error estimates of hybridizable interior penalty methods using a variable penalty for highly anisotropic diffusion problems

Grégory Etangsale^a, Marwan Fahs^b, Vincent Fontaine^{a,*}, Nalitiana Rajaonison^a

^aDepartment of Building and Environmental Sciences, University of La Réunion, France

^bUniversité de Strasbourg, CNRS, ENGEES, LHYGES UMR 7517, F-67000 Strasbourg, France

Abstract

In this paper, we derive improved *a priori* error estimates for families of hybridizable interior penalty discontinuous Galerkin (H-IP) methods using a variable penalty for second-order elliptic problems. The strategy is to use a penalization function of the form $\mathcal{O}(1/h^{1+\delta})$, where h denotes the mesh size and δ is a user-dependent parameter. We then quantify its direct impact on the convergence analysis, namely, the (strong) consistency, discrete coercivity and boundedness (with h^δ -dependency), and we derive updated error estimates for both discrete energy- and L^2 -norms. All theoretical results are supported by numerical evidence.

Keywords: Hybridizable discontinuous Galerkin, interior penalty methods, variable-penalty technique, convergence analysis, updated *a priori* error estimates

2020 MSC: 65N12, 65N15, 65N30, 65N38

1. Introduction

Hybridizable discontinuous Galerkin (HDG) methods were first introduced in the last decade by Cockburn *et al.* [1] (see, e.g., [2]) and have since received extensive attention from the research community. They are popular and very efficient numerical approaches for solving a large class of partial differential equations (see, e.g., [3, 4, 5, 6, 7] for a historical perspective). Indeed, they inherit attractive features from both (i) discontinuous Galerkin (DG) methods such as local conservation, *hp*-adaptivity and high-order polynomial approximation [8] and (ii) standard conforming Galerkin (CG) methods such as the Schur complement strategy [9]. One undeniable additional benefit of the HDG methods is their superconvergence property, obtained through the application of a local postprocessing technique on each element of the mesh [4]. In the hybrid formalism, additional unknowns are introduced along the mesh skeleton corresponding to discrete trace approximations. Thanks to the specific localization of its additional degrees of freedom (dofs), interior variables can be eliminated in favor of its Lagrange multipliers by only static condensation [10]. The resulting matrix system is significantly smaller and sparser than those associated with CG or DG methods for any given mesh and polynomial degree [9]. Several HDG formulations have been derived in the literature and can be classified into two main categories. The first is based on a primal form of the continuous problem, such as the class of interior penalty (IP) methods [11], whereas the second relies on a dual (often called mixed) form, such as local discontinuous Galerkin (LDG) methods [1, 4, 12]. In the latter formulation, the flux variable is introduced as an additional unknown of the problem.

Our focus is on families of hybridizable interior penalty discontinuous Galerkin (H-IP) methods [13]. They are hybridized counterparts of the well-known interior penalty DG (IPDG) methods [14, 15, 16] and have been analyzed until quite recently by several authors [11, 6]. Specifically, in our exposition, we considered the incomplete, non-symmetric and symmetric schemes denoted by H-IIP, H-NIP and H-SIP, respectively. The main difference between these schemes concerns the role of the *symmetrization* term in the discrete bilinear form [15]. Fabien *et al.* recently

*Corresponding author

Email address: vincent.fontaine@univ-reunion (Vincent Fontaine)

38 analyzed these schemes using a stabilization function of the form $O(1/h)$ for solving second-order elliptic problems
39 [11]. The authors conclude that H-IP methods inherit similar convergence properties to their IPDG equivalents.
40 Notably, they theoretically establish (i) optimal energy error estimates, and because of the lack of symmetry of the
41 associated discrete operator, (ii) only suboptimal L^2 -norm error estimates for H-IIP and H-NIP schemes. In addition,
42 they numerically conclude that the L^2 orders of convergence of both non-symmetric variants are suboptimal for only
43 even polynomial degrees and are optimal otherwise. Similar conclusions have also been suggested by Oikawa for
44 second-order elliptic problems [5].

45 To restore optimal L^2 -error estimates for the nonsymmetric IPDG method, Rivière *et al.* suggest using a sort of
46 superpenalty on the jumps [17, 18]. In the present paper, we explore a similar idea in the general context of H-IP
47 methods by using a variable penalty function of the form $\tau := O(1/h^{1+\delta})$, where $\delta \in \mathbb{R}$. Here, we analyze the direct
48 impact of the parameter δ on *a priori* error estimates in different norms. First, we propose a convergence analysis by
49 investigating three key properties: (strong) consistency, discrete coercivity and boundedness. One remarkable feature
50 of this strategy is the h^δ -dependency of the coercivity condition and the continuity (or boundedness) constant C_{bnd} ,
51 which consequently impacts the error estimates. Improved error estimates are then derived in the spirit of the second
52 Strang lemma [16], and we first prove that the order of convergence in the natural energy-norm is linear, δ -dependent,
53 and optimal when $\delta \geq 0$ for any scheme. Then, by using a duality argument, i.e., the so-called Aubin–Nitsche
54 technique, we also prove that the optimal convergence is theoretically reached as soon as $\delta \geq 0$ for the H-SIP scheme
55 only, and when $\delta \geq 2$ for both non-symmetric variants, i.e., H-NIP and H-IIP schemes. We recover some well-known
56 theoretical error estimates proposed in the literature for both the natural energy- and L^2 -norms in the particular case
57 of $\delta = 0$.

58 The rest of the material is organized as follows: Section 2 describes the model problem, mesh notation and
59 assumptions, and recalls some definitions and useful (trace) inequalities, while Section 3 derives the discrete H-IP
60 formulation and discusses its stability properties. In Section 4, optimal error estimates are provided for both the
61 energy- and L^2 -norms by using a standard duality argument. Section 5 concerns the numerical experiments that
62 validate our theoretical results. We briefly end with some remarks and perspectives.

63 2. Some preliminaries

64 2.1. The model problem

65 Let Ω be a bounded (polyhedron) domain in \mathbb{R}^d with Lipschitz boundary $\partial\Omega$ in spatial dimension $d \geq 2$. For
66 clarity, we consider the anisotropic diffusion problem with homogeneous Dirichlet boundary conditions:

$$-\nabla \cdot (\kappa \nabla u) = f \quad \text{in } \Omega \quad \text{and} \quad u = 0 \quad \text{on } \partial\Omega, \quad (1)$$

67 where $\kappa \in [L^\infty(\Omega)]^{d \times d}$ is a bounded, symmetric, uniformly positive-definite matrix-valued function and $f \in L^2(\Omega)$ is
68 a forcing term. Thus, the weak formulation of problem (1) is to find $u \in H_0^1(\Omega)$ such that

$$\int_{\Omega} \kappa \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx \quad \forall v \in H_0^1(\Omega). \quad (2)$$

69 It is well known that under elliptic regularity assumptions, the variational problem (2) is well posed.

70 2.2. Mesh notation and assumptions

71 Let h be a positive parameter; we assume without loss of generality that $h \leq 1$. We denote by $\{\mathcal{T}_h\}_{h>0}$ a family of
72 affine triangulations of the domain Ω , where h stands for the largest diameter: $h_E := \text{diam}(E)$. We also assume that
73 \mathcal{T}_h is *quasi-uniform*, meaning that for all $E \in \mathcal{T}_h$, there exists $0 < \rho_0 \leq 1$ independent of h such that $\rho_0 h \leq h_E \leq h$.
74 Following our notation, the generic term *interface* indicates a $(d-1)$ -dimensional geometric object, i.e., an edge, if
75 $d = 2$ and a face if $d = 3$. Thus, we denote by \mathcal{F}_h^i the set of interior interfaces; i.e., $F \in \mathcal{F}_h^i$ if there exist E_1 and E_2
76 in \mathcal{T}_h such that $F := \partial E_1 \cap \partial E_2$. The set of boundary interfaces is denoted by \mathcal{F}_h^b ; i.e., $F \in \mathcal{F}_h^b$ if there exists E in
77 \mathcal{T}_h such that $F := \partial E \cap \partial\Omega$. The set of all interfaces is often called the mesh skeleton and is denoted by \mathcal{F}_h , i.e.,
78 $\mathcal{F}_h := \mathcal{F}_h^i \cup \mathcal{F}_h^b$. We denote by $\partial\mathcal{T}_h := \{\cup \partial E, \forall E \in \mathcal{T}_h\}$, the collection of interfaces of all mesh elements. Let X be
79 a mesh element or an interface; we then denote by $|X|$ a positive d - or $(d-1)$ -dimensional Lebesgue measure of X ,

80 respectively. Moreover, for any mesh element $E \in \mathcal{T}_h$, we denote by $\mathcal{F}_E := \{F \in \mathcal{F}_h : F \subset \partial E\}$ the set of interfaces
 81 composing the boundary of E ; we define $\eta_E := \text{card}(\mathcal{F}_E)$ and $\eta_0 := \max_{E \in \mathcal{T}_h} (\eta_E)$.

82 2.3. Broken polynomial spaces

83 For any polyhedral domain $D \subset \mathbb{R}^d$ with $\partial D \subset \mathbb{R}^{d-1}$, we denote by $(\cdot, \cdot)_{0,D}$ (resp., $\langle \cdot, \cdot \rangle_{0,\partial D}$) the L^2 -inner product
 84 in $L^2(D)$ (resp., $L^2(\partial D)$) equipped with its natural norm $\|\cdot\|_{0,D}$ (resp., $\|\cdot\|_{0,\partial D}$). Let us now introduce some compact
 85 notation associated with the discrete L^2 -inner scalar product:

$$(\cdot, \cdot)_{0,\mathcal{T}_h} := \sum_{E \in \mathcal{T}_h} (\cdot, \cdot)_{0,E} \quad \text{and} \quad \langle \cdot, \cdot \rangle_{0,\partial\mathcal{T}_h} := \sum_{E \in \mathcal{T}_h} \langle \cdot, \cdot \rangle_{0,\partial E}. \quad (3)$$

86 We denote by $\|\cdot\|_{0,\mathcal{T}_h}$ and $\|\cdot\|_{0,\partial\mathcal{T}_h}$ the corresponding norms. Similarly, we denote by $H^s(D)$ the usual Hilbert space
 87 of index s on D equipped with its natural norm $\|\cdot\|_{s,D}$ and seminorm $|\cdot|_{s,D}$, respectively. If $s = 0$, then we set
 88 $H^0(D) = L^2(D)$. We denote by $H^s(\mathcal{T}_h)$ the usual broken Sobolev space and by $\mathbf{\nabla}_h$ the broken gradient operator acting
 89 on $H^s(\mathcal{T}_h)$ with $s \geq 1$. We assume an extended regularity requirement of the exact solution u of the weak problem (2),
 90 i.e., $u \in H_0^s(\Omega) \cap H^2(\mathcal{T}_h)$ with $s > 3/2$. We also introduce the additional unknown $\hat{u} \in L^2(\mathcal{F}_h)$ corresponding to the
 91 trace of u on the skeleton of the mesh. Let us now introduce the composite variable $\mathbf{u} := (u, \hat{u})$, which belongs to the
 92 continuous approximation space $\mathbf{V} := H_0^s(\Omega) \cap H^2(\mathcal{T}_h) \times L^2(\mathcal{F}_h)$; i.e., $\mathbf{u} \in \mathbf{V}$. As usual in HDG methods, we consider
 93 broken Sobolev spaces:

$$\mathbb{P}_k(\mathcal{T}_h) := \{v_h \in L^2(\mathcal{T}_h) : v_h|_E \in \mathbb{P}_k(E), \forall E \in \mathcal{T}_h\}, \quad (4)$$

94 and similarly for $\mathbb{P}_k(\mathcal{F}_h)$. Here, $\mathbb{P}_k(X)$ denotes the space of polynomials of at least degree k on X , where X corresponds
 95 to a generic element of \mathcal{T}_h or \mathcal{F}_h , respectively. For H-IP discretization, two types of discrete variables are necessary
 96 to approximate the weak solution u of problem (2). First, the discrete variable $u_h \in V_h$ is defined within each mesh
 97 element, and its trace $\hat{u}_h \in \hat{V}_h$ is defined on the mesh skeleton with respect to the imposed homogeneous Dirichlet
 98 boundary conditions. Thus, we set $V_h := \mathbb{P}_k(\mathcal{T}_h)$ and $\hat{V}_h := \mathbb{P}_k^0(\mathcal{F}_h)$, where

$$\mathbb{P}_k^0(\mathcal{F}_h) := \{\hat{v}_h \in \mathbb{P}_k(\mathcal{F}_h) : \hat{v}_h|_F = 0, \forall F \in \mathcal{F}_h^b\}. \quad (5)$$

99 Throughout the manuscript, we use the following compact notation: Let $\mathbf{V}_h := V_h \times \hat{V}_h$ denote the composite approx-
 100 imation space and a generic element of \mathbf{V}_h be denoted by $\mathbf{v}_h := (v_h, \hat{v}_h)$. For all $E \in \mathcal{T}_h$ and $F \in \mathcal{F}_E$, we define the
 101 jump of $v_h \in V_h$ across F as $\llbracket v_h \rrbracket_{E,F} := (v_h|_E - \hat{v}_h|_F) \mathbf{n}_F$, where \mathbf{n}_F denotes the unit normal vector to F pointing
 102 out of E . When confusion cannot arise, we omit the subscripts E and F from the definition, and we simply write
 103 $\llbracket v_h \rrbracket := (v_h - \hat{v}_h) \mathbf{n}$. Finally, we introduce the space $\mathbf{V}(h) := \mathbf{V} + \mathbf{V}_h$ to analyze the boundedness of the discrete bilinear
 104 form.

105 2.4. Useful inequalities

106 We recall here some useful inequalities that will be used extensively later on (see, e.g., [19, 16, 15]). For clarity,
 107 C denotes a generic constant that is independent of h , h_E and κ in the rest of the manuscript. Owing to the shape
 108 regularity of \mathcal{T}_h , we now introduce multiplicative trace inequalities. Let $E \in \mathcal{T}_h$ and $F \in \mathcal{F}_E$. For all $v \in H^2(E)$, there
 109 exists a positive constant C_M independent of h_E , v and E such that

$$\|v\|_{0,F}^2 \leq C_M (\|v\|_{0,E} \|v\|_{1,E} + h_E^{-1} \|v\|_{0,E}^2), \quad (6a)$$

$$\|\mathbf{\nabla}_h v\|_{0,F}^2 \leq C_M (\|v\|_{1,E} \|v\|_{2,E} + h_E^{-1} \|v\|_{1,E}^2). \quad (6b)$$

110 On broken polynomial spaces $v_h \in V_h$, we obtain the discrete and inverse trace inequalities, respectively:

$$\|v_h\|_{0,F} \leq C_{\text{tr}} h_E^{-1/2} \|v_h\|_{0,E}, \quad (7a)$$

$$\|\mathbf{\nabla}_h v_h\|_{0,E} \leq C_{\text{inv}} h_E^{-1} \|v_h\|_{0,E}, \quad (7b)$$

111 where C_{tr} and C_{inv} are positive constants independent of h_E .

112 **Remark 2.1.** Following Rivière [15] (see Section 2.1.3, p.24), one can obtain an exact expression of the constant C_{tr}
 113 used in the discrete trace inequality (7a) for a d -simplex mesh element:

$$C_{\text{tr}} := \sqrt{\frac{(k+1)(k+d)}{d}}, \quad (8)$$

114 where k denotes the polynomial degree of V_h and d denotes the spatial dimension. This expression is particularly
 115 important in our analysis since it will be used later in the definition of the penalty parameter.

116 We are now in a position to introduce the energy-norm used in the stability analysis and error estimations [13, 10].
 117 For any given composite function $\mathbf{v}_h \in \mathbf{V}_h$, we consider the jump seminorm:

$$|\mathbf{v}_h|_{\gamma}^2 := \sum_{E \in \mathcal{T}_h} |\mathbf{v}_h|_{\gamma, \partial E}^2 \quad \text{with} \quad |\mathbf{v}_h|_{\gamma, \partial E} := \sum_{F \in \mathcal{F}_E} \|\gamma_F^{1/2} \llbracket \mathbf{v}_h \rrbracket\|_{0,F}^2, \quad (9)$$

118 where $\gamma_F \geq 0$ is an arbitrary positive constant associated with $F \in \mathcal{F}_E$. The natural energy-norm equipping the
 119 discrete approximation space \mathbf{V}_h is given by

$$\|\mathbf{v}_h\|_*^2 := \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v}_h\|_{0, \mathcal{T}_h}^2 + |\mathbf{v}_h|_{\gamma}^2, \quad (10)$$

120 which clearly depends on $\boldsymbol{\kappa}$.

121 3. Hybridizable interior penalty methods

122 The discrete H-IP problem is to find $\mathbf{u}_h \in \mathbf{V}_h$ such that

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{u}_h, \mathbf{v}_h) = l(\mathbf{v}_h), \quad \forall \mathbf{v}_h \in \mathbf{V}_h, \quad (11)$$

123 where $l(\mathbf{v}_h) := (f, \mathbf{v}_h)_{0, \mathcal{T}_h}$ and the bilinear form $\mathcal{B}_h^{(\epsilon)} : \mathbf{V}_h \times \mathbf{V}_h \rightarrow \mathbb{R}$ is given by

$$\begin{aligned} \mathcal{B}_h^{(\epsilon)}(\mathbf{u}_h, \mathbf{v}_h) := & (\boldsymbol{\kappa} \nabla_h \mathbf{u}_h, \nabla_h \mathbf{v}_h)_{0, \mathcal{T}_h} - \langle \boldsymbol{\kappa} \nabla_h \mathbf{u}_h, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h} \\ & - \epsilon \langle \boldsymbol{\kappa} \nabla_h \mathbf{v}_h, \llbracket \mathbf{u}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h} + \langle \tau \llbracket \mathbf{u}_h \rrbracket, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h}, \end{aligned} \quad (12)$$

124 where $\epsilon \in \{0, \pm 1\}$. The second, third and fourth terms on the right-hand side of (12) are called the consistency,
 125 symmetry and penalty terms, respectively. The discrete bilinear operator $\mathcal{B}_h^{(\epsilon)}$ is symmetric iff $\epsilon = 1$ and is nonsym-
 126 metric otherwise. We obtain the symmetric scheme (H-SIP) if $\epsilon = 1$, the incomplete scheme (H-IIP) if $\epsilon = 0$ and the
 127 nonsymmetric scheme (H-NIP) if $\epsilon = -1$. For all $E \in \mathcal{T}_h$ and $F \in \mathcal{F}_E$, the penalty term is chosen as follows:

$$\tau_F := \frac{\gamma_0 C_{\text{tr}}^2 \kappa_F}{h_F^{1+\delta}} \quad \text{with} \quad \delta \in \mathbb{R}, \quad (13)$$

128 where γ_0 is a user-dependent parameter, C_{tr} is given by (8) and results from the discrete trace inequality (7a), h_F
 129 is a local length scale associated with the interface F , and $\kappa_F := \mathbf{n}_F \boldsymbol{\kappa}_E \mathbf{n}_F$ denotes the normal diffusivity. We then
 130 assume that the quantity h_F satisfies the following *equivalence condition*, where for all $E \in \mathcal{T}_h$ and $F \in \mathcal{F}_E$, there
 131 exist positive constants ρ_1 and ρ_2 independent of h_E such that

$$\rho_1 h_E \leq h_F \leq \rho_2 h_E. \quad (14)$$

132 **Remark 3.1.** Different choices of the local length scale h_F have been suggested in the literature, i.e., $h_F := \text{diam}(F)$,
 133 $h_F := h_E$ (the diameter of E), $h_F := |F|$ (the Lebesgue measure of F) and $h_F := |E|/|F|$ (the Hausdorff measure of
 134 F) (see, e.g., [16]). For simplicity, we assume that $\boldsymbol{\kappa}$ is approximated by piecewise constants on the mesh element \mathcal{T}_h ;
 135 i.e., $\boldsymbol{\kappa}|_E \in \mathbb{R}^{d \times d}$ for all $E \in \mathcal{T}_h$.

136 **Lemma 3.1** (Consistency). *Let $\mathbf{u} = (u, \hat{u}) \in \mathbf{V}$, where $u \in H^s(\Omega)$ is a solution of the weak problem (2) with $s > 3/2$.*
 137 *Then, the following holds:*

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{u}, \mathbf{v}_h) = l(\mathbf{v}_h), \quad \forall \mathbf{v}_h \in \mathbf{V}_h. \quad (15)$$

138 *Proof.* The regularity of the weak solution implies that the quantities u and $\boldsymbol{\kappa} \nabla_h u \cdot \mathbf{n}$ are single-valued fields on the
 139 mesh skeleton; i.e., $\llbracket \mathbf{u} \rrbracket = 0$ for all $E \in \mathcal{T}_h$ and $F \in \mathcal{F}_E$, and $\llbracket \boldsymbol{\kappa} \nabla_h u \rrbracket = 0$ for all $F \in \mathcal{F}_h^i$, where $\llbracket \cdot \rrbracket$ denotes the
 140 standard jump operator as used in the DG method [8]. After integrating by parts, the bilinear form $\mathcal{B}_h^{(\epsilon)}$ yields

$$\begin{aligned} \mathcal{B}_h^{(\epsilon)}(\mathbf{u}, \mathbf{v}_h) &= (\boldsymbol{\kappa} \nabla_h u, \nabla_h v_h)_{0, \mathcal{T}_h} - \langle \boldsymbol{\kappa} \nabla_h u, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h}, \\ &= \sum_{E \in \mathcal{T}_h} (\nabla_h \cdot (-\boldsymbol{\kappa} \nabla_h u), v_h)_{0, E} + \sum_{F \in \mathcal{F}_h^i} \underbrace{\langle \llbracket \boldsymbol{\kappa} \nabla_h u \rrbracket, \hat{v}_h \rangle}_{=0} = \sum_{E \in \mathcal{T}_h} (f, v_h)_{0, E} \quad \forall \mathbf{v}_h \in \mathbf{V}_h, \end{aligned}$$

141 since \hat{v}_h vanishes on the boundary skeleton \mathcal{F}_h^b . This completes the proof. \square

142 A straightforward consequence of the consistency property is the Galerkin orthogonality.

143 **Proposition 3.1** (Galerkin orthogonality). *Let $\mathbf{u} = (u, \hat{u}) \in \mathbf{V}$, where $u \in H^s(\Omega)$ a solution of the weak problem (2)*
 144 *with $s > 3/2$. We denote by $\mathbf{u}_h \in \mathbf{V}_h$ the approximate solution of the discrete problem (11). Then,*

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{u} - \mathbf{u}_h, \mathbf{v}_h) = 0 \quad \forall \mathbf{v}_h \in \mathbf{V}_h. \quad (16)$$

145 *Proof.* Subtracting (15) and (11) yields the assertion. \square

146 3.1. Coercivity and well-posedness

147 The next step is to prove the key property, i.e., the discrete coercivity of the bilinear form $\mathcal{B}_h^{(\epsilon)}$, to ensure the well-
 148 posedness of the discrete problem (11). To this end, we first need to establish an upper bound of the consistency term
 149 using the jump seminorm $|\cdot|_\tau$.

150 **Lemma 3.2** (Bound on consistency term). *Let $(\mathbf{w}_h, \mathbf{v}_h) \in \mathbf{V}_h \times \mathbf{V}_h$; then, there exists a constant $C_\delta > 0$ such that*

$$\left| \langle \boldsymbol{\kappa} \nabla_h \mathbf{w}_h, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h} \right| \leq C_\delta^{1/2} \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}_h\|_{0, \mathcal{T}_h} |\mathbf{v}_h|_\tau, \quad (17)$$

151 where $C_\delta := C_0 h^\delta$ and $C_0 := C \eta_0 \gamma_0^{-1}$ is a positive constant independent of h .

152 *Proof.* The decomposition of the consistency term yields

$$\langle \boldsymbol{\kappa} \nabla_h \mathbf{w}_h, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h} = \sum_{E \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_E} \langle \boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}_h, \boldsymbol{\kappa}^{1/2} \llbracket \mathbf{v}_h \rrbracket \rangle_{0, F}. \quad (18)$$

153 For any $F \in \mathcal{F}_h$, successively applying the Cauchy–Schwarz and the discrete trace inequalities (7a), using the defini-
 154 tion (13) and the equivalence condition (14), we infer that

$$\begin{aligned} \left| \langle \boldsymbol{\kappa} \nabla_h \mathbf{w}_h, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, F} \right| &\leq \left[\frac{h_F^{1+\delta}}{\gamma_0 C_{\text{tr}}^2} \right]^{1/2} (\|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}_h\|_{0, F}) (\tau_F^{1/2} \|\llbracket \mathbf{v}_h \rrbracket\|_{0, F}), \\ &\leq \left[\frac{C h_E^\delta}{\gamma_0} \right]^{1/2} \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}_h\|_{0, E} |\mathbf{v}_h|_{\tau, F}. \end{aligned}$$

155 By summing over all interfaces $F \in \mathcal{F}_E$ and then over all mesh elements $E \in \mathcal{T}_h$ and by using the quasi-uniformity
 156 property of the mesh \mathcal{T}_h , we obtain the assertion

$$\left| \langle \boldsymbol{\kappa} \nabla_h \mathbf{w}_h, \llbracket \mathbf{v}_h \rrbracket \rangle_{0, \partial \mathcal{T}_h} \right| \leq \left[\frac{C \eta_0}{\gamma_0} h^\delta \right]^{1/2} \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}_h\|_{0, \mathcal{T}_h} |\mathbf{v}_h|_\tau, \quad (19)$$

157 which concludes the proof. \square

158 **Lemma 3.3** (Coercivity). *For a penalty parameter γ_0 that is large enough—i.e., $\gamma_0 > 4C\eta_0h^\delta$ —the discrete bilinear*
 159 *form $\mathcal{B}_h^{(\epsilon)}$ is V_h -coercive with respect to the energy-norm $\|\cdot\|_*$; i.e.,*

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{v}_h, \mathbf{v}_h) \geq \frac{1}{2}\|\mathbf{v}_h\|_*^2, \quad (20)$$

160 *for all $\mathbf{v}_h \in V_h$ and for any value of the parameter ϵ .*

161 *Proof.* Setting $\mathbf{u}_h = \mathbf{v}_h$ in the definition of the bilinear form (12), we obtain

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{v}_h, \mathbf{v}_h) = \|\kappa^{1/2}\nabla_h \mathbf{v}_h\|_{0,\mathcal{T}_h}^2 - (1+\epsilon)\langle \kappa\nabla_h \mathbf{v}_h, \llbracket \llbracket \mathbf{v}_h \rrbracket \rrbracket \rangle_{0,\partial\mathcal{T}_h} + |\mathbf{v}_h|_\tau^2. \quad (21)$$

162 Thus, owing to Lemmata 3.2 and using Young's inequality, for any $\zeta > 0$, there exists a constant $C_\zeta^{(\epsilon)} > 0$ such that

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{v}_h, \mathbf{v}_h) \geq \left[1 - \frac{1+\epsilon}{2} \frac{C_\delta}{\zeta}\right] \|\kappa^{1/2}\nabla_h \mathbf{v}_h\|_{0,\mathcal{T}_h}^2 + \left[1 - \frac{1+\epsilon}{2} \zeta\right] |\mathbf{v}_h|_\tau^2 \geq C_\zeta^{(\epsilon)} \|\mathbf{v}_h\|_*^2,$$

163 where $C_\zeta^{(\epsilon)}$ is given by

$$C_\zeta^{(\epsilon)} := 1 - \frac{1+\epsilon}{2} \max(C_\delta/\zeta, \zeta).$$

164 We now select γ_0 such that $C_\delta < \zeta^2$; i.e., $\gamma_0 > \zeta^{-2}C\eta_0h^\delta$. Setting $\zeta = 1/2$, we easily bound $C_{1/2}^{(\epsilon)} \geq 1/2$ for any value
 165 of the parameter ϵ , thus completing the proof. \square

166 **Remark 3.2.** *Note the h^δ -dependency of the coercivity condition. A straightforward consequence of the consistency*
 167 *and coercivity requirements via the Lax–Milgram Theorem is the well-posedness of the weak problem (11); i.e., the*
 168 *existence and uniqueness of the discrete solution $\mathbf{u}_h \in V_h$ are ensured.*

169 3.2. Boundedness

170 We now assume that the discrete bilinear form $\mathcal{B}_h^{(\epsilon)}$ can be extended to $V(h) \times V(h)$, and we assert the boundedness
 171 of the product space. To this end, we introduce the enriched energy-norm on $V(h)$ denoted by $\|\cdot\|$ (which is also a
 172 natural norm on V_h) to bound the (normal) derivative terms [10]:

$$\|\mathbf{v}\|^2 := \|\mathbf{v}\|_*^2 + \sum_{E \in \mathcal{T}_h} h_E \|\kappa^{1/2}\nabla_h \mathbf{v}\|_{0,\partial E}^2, \quad \forall \mathbf{v} \in V(h). \quad (22)$$

173 **Lemma 3.4** (Equivalency of $\|\cdot\|_*$ - and $\|\cdot\|$ -norms). *For all $\mathbf{v} \in V(h)$, the norms $\|\mathbf{v}\|_*$ and $\|\mathbf{v}\|$ are equivalent; i.e., there*
 174 *exists a constant $\rho > 0$ such that*

$$\rho^{-1}\|\mathbf{v}\| \leq \|\mathbf{v}\|_* \leq \|\mathbf{v}\|, \quad (23)$$

175 where $\rho := (1 + \eta_0 C_{\text{tr}}^2)^{\frac{1}{2}}$ depends only on the element shape.

176 *Proof.* Following the definition (22), we notice that $\|\mathbf{v}\|_* \leq \|\mathbf{v}\|$. We now can easily bound the difference of both norms
 177 by using the discrete trace inequality (7a)

$$\|\mathbf{v}\|^2 - \|\mathbf{v}\|_*^2 \leq \eta_0 C_{\text{tr}}^2 \|\kappa^{1/2}\nabla_h \mathbf{v}\|_{0,\mathcal{T}_h}^2 \leq \eta_0 C_{\text{tr}}^2 \|\mathbf{v}\|_*^2, \quad (24)$$

178 which yields the assertion. \square

179 **Lemma 3.5** (Boundedness with h^δ -dependency). *For all $(\mathbf{w}, \mathbf{v}) \in V(h) \times V(h)$, there exists a constant $C_{\text{bnd}} > 0$ such*
 180 *that*

$$\mathcal{B}_h^{(\epsilon)}(\mathbf{w}, \mathbf{v}) \leq C_{\text{bnd}} \|\mathbf{w}\| \cdot \|\mathbf{v}\|, \quad (25)$$

181 where $C_{\text{bnd}} := 2 + C_1 h^\delta$ and $C_1 := (\gamma_0 C_{\text{tr}}^2)^{-1}$ is a positive constant independent of h .

182 *Proof.* Following the definition of the bilinear form (12), we deduce that

$$\begin{aligned} |\mathcal{B}_h^{(\epsilon)}(\mathbf{w}, \mathbf{v})| &\leq |(\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}, \boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v})_{0, \mathcal{T}_h}| + \langle \boldsymbol{\tau}^{1/2} \llbracket \mathbf{w} \rrbracket, \boldsymbol{\tau}^{1/2} \llbracket \mathbf{v} \rrbracket \rangle_{0, \partial \mathcal{T}_h} + \\ &\quad \left| \langle \boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}, \boldsymbol{\kappa}^{1/2} \llbracket \mathbf{v} \rrbracket \rangle_{0, \partial \mathcal{T}_h} \right| + |\epsilon| \left| \langle \boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v}, \boldsymbol{\kappa}^{1/2} \llbracket \mathbf{w} \rrbracket \rangle_{0, \partial \mathcal{T}_h} \right| \\ &\leq |\mathcal{T}_1 + \mathcal{T}_2| + |\mathcal{T}_3| + |\epsilon| |\mathcal{T}_4|. \end{aligned}$$

183 Applying the Cauchy–Schwarz inequality, the first two terms can be bounded as follows:

$$|\mathcal{T}_1 + \mathcal{T}_2| \leq [\|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}\|_{0, \mathcal{T}_h}^2 + |\mathbf{w}|_{\boldsymbol{\tau}}^2]^{1/2} [\|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v}\|_{0, \mathcal{T}_h}^2 + |\mathbf{v}|_{\boldsymbol{\tau}}^2]^{1/2} = \|\mathbf{w}\|_* \|\mathbf{v}\|_*. \quad (26)$$

184 Proceeding as in the proof of Lemmata 3.2, the third and fourth terms can also be bounded as follows:

$$|\mathcal{T}_3| \leq \left[C_1 h^\delta \sum_{E \in \mathcal{T}_h} h_E \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}\|_{0, \partial E}^2 \right]^{1/2} \|\mathbf{v}\|_*, \quad (27a)$$

$$|\mathcal{T}_4| \leq \left[C_1 h^\delta \sum_{E \in \mathcal{T}_h} h_E \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v}\|_{0, \partial E}^2 \right]^{1/2} \|\mathbf{w}\|_*, \quad (27b)$$

185 where $C_1 := (\gamma_0 C_{\text{tr}}^2)^{-1}$. By combining these estimates via the Cauchy–Schwarz inequality, we obtain

$$\begin{aligned} |\mathcal{B}_h^{(\epsilon)}(\mathbf{w}, \mathbf{v})| &\leq \left[(1 + |\epsilon|) \|\mathbf{w}\|_*^2 + C_1 h^\delta \sum_{E \in \mathcal{T}_h} h_E \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{w}\|_{0, \partial E}^2 \right]^{1/2} \left[2 \|\mathbf{v}\|_*^2 + |\epsilon| C_1 h^\delta \sum_{E \in \mathcal{T}_h} h_E \|\boldsymbol{\kappa}^{1/2} \nabla_h \mathbf{v}\|_{0, \partial E}^2 \right]^{1/2} \\ &\leq \max(2, C_1 h^\delta) \|\mathbf{w}\|_* \|\mathbf{v}\|_*, \end{aligned}$$

186 which yields the assertion. \square

187 **Remark 3.3.** Let us emphasize that $C_{\text{bnd}} \leq Ch^r$, where $r = \min(0, \delta)$ and $C := 2 \max(2, C_1)$ and is a positive constant
188 independent of h .

189 4. A priori error analysis

190 We now derive *a priori* error estimates in both the discrete $\|\cdot\|_*$ - and $\|\cdot\|_{0, \mathcal{T}_h}$ -norms to show the accuracy of the
191 H-IP method. To this end, we first recall some definitions such as the continuous interpolant and derive standard
192 interpolation estimates that will be used extensively in the rest of the document (for more details, we refer the reader
193 to [16, 19]). Let us introduce π_h^i and π_h^b , the standard L^2 -orthogonal projectors on the discrete approximation spaces
194 V_h and \hat{V}_h , respectively. Then, if $\phi \in H^s(\Omega)$ with $s \geq 2$, the standard interpolation estimate is written as

$$|\phi - \pi_h^i \phi|_{q, \mathcal{T}_h} \leq Ch^{\mu-q} |\phi|_{\mu, \mathcal{T}_h}, \quad \forall q \in \{0, \dots, s-1\}, \quad (28a)$$

$$\left[\sum_{E \in \mathcal{T}_h} h_E^\alpha \|\nabla_h(\phi - \pi_h^i \phi)\|_{0, \partial E}^2 \right]^{1/2} \leq Ch^{\mu+\frac{\alpha-3}{2}} |\phi|_{\mu, \mathcal{T}_h}, \quad (28b)$$

195 where $\mu := \min(k+1, s)$ and k denote the polynomial degrees of approximation spaces V_h and \hat{V}_h , respectively.

196 **Lemma 4.1** (Optimal error estimates). Let $\mathbf{u} := (u, \hat{u}) \in H^s(\mathcal{T}_h) \times L^2(\mathcal{F}_h)$, where u is the weak solution of (2) and
197 $s > 3/2$. We denote by $\pi_h \mathbf{u} := (\pi_h^i u, \pi_h^b \hat{u})$ the continuous interpolant of the composite variable \mathbf{u} , which is contained
198 in V_h ; i.e., $\pi_h \mathbf{u} \in V_h$. Then,

$$\|\mathbf{u} - \pi_h \mathbf{u}\|_* \leq \|\mathbf{u} - \pi_h \mathbf{u}\| \leq C_\kappa h^{\mu_0} |u|_{\mu, \mathcal{T}_h}, \quad (29)$$

199 where $\mu_0 := \min(k, s-1)$ and $C_\kappa := C \|\boldsymbol{\kappa}^{1/2}\|_{\infty, \Omega}$.

200 *Proof.* Successively using the definition of the $\|\cdot\|$ -norm (22), the Cauchy–Schwarz inequality, and the interpolation
 201 estimates (28) yields

$$\begin{aligned} \|\mathbf{u} - \pi_h \mathbf{u}\|_*^2 &\stackrel{(23)}{\leq} \|\mathbf{u} - \pi_h \mathbf{u}\|^2 \stackrel{(22)}{=} \|\mathbf{u} - \pi_h \mathbf{u}\|_*^2 + \sum_{E \in \mathcal{T}_h} h_E \|\boldsymbol{\kappa}^{1/2} \nabla_h (\mathbf{u} - \pi_h^i \mathbf{u})\|_{0,\partial E}^2, \\ &\stackrel{(28)}{\leq} \|\boldsymbol{\kappa}^{1/2}\|_{\infty,\Omega}^2 \sum_{E \in \mathcal{T}_h} (|\mathbf{u} - \pi_h^i \mathbf{u}|_{1,E}^2 + h_E \|\nabla_h (\mathbf{u} - \pi_h^i \mathbf{u})\|_{0,\partial E}^2), \\ &\leq C^2 \|\boldsymbol{\kappa}^{1/2}\|_{\infty,\Omega}^2 h^{2\mu-2} |\mathbf{u}|_{\mu,\mathcal{T}_h}^2, \end{aligned}$$

202 which concludes the proof. \square

203 4.1. Energy-norm error estimates

204 We now derive an error estimation of the discrete composite variable \mathbf{u}_h in the natural $\|\cdot\|_*$ -norm.

205 **Theorem 4.1** ($\|\cdot\|_*$ -norm estimate and optimal convergence rate). *Let $\mathbf{u} := (u, \hat{u}) \in H^s(\Omega) \times L^2(\mathcal{F}_h)$, where u is a
 206 solution of (2) with $s > 3/2$. We denote by $\mathbf{u}_h \in \mathbf{V}_h$ the approximate solution of the discrete problem (11). Then, for
 207 any value of the parameter δ , the following estimate holds:*

$$\|\mathbf{u} - \mathbf{u}_h\|_* \leq \|\mathbf{u} - \mathbf{u}_h\| \leq C_\kappa h^{\mu_0+r} |\mathbf{u}|_{\mu,\mathcal{T}_h}, \quad (30)$$

208 where $\mu_0 := \min(k, s-1)$, $r := \min(0, \delta)$, and $C_\kappa := C \|\boldsymbol{\kappa}^{1/2}\|_{\infty,\Omega}$.

209 *Proof.* We decompose this quantity as $\mathbf{u} - \mathbf{u}_h = \mathbf{u} - \pi_h \mathbf{u} + \pi_h \mathbf{u} - \mathbf{u}_h$. By using the triangle inequality, we easily infer
 210 that

$$\|\mathbf{u} - \mathbf{u}_h\| \leq \|\mathbf{u} - \pi_h \mathbf{u}\| + \|\pi_h \mathbf{u} - \mathbf{u}_h\|. \quad (31)$$

211 Only an upper bound on the last term of (31) remains to be established. Successively using the coercivity, energy-norm
 212 equivalency, Galerkin orthogonality, and boundedness, we deduce that

$$\begin{aligned} \frac{1}{2\rho^2} \|\pi_h \mathbf{u} - \mathbf{u}_h\|^2 &\stackrel{(23)}{\leq} \frac{1}{2} \|\pi_h \mathbf{u} - \mathbf{u}_h\|_*^2 \stackrel{(20)}{\leq} \mathcal{B}_h^{(\epsilon)}(\pi_h \mathbf{u} - \mathbf{u}_h, \pi_h \mathbf{u} - \mathbf{u}_h), \\ &\stackrel{(16)}{\leq} \mathcal{B}_h^{(\epsilon)}(\pi_h \mathbf{u} - \mathbf{u}, \pi_h \mathbf{u} - \mathbf{u}_h) \stackrel{(25)}{\leq} C_{\text{bnd}} \|\mathbf{u} - \pi_h \mathbf{u}\| \|\pi_h \mathbf{u} - \mathbf{u}_h\|, \end{aligned}$$

213 and then we insert $\|\pi_h \mathbf{u} - \mathbf{u}_h\| \leq 2\rho^2 C_{\text{bnd}} \|\mathbf{u} - \pi_h \mathbf{u}\|$ into (31) to obtain

$$\|\mathbf{u} - \mathbf{u}_h\|_* \stackrel{(23)}{\leq} \|\mathbf{u} - \mathbf{u}_h\| \leq (1 + 2\rho^2 C_{\text{bnd}}) \|\mathbf{u} - \pi_h \mathbf{u}\|.$$

214 Proceeding as in Remark 3.3, we can conclude that there exists a positive constant C such that $1 + 2\rho^2 C_{\text{bnd}} \leq Ch^r$,
 215 which yields the assertion. \square

216 **Corollary 4.1** (Estimate for strong-regularity solutions). *Assume that $s \geq k+1$ with $u \in H_0^{k+1}(\mathcal{T}_h)$ and $\delta \in \mathbb{R}$. Then,
 217 we have the estimate*

$$\|\mathbf{u} - \mathbf{u}_h\|_* \leq C_{u,\kappa} h^{k+r}, \quad (32)$$

218 where $r := \min(0, \delta)$, $C_{u,\kappa} := C_\kappa |u|_{k+1,\mathcal{T}_h}$ and $C_\kappa := C \|\boldsymbol{\kappa}^{1/2}\|_{\infty,\Omega}$.

219 *Proof.* (Evident) \square

220 **Remark 4.1.** *Following Di Pietro and Ern, since C in Theorem 4.1 is independent of $\boldsymbol{\kappa}$, the discrete method is said
 221 to be robust with respect to the anisotropy and heterogeneity of the diffusion tensor. The given estimate (32) indicates
 222 that the order of convergence in the $\|\cdot\|_*$ -norm, or equivalently, $\|\cdot\|$ -norm, is linear and δ -dependent, i.e., suboptimal
 223 if $\delta < 0$ and optimal otherwise.*

224 4.2. L^2 -norm error estimate

225 Using a standard Aubin–Nitsche duality argument, we now derive an improved L^2 -error estimate of the H-IP
 226 method in terms of the parameter δ . To this end, we define an auxiliary function ψ as the solution of the adjoint
 227 problem:

$$-\nabla \cdot (\kappa \nabla \psi) = u - u_h \quad \text{in } \Omega, \quad \text{and } \psi = 0 \quad \text{on } \partial\Omega. \quad (33)$$

228 By assuming elliptic regularity, the following estimate holds:

$$\|\psi\|_{2,\Omega} \leq C_\kappa \|u - u_h\|_{0,\Omega}, \quad (34)$$

229 where C_κ depends on the shape regularity (i.e., the convexity) of Ω and the distribution of κ inside it [20]. The weak
 230 adjoint problem is to find $\psi \in H^2(\Omega) \cap H_0^1(\Omega)$ such that

$$(\kappa \nabla_h \psi, \nabla_h v)_{0,\mathcal{T}_h} - \langle \kappa \nabla_h \psi \cdot \mathbf{n}, v \rangle_{0,\partial\mathcal{T}_h} = (u - u_h, v)_{0,\mathcal{T}_h}, \quad \forall v \in H_0^1(\Omega). \quad (35)$$

231 By setting $v = u - u_h$ in (35), we obtain

$$\|u - u_h\|_{0,\mathcal{T}_h}^2 = (\kappa \nabla_h \psi, \nabla_h(u - u_h))_{0,\mathcal{T}_h} - \langle \kappa \nabla_h \psi, (u - u_h)\mathbf{n} \rangle_{0,\partial\mathcal{T}_h}. \quad (36)$$

232 Let us now introduce the composite error variable $\mathbf{e}_h^\mu := \mathbf{u} - \mathbf{u}_h = (e_h^\mu, \hat{e}_h^\mu)$. From the regularity of the variables \hat{u} ,
 233 \hat{u}_h and ψ , we deduce that $\langle \kappa \nabla_h \psi, (\hat{u} - \hat{u}_h)\mathbf{n} \rangle_{0,\partial\mathcal{T}_h} = 0$. By embedding this condition in (36), we obtain an equivalent
 234 reformulation of the weak adjoint problem in terms of the discrete bilinear operator $\mathcal{B}_h^{(\epsilon)}$:

$$\|e_h^\mu\|_{0,\mathcal{T}_h}^2 = (\kappa \nabla \psi, \nabla e_h^\mu)_{0,\mathcal{T}_h} - \langle \kappa \nabla \psi, \llbracket e_h^\mu \rrbracket \rangle_{0,\partial\mathcal{T}_h} = \mathcal{B}_h^{(\epsilon)}(\psi, \mathbf{e}_h^\mu), \quad (37)$$

where $\psi := (\psi, \hat{\psi})$. Following the definition of the bilinear form $\mathcal{B}_h^{(\epsilon)}$ (12) and using the Galerkin orthogonality
 $\mathcal{B}_h^{(\epsilon)}(\mathbf{e}_h^\mu, \pi_h \psi) = 0$, since $\pi_h \psi \in \mathbf{V}_h$ (see Proposition 3.1), we easily infer

$$\mathcal{B}_h^{(\epsilon)}(\psi, \mathbf{e}_h^\mu) = \mathcal{B}_h^{(\epsilon)}(\mathbf{e}_h^\mu, \mathbf{e}_\pi^\psi) - (1 - \epsilon) \langle \kappa \nabla \psi, \llbracket e_h^\mu \rrbracket \rangle_{0,\partial\mathcal{T}_h} := \mathcal{T}_1^{(\epsilon)} - (1 - \epsilon) \mathcal{T}_2, \quad (38)$$

235 where $\mathbf{e}_\pi^\psi := \psi - \pi_h \psi$. We will now determine an upper bound of the quantity $\|e_h^\mu\|_{0,\mathcal{T}_h}^2$. Owing to Lemmas 3.5 and 4.1
 236 and using the regularity assumption $\psi \in H^2(\Omega)$, we can bound the first term \mathcal{T}_1 :

$$|\mathcal{T}_1^{(\epsilon)}| \leq C_{\text{bnd}} \|\mathbf{e}_\pi^\psi\| \|\mathbf{e}_h^\mu\| \leq C_\kappa C_{\text{bnd}} h \|\psi\|_{2,\Omega} \|\mathbf{e}_h^\mu\|. \quad (39)$$

237 Using the trace inequality $\|\nabla_h \psi\|_{0,\partial\mathcal{T}_h} \leq Ch^{-1/2} \|\psi\|_{2,\Omega}$ [19], the second term \mathcal{T}_2 can be bounded as follows:

$$|\mathcal{T}_2| \leq C_\kappa h^{\frac{1+\delta}{2}} \|\nabla_h \psi\|_{0,\partial\mathcal{T}_h} |e_h^\mu|_\tau \leq C_\kappa h^{\frac{\delta}{2}} \|\psi\|_{2,\Omega} \|\mathbf{e}_h^\mu\|. \quad (40)$$

238 Combining (39) and (40), we obtain the estimate

$$\|u - u_h\|_{0,\mathcal{T}_h} \leq C_\kappa (C_{\text{bnd}} h + (1 - \epsilon) h^{\frac{\delta}{2}}) \|\mathbf{e}_h^\mu\|, \quad (41)$$

239 and we can assert the theorem below.

240 **Theorem 4.2** (L^2 -norm estimate). *Let $\mathbf{u} := (u, \hat{u}) \in H^s(\Omega) \times L^2(\mathcal{F}_h)$, where u is a solution of (2) with $s > 3/2$. We
 241 denote by $\mathbf{u}_h \in \mathbf{V}_h$ the approximate solution of the discrete problem (11). Then, for any value of the parameters δ and
 242 $\epsilon \in \{0, \pm 1\}$, the following estimate holds for the H-IP method:*

$$\|u - u_h\|_{0,\mathcal{T}_h} \leq C_\kappa h^{\mu_0 + s_\delta^{(\epsilon)}} |u|_{s,\mathcal{T}_h}, \quad (42)$$

243 where $C_\kappa := C \|\kappa^{1/2}\|_{\infty,\Omega}$, $\mu_0 := \min(k, s - 1)$, and the parameter $s_\delta^{(\epsilon)}$ is only dependent on ϵ and δ and is given by

$$s_\delta^{(\epsilon)} := \begin{cases} \min(1, 1 + 2\delta) & \text{if } \epsilon = 1, \\ \min(1, \delta/2) & \text{if } \epsilon \neq 1 \text{ and } \delta \geq 0, \\ \min(1 + 2\delta, 3\delta/2) & \text{if } \epsilon \neq 1 \text{ and } \delta < 0. \end{cases} \quad (43)$$

244 *Proof.* The estimate (42) using (43) follows after some algebraic manipulations from the previous equation (41), the
 245 definition of C_{bnd} given in Lemma 3.5 and the optimal error estimate given in Lemma 4.1. \square

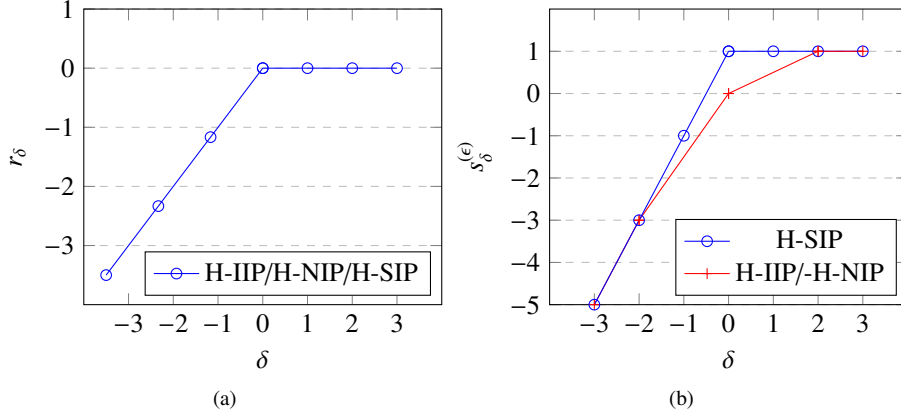


Figure 1: Representation of the quantities r_δ and $s_\delta^{(\epsilon)}$ vs. δ given in Theorems 4.1 and 4.2, respectively.

246 **Remark 4.2.** The authors are certain that the estimates given in Theorems 4.1 and 4.2 have already been established
 247 in the literature, but we have not been able to find them.

248 5. Numerical experiments

249 In the previous sections, we built families of hybridizable interior penalty methods based on an adaptive definition
 250 of the penalty parameter that depends on several coefficients. This section highlights the benefit these methods provide
 251 in the approximation of diffusion problems with anisotropic and/or discontinuous coefficients and in the validation of a
 252 priori error estimates. In the rest of the document, we assume that the local length scale h_F in (13) is chosen to be equal
 253 to the diameter of the associated element, i.e., $h_F := h_E$, for all $E \in \mathcal{T}_h$ and for all $F \in \mathcal{F}_E$. All numerical experiments
 254 are performed using the high-performance finite element library NGSOLVE [21]. Then, the physical domain is taken
 255 to be a unit square—i.e., $\Omega := [0, 1]^2 \subset \mathbb{R}^2$ —and the right-hand-side f is chosen such that the given exact solution
 256 u respecting the homogeneous boundary conditions is verified. We use a sequence of subdivisions \mathcal{T}_h , where regular
 257 triangles or squares form each partition (see, e.g., Figure 2). Standard h - and k -refinement strategies are used
 258 to compute the numerical errors and estimated convergence rates (ECRs). To pursue our quantitative analysis, we first
 259 measure the impact of the parameter δ on the a posteriori error estimates. Second, we point out the crucial role of the
 260 factor κ_n arising in (12) for the robustness of the H-IP methods when the medium becomes highly anisotropic and/or
 261 discontinuous. Finally, we complete our experiments by pointing out some unexpected benefits of the value of γ_0 for
 262 the ECRs of the H-SIP scheme.

263 5.1. Test A: Influence of the parameter δ

264 We consider the following test case, which was previously proposed in Fabien *et al.* [11]: the diffusion tensor is
 265 homogeneous and isotropic— $\kappa = \mathbf{I}_2$ (identity matrix)—and the exact smooth solution is given by $u(x, y) = xy(1 -$
 266 $x)(1 - y) \exp(-x^2 - y^2)$. Then, for all $E \in \mathcal{T}_h$ and for all $F \in \mathcal{F}_E$, we assume that the penalty parameter has the
 267 following simplified form:

$$\tau_F := \frac{\tau_0}{h_E^{1+\delta}}, \quad (44)$$

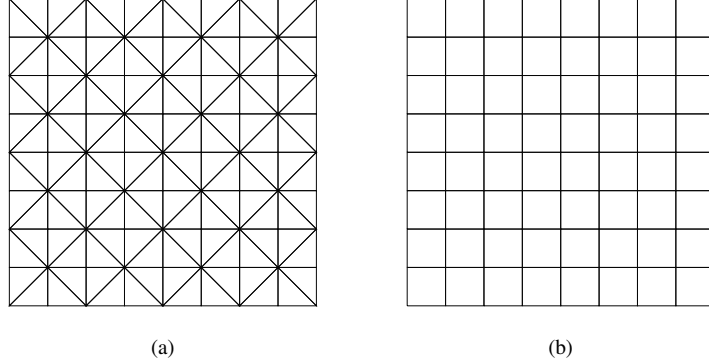


Figure 2: Uniform triangular (a) and square (b) meshes with $h = 1/8$, respectively.

where $\tau_0 > 0$ is a positive constant chosen to be large enough in accordance with Lemma 3.3. The objective here is to measure the impact of the parameter δ on the ECRs in both the L^2 and energy-norms. A history of convergence is shown in Figures 3 ($\|\cdot\|$ -norm) and 4 ($\|\cdot\|_{0,\mathcal{T}_h}$ -norm) for uniform triangular meshes and for polynomial degrees $k \in \{1, \dots, 3\}$, and Table 1 summarizes our numerical observations.

Degree	Norm	H-IIP			H-NIP			H-SIP	
		$\delta = -1$	$\delta = 0$	$\delta \geq 2$	$\delta = -1$	$\delta = 0$	$\delta \geq 2$	$\delta = -1$	$\delta = 0$
$k = 1$	$\ \cdot\ _{0,\mathcal{T}_h}$	1.0	2.0	–	1.0	2.0	–	1.0	2.0
	$\ \cdot\ $	1.0	1.0	–	1.0	1.0	–	1.0	1.0
odd k	$\ \cdot\ _{0,\mathcal{T}_h}$	k	$k+1$	–	$k+1$	–	–	$k+1$	–
	$\ \cdot\ $	$k-1$	k	–	k	–	–	k	–
even k	$\ \cdot\ _{0,\mathcal{T}_h}$	$k-1$	k	$k+1$	k	k	$k+1$	$k+1$	–
	$\ \cdot\ $	$k-1$	k	k	k	k	k	k	–

Table 1: Test A: a summary of the ECRs in the L^2 - and energy-norm of H-IP methods in terms of the parameter δ and the polynomial parity k .

As expected, these observations are in agreement with theoretical estimates and underline that the stabilization parameter δ influences the convergence rate. In particular, we recover some well-known estimates if $\delta = 0$. First, we notice that the convergence of the H-IP method in the energy-norm is linearly δ -dependent if $\delta \leq 0$ and optimal if $\delta \geq 0$, which is in accordance with Lemma 4.1 (see Figure 3). A brief analysis of the convergence in the L^2 -norm indicates that both the H-IIP and H-NIP schemes behave differently from the H-SIP scheme. Nonsymmetric variants are strongly influenced by the polynomial parity of k and by the penalty parameter δ . We observe that the convergence rate increases linearly and optimally if $\delta \geq 0$ for odd k and $\delta \geq 2$ for even k . In this last case, let us point out that the optimal convergence is nearly reached once $\delta \geq 1$. As expected, the symmetric scheme converges optimally when $\delta \geq 0$. These results agree with the theoretical results established in Theorem 4.2.

5.2. Test B: Influence of the parameter κ_F

In the second experiment, we analyze the behavior of the discretization method in the context of genuine anisotropic and heterogeneous properties. Then, the unit square Ω is split into four subdomains $\Omega_1 = [0, 1/2]^2$, $\Omega_2 = [1/2, 1] \times [0, 1/2]$, $\Omega_3 = [1/2, 1]^2$ and $\Omega_4 = [0, 1/2] \times [1/2, 1]$, such that $\Omega := \cup_{i=1}^4 \Omega_i$. The exact solution on the whole domain Ω is given by $u(x, y) = \sin(\pi x) \sin(\pi y)$, and the diffusivity tensor takes different values in each subregion:

$$\kappa = \begin{bmatrix} 1 & 0 \\ 0 & \lambda \end{bmatrix} \quad \text{for } (x, y) \in \Omega_1, \Omega_3, \quad \text{and} \quad \kappa = \begin{bmatrix} 1/\lambda & 0 \\ 0 & 1 \end{bmatrix} \quad \text{for } (x, y) \in \Omega_2, \Omega_4, \quad (45)$$

where the parameter $\lambda > 0$ simultaneously controls both the anisotropy and the medium heterogeneity. Here, we focus on the influence of the parameter κ_F on the robustness of the discretization method in the context of highly anisotropic and heterogeneous coefficients, and we choose $\lambda = 10^{-3}$. In this context, the anisotropy and heterogeneity

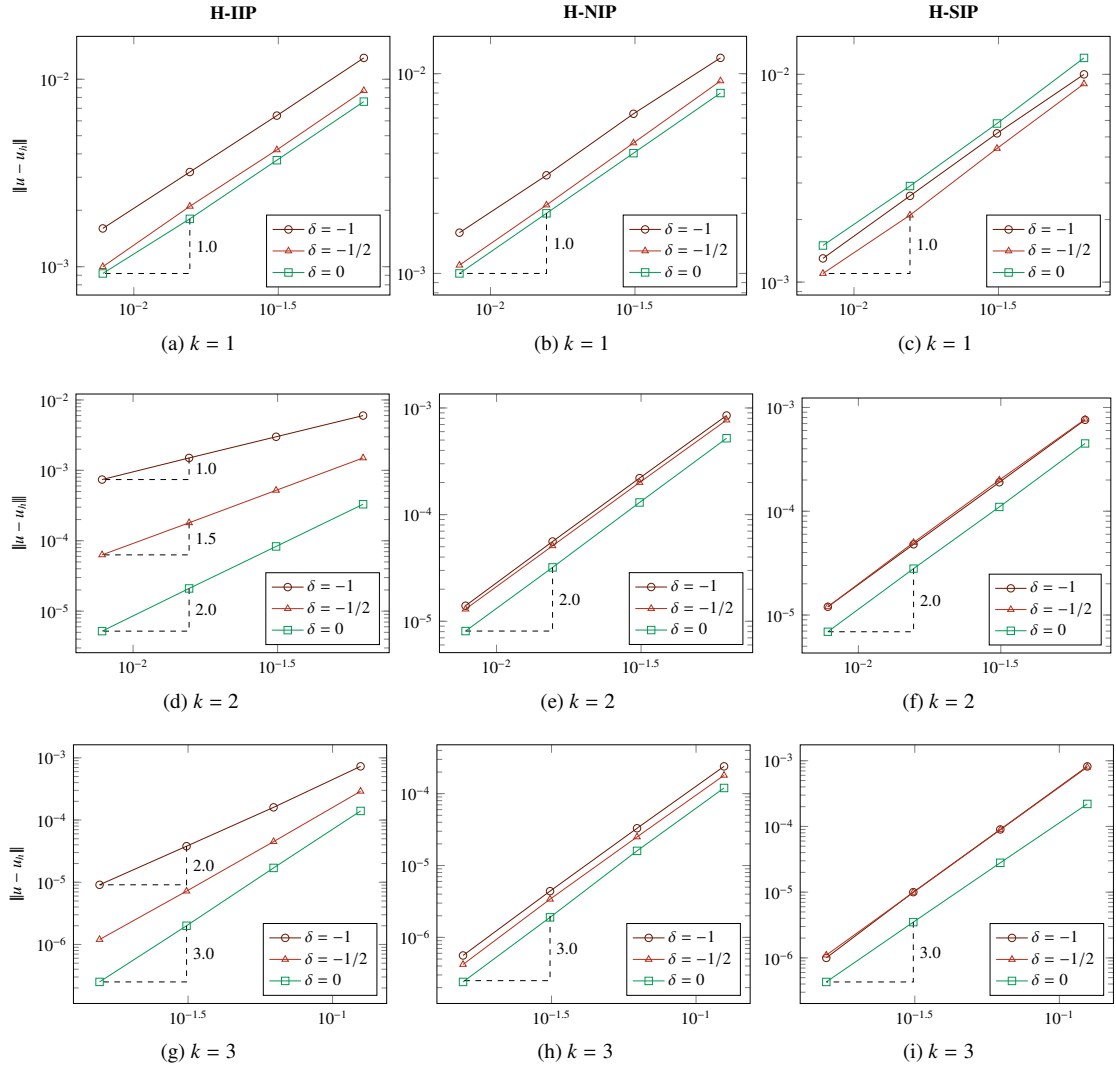


Figure 3: Test A: history of convergence of the H-IP methods with $-1 \leq \delta \leq 0$: $\|u - u_h\|$ vs. h for the three H-IP variants and various polynomial degrees ($1 \leq k \leq 3$) on uniform triangular meshes.

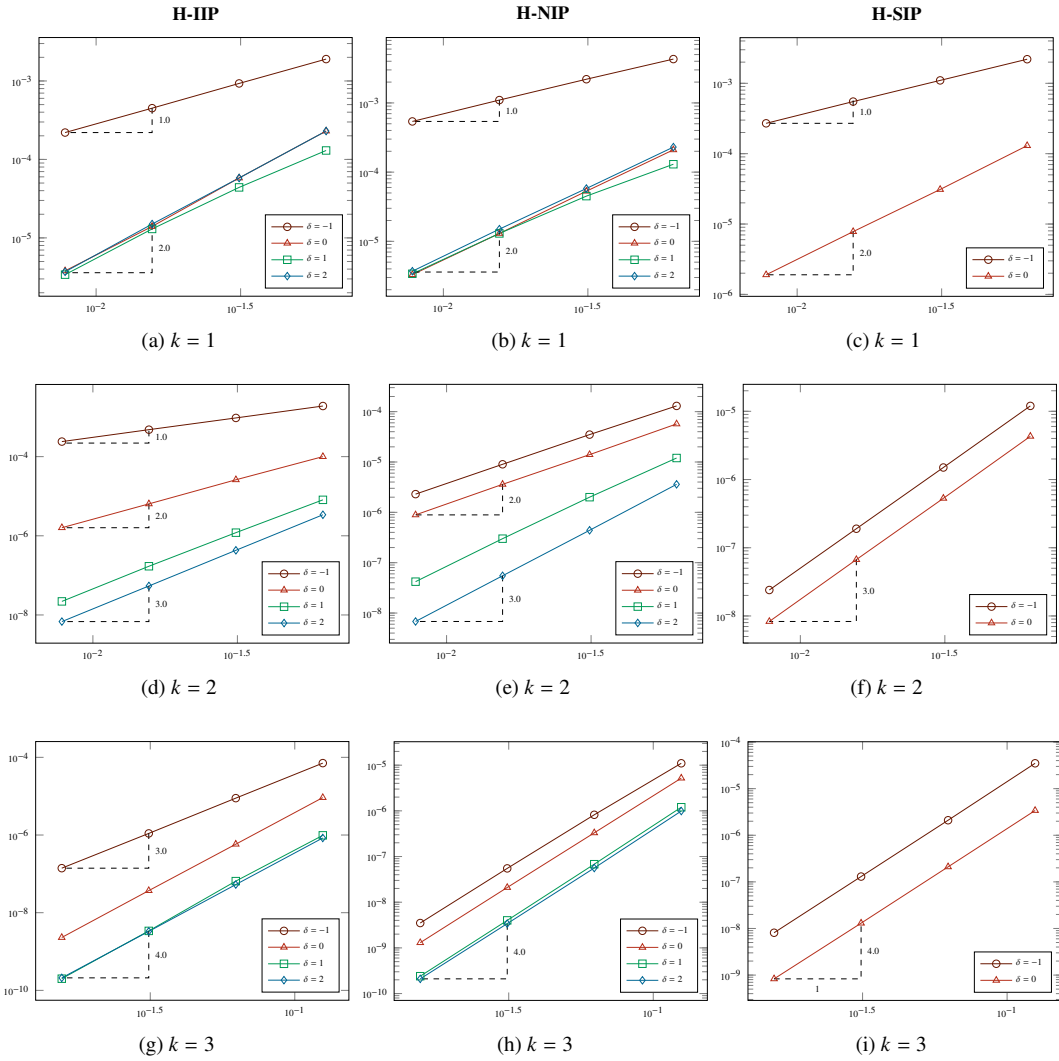


Figure 4: Test A: history of convergence of the H-IP methods with $-1 \leq \delta \leq 2$: $\|u - u_h\|_{0, \mathcal{T}_h}$ vs. h for the three H-IP variants and various polynomial degrees ($1 \leq k \leq 3$) on uniform triangular meshes.

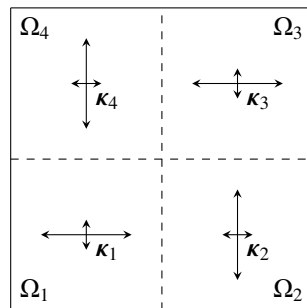


Figure 5: Description of test case B with genuine anisotropic and heterogeneous properties.

289 ratios are approximately 10^3 and 10^6 , respectively. For the simulations, we consider a conforming triangular mesh
 290 ($h = 1/32$) respecting the discontinuities of κ , we use piecewise linear approximations of the discrete variable u_h , and
 291 we set $\delta = 0$ in the definition of the penalty parameter (13). Here, the comparisons are only graphical (Figure 6).
 292 We depict the discrete solutions u_h obtained successively using $\kappa_F := 1$ (Case 1) and $\kappa_F := \mathbf{n}_F \kappa_E \mathbf{n}_F$ (Case 2) for
 293 all variations of $\epsilon \in \{0, \pm 1\}$. In the first situation (Figures 6-a, 6-b and 6-c), the discrete solutions exhibit spurious
 294 oscillations and erratic behavior, thus violating the discrete maximum principle (see, e.g., Table 2). This can be easily
 295 explained by observing that the first formulation does not distinguish between the principal directions of the diffusivity
 296 tensor. Consequently, a misestimated penalty is applied in directions of low or high diffusivity. In the second situation
 297 (Figures 6-d, 6-e and 6-f), the jumps in diffusivity are better captured at the interfaces of discontinuities, and the
 298 discrete solutions are significantly more robust, i.e., exhibit less erratic behavior.

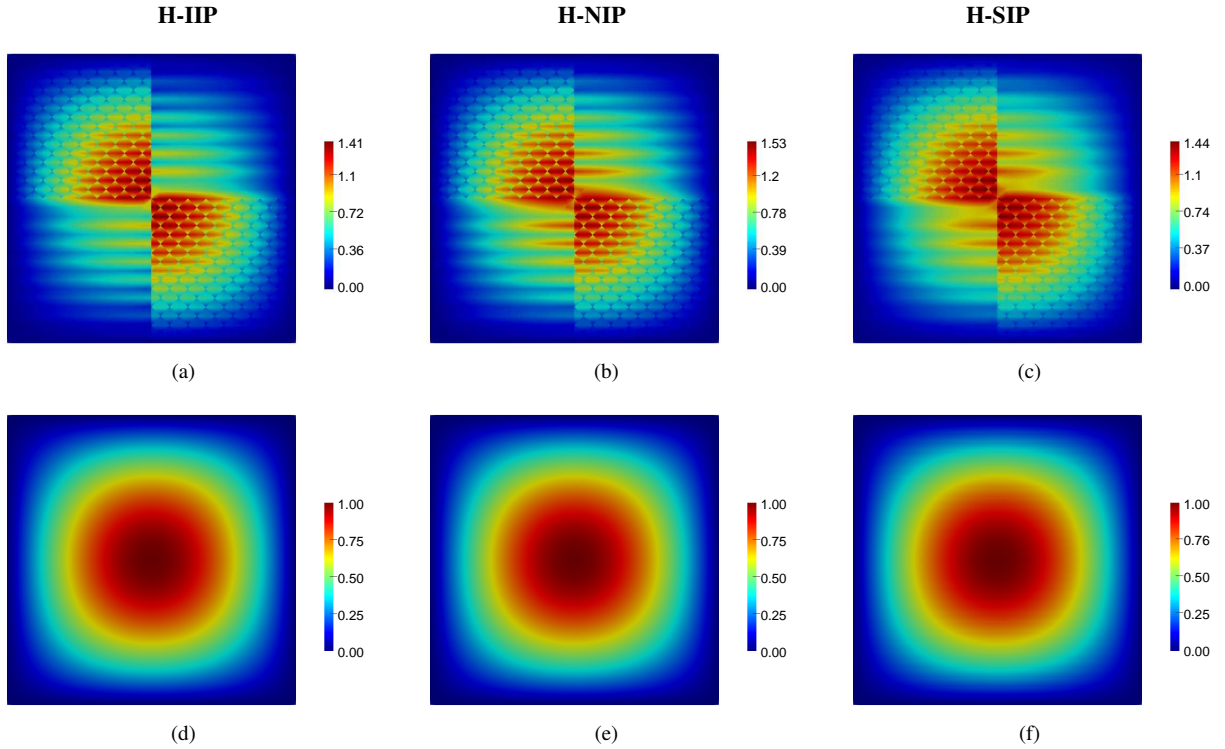


Figure 6: Test B: representation of the discrete solution u_h obtained by the H-IIP, H-NIP and H-SIP schemes, respectively, on the structured triangular mesh ($h = 1/32$). In the top images, the parameter κ_F in (13) is chosen as $\kappa_F := 1$, and in the bottom images, $\kappa_F := \mathbf{n}_F \kappa_E \mathbf{n}_F$.

	H-IIP		H-NIP		H-SIP	
	Case 1	Case 2	Case 1	Case 2	Case 1	Case 2
$\min(u_h)$	$1.54e-03$	$2.14e-03$	$2.79e-03$	$2.09e-03$	$2.68e-03$	$2.12e-03$
$\max(u_h)$	$1.25e+00$	$9.97e-01$	$1.33e+00$	$9.97e-01$	$1.30e+00$	$9.97e-01$
$\ u - u_h\ _{0, \mathcal{T}_h}$	$1.31e-01$	$4.33e-04$	$1.39e-01$	$5.43e-04$	$1.21e-01$	$1.96e-03$

Table 2: Test B: comparison of H-IP methods using a piecewise linear approximation ($u_h \in \mathbb{P}_1(\mathcal{T}_h)$) and two distinct definitions of the coefficient κ_F for highly anisotropic and heterogeneous media ($\lambda = 10^{-3}$). In Case 1, $\kappa_F := 1$, and in Case 2, $\kappa_F := \mathbf{n}_F \kappa_E \mathbf{n}_F$.

299 5.3. Test C: Influence of the parameter γ_0

300 To conclude the sequence of numerical tests, we analyze the influence of the parameter γ_0 on the convergence
 301 of the H-SIP method for κ -orthogonal grids only. For simplicity, we consider the same test case as Test B, (5.2),
 302 and we set two values of the parameter λ : (i) $\lambda = 1$ for a homogeneous and isotropic media and (ii) $\lambda = 0.1$ for a

303 heterogeneous and anisotropic media. We plot the computed L^2 -error of the H-SIP method for a wide range of values
304 of the parameter γ_0 —i.e., $1 \leq \gamma_0 \leq 6$ —using a uniform square mesh ($h = 1/32$). The analysis is done for polynomial
305 degrees $1 \leq k \leq 4$, but the results are presented for $k = 1, 2$ only. Analyzing Figure 7, we observe that there exists an
306 optimal value of the parameter $\gamma_0 := \gamma_{\text{opt}}$ that minimizes the L^2 -error of the scheme. In the context of κ -orthogonal
307 grids, this optimal value ($\gamma_{\text{opt}} = 2$) is insensitive to the mesh form, the mesh size h , the polynomial degree k , and the
308 heterogeneity and/or anisotropy of the media λ . A history of the convergence of the H-SIP method using $\gamma_{\text{opt}} = 2$
309 is then given in Table 7, and we note the surprising superconvergence of u_h ($k + 2$) in the discrete L^2 -norm obtained
310 without any postprocessing. We emphasize that the *superconvergence* property is not achieved for any triangular mesh
311 or any value of the parameter $\epsilon \neq 1$, even using the optimal parameter γ_{opt} in (13).

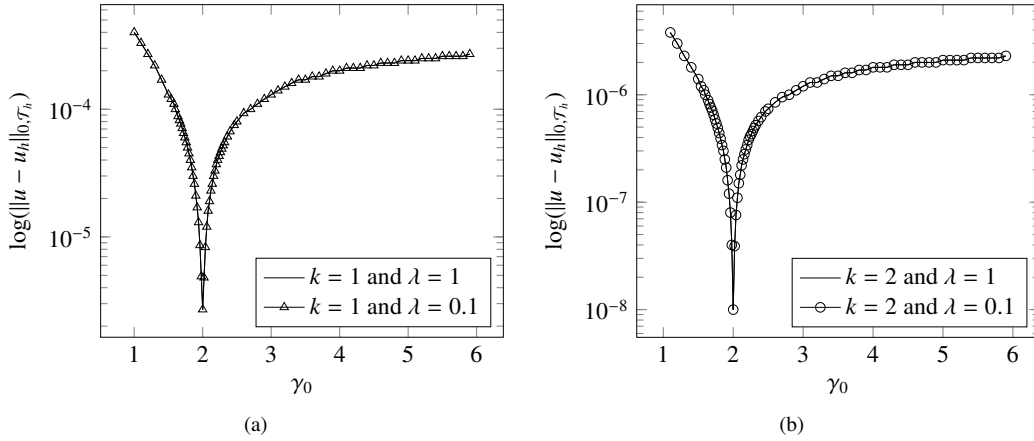


Figure 7: Test C: the L^2 -error of the H-SIP method vs. γ_0 for a uniform square mesh using piecewise linear (a) and quadratic (b) approximations.

	H-SIP ($k = 1$)				H-SIP ($k = 2$)			
	$\lambda = 1$		$\lambda = 0.1$		$\lambda = 1$		$\lambda = 0.1$	
h^{-1}	$\ u - u_h\ _{0, \mathcal{T}_h}$	ECR	$\ u - u_h\ _{0, \mathcal{T}_h}$	ECR	$\ u - u_h\ _{0, \mathcal{T}_h}$	ECR	$\ u - u_h\ _{0, \mathcal{T}_h}$	ECR
8	$1.7e-04$	—	$1.7e-04$	—	$2.6e-06$	—	$2.6e-06$	—
16	$2.1e-05$	3.00	$2.1e-05$	3.00	$1.6e-07$	3.99	$1.6e-07$	3.99
32	$2.7e-06$	3.00	$2.7e-06$	3.00	$1.0e-08$	4.00	$1.0e-08$	4.00
64	$3.4e-07$	3.00	$3.4e-07$	3.00	$6.4e-10$	4.00	$6.4e-10$	4.00

Table 3: Test C: history of the convergence $\|u - u_h\|_{0, \mathcal{T}_h}$ of the H-SIP method using the optimal parameter γ_{opt} on uniform square meshes

312 6. Conclusion

313 We derive improved *a priori* error estimates of families of hybridizable interior penalty discontinuous Galerkin
314 methods using a variable penalty to solve highly anisotropic diffusion problems. The convergence analysis highlights
315 the h^δ -dependency of the coercivity condition and the boundedness requirement that strongly impacts the derived
316 error estimates in terms of both energy- and L^2 -norms. The optimal convergence of the energy-norm is proven for
317 any penalty parameter $\delta \geq 0$ and $\epsilon \in \{0, \pm 1\}$. The situation is somewhat different in L^2 , and distinctive features
318 can be found between the three schemes. Indeed, the symmetric method theoretically converges optimally if $\delta \geq 0$,
319 and non-symmetric variants converge only if $\delta \geq 2$ independently of the polynomial parity. All of these estimates are
320 corroborated by numerical evidence. Notably, the superconvergence of the H-SIP scheme is achieved for κ -orthogonal
321 grids without any postprocessing but only if an appropriate γ_0 is selected.

322 Acknowledgments

323 By convention, the names of the authors are listed in alphabetical order. The corresponding author is grateful to
324 Sander Rhebergen for his invitation to the Department of Applied Mathematics at the University of Waterloo (UW) in
325 April 2019. Thank you very much for introducing me to HDG methods and to the subtle mechanisms of the NGSolve
326 library. Special thanks to UW for the kind hospitality. He also would like to thank Béatrice Rivière at Rice University
327 for her insightful suggestions and remarks concerning *a priori* error estimates. Our fruitful discussions of hybridizable
328 interior penalty methods using superpenalties were the source of inspiration of the present work.

- 329 [1] B. Cockburn, J. Gopalakrishnan, R. Lazarov, Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for
330 second order elliptic problems, *SIAM Journal on Numerical Analysis* 47 (2) (2009) 1319–1365.
- 331 [2] H. Egger, J. Schöberl, A mixed-hybrid-discontinuous Galerkin finite element method for convection-diffusion problems, *IMA J. Numer. Anal*
332 30 (2009) 1–2.
- 333 [3] B. Cockburn, B. Dong, J. Guzmán, M. Restelli, R. Sacco, A hybridizable discontinuous Galerkin method for steady-state convection-
334 diffusion-reaction problems, *SIAM Journal on Scientific Computing* 31 (5) (2009) 3827–3846.
- 335 [4] N. C. Nguyen, J. Peraire, B. Cockburn, An implicit high-order hybridizable discontinuous Galerkin method for linear convection-diffusion
336 equations, *Journal of Computational Physics* 228 (9) (2009) 3232–3254.
- 337 [5] I. Oikawa, HDG METHODS FOR SECOND-ORDER ELLIPTIC PROBLEMS (Numerical Analysis: New Developments for Elucidating
338 Interdisciplinary Problems II), *RIMS Kokyuroku* 2037 (2017) 61–74.
- 339 [6] K. L. Kirk, S. Rhebergen, Analysis of a Pressure-Robust Hybridized Discontinuous Galerkin Method for the Stationary Navier–Stokes
340 Equations, *Journal of Scientific Computing* 81 (2) (2019) 881–897.
- 341 [7] M. S. Fabien, M. Knepley, B. Riviere, A high order hybridizable discontinuous Galerkin method for incom-
342 pressible miscible displacement in heterogeneous media, *Results in Applied Mathematics* (2020) 100089doi:
343 [\\let\@tempa\bibinfo@X@doihttps://doi.org/10.1016/j.rinam.2019.100089](https://doi.org/10.1016/j.rinam.2019.100089).
- 344 [8] D. N. Arnold, F. Brezzi, B. Cockburn, L. D. Marini, Unified analysis of discontinuous Galerkin methods for elliptic problems, *SIAM journal*
345 *on numerical analysis* 39 (5) (2002) 1749–1779.
- 346 [9] R. M. Kirby, S. J. Sherwin, B. Cockburn, To CG or to HDG: a comparative study, *Journal of Scientific Computing* 51 (1) (2012) 183–212.
- 347 [10] C. Lehrenfeld, Hybrid discontinuous Galerkin methods for solving incompressible flow problems, *Rheinisch-Westfälischen Technischen*
348 *Hochschule Aachen* (2010) 111.
- 349 [11] M. S. Fabien, M. G. Knepley, B. M. Riviere, Families of Interior Penalty Hybridizable discontinuous Galerkin methods for second order
350 elliptic problems, *Journal of Numerical Mathematics* (0), doi:[\\let\@tempa\bibinfo@X@doihttps://doi.org/10.1515/jnma-2019-0027](https://doi.org/10.1515/jnma-2019-0027).
- 351 [12] L. Dijoux, V. Fontaine, T. A. Mara, A projective hybridizable discontinuous Galerkin mixed method for
352 second-order diffusion problems, *Applied Mathematical Modelling* 75 (2019) 663–677, ISSN 0307-904X, doi:
353 [\\let\@tempa\bibinfo@X@doihttps://doi.org/10.1016/j.apm.2019.05.054](https://doi.org/10.1016/j.apm.2019.05.054).
- 354 [13] G. N. Wells, Analysis of an interface stabilized finite element method: the advection-diffusion-reaction equation, *SIAM Journal on Numerical*
355 *Analysis* 49 (1) (2011) 87–109.
- 356 [14] D. N. Arnold, An interior penalty finite element method with discontinuous elements, *SIAM journal on numerical analysis* 19 (4) (1982)
357 742–760.
- 358 [15] B. Riviere, *Discontinuous Galerkin methods for solving elliptic and parabolic equations: theory and implementation*, SIAM, 2008.
- 359 [16] D. A. Di Pietro, A. Ern, *Mathematical aspects of discontinuous Galerkin methods*, vol. 69, Springer Science & Business Media, 2011.
- 360 [17] B. Rivière, M. Wheeler, V. Girault, Improved energy estimates for interior penalty, constrained and discontinuous Galerkin methods for elliptic
361 problems, *Computational Geosciences* 3 (3) (1998) 337–360, doi:[\\let\@tempa\bibinfo@X@doihttps://doi.org/10.1023/A:1011591328604](https://doi.org/10.1023/A:1011591328604).
- 362 [18] J. Guzmán, B. Rivière, Sub-optimal Convergence of Non-symmetric Discontinuous Galerkin Methods for Odd Polynomial Approximations,
363 *Journal of Scientific Computing* 40 (1) (2009) 273–280.
- 364 [19] P. G. Ciarlet, Basic error estimates for elliptic problems, *Handbook of Numerical Analysis* 2 (1991) 17–351.
- 365 [20] A. Ern, A. F. Stephansen, P. Zunino, A discontinuous Galerkin method with weighted averages for advection–diffusion equations with locally
366 small and anisotropic diffusivity, *IMA Journal of Numerical Analysis* 29 (2) (2009) 235–256.
- 367 [21] J. Schöberl, C++ 11 implementation of finite elements in NGSolve, Institute for Analysis and Scientific Computing, Vienna University of
368 Technology .