



**HAL**  
open science

# A DPG FRAMEWORK FOR STRONGLY MONOTONE OPERATORS

Pierre Cantin, Norbert Heuer

► **To cite this version:**

Pierre Cantin, Norbert Heuer. A DPG FRAMEWORK FOR STRONGLY MONOTONE OPERATORS. 2018. hal-01690281

**HAL Id: hal-01690281**

**<https://hal.science/hal-01690281>**

Preprint submitted on 22 Jan 2018

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

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

# A DPG FRAMEWORK FOR STRONGLY MONOTONE OPERATORS\*

PIERRE CANTIN<sup>†</sup> AND NORBERT HEUER<sup>‡</sup>

**Abstract.** We present and analyze a hybrid technique to numerically solve strongly monotone nonlinear problems by the discontinuous Petrov–Galerkin method with optimal test functions (DPG). Our strategy is to relax the nonlinear problem to a linear one with additional unknown and to consider the nonlinear relation as a constraint. We propose to use optimal test functions only for the linear problem and to enforce the nonlinear constraint by penalization. In fact, our scheme can be seen as a minimum residual method with nonlinear penalty term. We develop an abstract framework of the relaxed DPG scheme and prove under appropriate assumptions the well-posedness of the continuous formulation and the quasi-optimal convergence of its discretization. As an application we consider an advection-diffusion problem with nonlinear diffusion of strongly monotone type. Some numerical results in the lowest-order setting are presented to illustrate the predicted convergence.

**Key words.** Discontinuous Petrov–Galerkin method, optimal test functions, strongly monotone operator, advection-diffusion, nonlinear penalty.

**AMS subject classifications.** 65N30, 65J15, 65N12, 47H05

**1. Introduction.** In recent years, the discontinuous Petrov–Galerkin method with optimal test functions (“DPG method” in the following) has proved to be an attractive strategy to produce inf-sup stable approximations for a wide class of problems. The basic setting stems from Demkowicz and Gopalakrishnan [14, 13] and has been extended, e.g., to linear elasticity [1, 18], the Stokes and Maxwell equations [28, 7], the Schrödinger equation [15], boundary integral and fractional equations [24, 17]. Another promising application area is singularly perturbed problems [16, 9, 3, 4, 25].

All the cited references, however, deal with linear problems. An extension of the DPG technology to nonlinear problems, on the other hand, is a delicate issue. Principal problem is that the calculation (or approximation) of optimal test functions involves an application of the underlying operator (the DPG method is a minimum residual method). For nonlinear problems this step thus becomes nonlinear, i.e., expensive. One way to circumvent the nonlinearity is, of course, to linearize the underlying problem. This has been the approach in [8, 29]. A different idea is to apply the minimum residual technique in product or “broken” spaces to the nonlinear problem. Bui-Thanh and Ghattas [5] did this by considering the entire nonlinear problem as a constraint, and Carstensen *et al.* [6] developed a representation of the DPG scheme by a nonlinear mixed form and analyzed the case of lowest order approximations. DPG for contact problems has been studied in [21], though in this case the nonlinearity is due to the contact condition which is treated by a variational inequality. We also note that Muga and van der Zee [27] study problems posed in Banach spaces. In those cases the calculation of optimal test functions becomes nonlinear even though the underlying PDE is linear.

In this paper we propose a combined scheme that employs the DPG technique to a *linear relaxation* of the nonlinear problem and where the nonlinearity is added as a constraint. Specifically, we relax the nonlinear problem by introducing an additional variable which then has a nonlinear relation with the original variables. This nonlinear relation is dealt with outside the DPG framework which can therefore develop its full potential, e.g., for singularly perturbed problems. In fact, although here we consider continuous and strongly monotone operators, we claim that our technique is applicable to

---

\***Fundings** :This work has been supported by CONICYT-Chile through FONDECYT projects 3170170 and 1150056.

<sup>†</sup>Facultad de Matemáticas, Pontificia Universidad Católica de Chile, Avenida Vicuña Mackenna 4860, Santiago, Chile. pierre.cantin@mat.uc.cl.

<sup>‡</sup>Facultad de Matemáticas, Pontificia Universidad Católica de Chile, Avenida Vicuña Mackenna 4860, Santiago, Chile. nheuer@mat.uc.cl.

singularly perturbed strongly monotone operators. This, and extensions to more general nonlinear problems, is ongoing research. In contrast, we do not see an obvious extension of [5] or [6] to singularly perturbed nonlinear problems, except for introducing a linear relaxation as we propose.

In the context of the nonlinear DPG scheme from [6] we mentioned their extension to a nonlinear mixed form. In fact, in the linear case it is well known that there is a mixed form of the DPG scheme, and this is precisely the method proposed (for a specific model problem) by Cohen *et al.* [11]. As we will see, our scheme can also be viewed as an extension of this mixed form. To be specific at this point, let us consider a (linear) continuously invertible operator  $A : U \rightarrow V'$  with Banach space  $U$ , Hilbert space  $V$ , and dual  $V'$ . A mixed (or saddle-point) form of  $A\mathbf{u} = F$  is

$$\begin{pmatrix} 0 & A^* \\ A & -R \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} = \begin{pmatrix} 0 \\ F \end{pmatrix} \quad \text{in } U' \times V'$$

with solution  $(\mathbf{u}, \mathbf{v}) = (\mathbf{u}, 0)$ . Here,  $R : V \rightarrow V'$  is the Riesz operator and  $A^* : V \rightarrow U'$  the adjoint of  $A$ . The (practical) DPG method [23] can be seen as a conforming discretization of this saddle point problem. In the nonlinear case, our method is equivalent to replacing the operator  $A$  by a linear relaxation  $B$  and the zero block by a nonlinear operator  $C$  (and, of course, redefining spaces and variables). This yields an operator of the form

$$(1) \quad \begin{pmatrix} C & B^* \\ B & -R \end{pmatrix} : U \times V \rightarrow U' \times V',$$

whose stability relies (among other properties) on the boundedness below of  $B$ . In our case, the nonlinearity is outsourced to  $C$  so that all the (linear) DPG strategies, aiming precisely at the boundedness-below property, can be employed. In our analysis it will be necessary to weight the operator  $R$ , though in specific applications there are precise bounds for this weighting parameter. We stress the fact that, since both  $B$  and  $C$  are acting on the unknowns of interest (represented by  $\mathbf{u}$ ), it is not necessary to consider the variable  $\mathbf{v}$ . Indeed, in the numerical scheme we will be dealing with the Schur complement of  $-R$  only. Standard DPG feature is to use product (broken) spaces  $V$  so that the numerical inversion of a discretization of  $R$  can be done locally and is, thus, cheap.

In the linear context, similar ideas have been used to deal with boundary, transmission and contact conditions outside the DPG framework, [19, 20, 21]. In fact, our abstract framework includes the analysis of linear boundary and transmission problems presented in [19, 20] as special cases. Though, differently from before, we decompose (or extend) nonlinear operators and develop an analysis based on the saddle point structure. Furthermore, we present an analysis that includes the approximation of optimal test functions whereas in [19, 20], these functions were assumed to be known exactly.

In this paper we consider, as a model, an advection-diffusion problem with nonlinear strongly monotone diffusivity,

$$(2) \quad -\nabla \cdot (\boldsymbol{\lambda}(|\nabla u|)\nabla u + \boldsymbol{\beta}u) = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$

with connected Lipschitz domain  $\Omega \subset \mathbb{R}^d$  and  $d \geq 2$ . Of course, there is extensive literature on the numerical analysis of advection-diffusion problems, going back at least to Ciarlet, Schultz and Varga [10] when considering monotone operators. We do not start to discuss all the options as there are too many and since, more importantly, we use this problem only as a model to illustrate our idea and to show its applicability.

Considering the model problem (2), our relaxed linearized problem will be

$$-\nabla \cdot (\boldsymbol{\rho} + \boldsymbol{\beta}u) = f \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega$$

with the nonlinear flux  $\boldsymbol{\rho} = \boldsymbol{\rho}(\nabla u)$  as additional variable. The solution of this problem is not unique so that the associated operator, denoted by  $B$ , has a non-trivial kernel. The missing nonlinear closure relation  $\boldsymbol{\rho} = \boldsymbol{\lambda}(|\nabla u|)\nabla u$  will be represented by the kernel of a nonlinear operator  $C$ .

The rest of this paper is organized as follows. In the next section we present an abstract framework for operators of the form (1). To the best of our knowledge, such kind of operator has not been analyzed before, although the particular case of  $R = 0$  has the typical structure of a mixed method for nonlinear problems. Under appropriate assumptions we prove its invertibility by using the Schur complement (Theorem 2.3 in §2.1). In §2.2 we present our discrete scheme in abstract form and prove its quasi-optimal convergence (Theorem 2.8). The remainder of the paper is devoted to applying the abstract framework to the model problem (2). In Section 3, we precisely define the model problem, state necessary assumptions, and introduce spaces and norms. The introduction of meshes and corresponding product spaces is necessary for the DPG approximation, i.e., to localize the calculation of optimal test functions. In Section 4 we develop a variational formulation of our model problem with resulting operator of the type (1). The relaxed linear part is developed in §4.1. Here, any well-posed variational formulation of the linear problem will do, but for illustration we focus on an ultra-weak variant. This is by no means mandatory and, indeed, different formulations are equivalent, cf. [7]. Though in complicated cases like singular perturbations ultra-weak formulations are easier to analyze (current state of the art) and give the option of direct access to field variables, cf. [16, 25]. The nonlinear closure relation is studied in §4.2 and afterwards, in §4.3, the combined variational formulation is presented and its well-posedness is proved (Theorem 4.7). With all these preparations at hand, the presentation of our *relaxed DPG scheme* for the model problem is brief and a proof of its quasi-optimal convergence (Theorem 5.1) is immediate. This is the contents of Section 5. In Section 6 we present a numerical realization of our relaxed DPG scheme for the model problem and report on results for the cases with and without advection.

To alleviate notation, the expression  $|\cdot|$  is context-dependent and denotes either the Lebesgue measure of a set, the absolute value of a real number or the Euclidean norm of a vector. We use boldface letters for vector and tensor valued quantities. In the calculation of norms via duality, suprema are taken over non-zero elements without further notice.

**2. Abstract nonlinear penalized mixed problem.** In this section we present the abstract framework of our DPG scheme. We first discuss specific continuous formulations, as an operator system similar to a saddle point problem and its Schur complement. In the second part we present two discretizations. The first is a conforming scheme based on the Schur complement and amounts to a DPG method with exactly optimal test functions. The second discretization uses the operator system and amounts to approximating the optimal test functions. Under appropriate assumptions we prove the quasi-optimal convergence of both methods (Theorems 2.5 and 2.8).

**2.1. Continuous setting.** Let  $\mathbf{U}$  and  $\mathbf{V}$  be two real Hilbert spaces with topological duals  $\mathbf{U}'$  and  $\mathbf{V}'$ , respectively. We consider a bounded linear operator  $B : \mathbf{U} \rightarrow \mathbf{V}'$ , an isomorphism  $R : \mathbf{V} \rightarrow \mathbf{V}'$ , and a continuous nonlinear operator  $C : \mathbf{U} \rightarrow \mathbf{U}'$ . (Later,  $C$  will be assumed to be Lipschitz continuous and strongly monotone in a certain sense.) Then we define the block operator  $\mathbf{T}_\kappa : \mathbf{U} \times \mathbf{V} \rightarrow \mathbf{U}' \times \mathbf{V}'$ , with  $\kappa > 0$ , as

$$\mathbf{T}_\kappa = \begin{pmatrix} C & B^* \\ B & -\frac{1}{\kappa}R \end{pmatrix}.$$

Here,  $B^* : \mathbf{V} \rightarrow \mathbf{U}'$  is the adjoint operator of  $B$ , i.e.,  $\langle B\mathbf{u}, \mathbf{v} \rangle = \langle B^*\mathbf{v}, \mathbf{u} \rangle$  for all  $\mathbf{u} \in \mathbf{U}$  and  $\mathbf{v} \in \mathbf{V}$ . To alleviate the notation,  $\langle \cdot, \cdot \rangle$  is used generically to denote the duality between two arbitrary dual spaces. Denoting  $\mathcal{N}(B)$  and  $\mathcal{R}(B)$  the kernel and the range of  $B$ , respectively, we define  $P_B : \mathbf{U} \rightarrow \mathcal{N}(B)$  the

$U$ -orthogonal projector of  $U$  on  $\mathcal{N}(B)$ .

The following theorem provides sufficient conditions on the continuous operators  $B$ ,  $C$  and  $R$  such that the operator  $\mathbf{T}_\kappa$  is bijective, with Lipschitz continuous inverse.

PROPOSITION 2.1. *Additionally to the assumptions made above, assume that*

- (i)  $R^{-1}$  is coercive on  $\mathcal{R}(B)$  with coercivity constant  $c_R^{-1} > 0$ ,
- (ii) there exists  $c_B > 0$  such that  $\|\mathbf{u} - P_B(\mathbf{u})\|^2 \leq c_B \|B\mathbf{u}\|_{\mathbf{V}'}^2$ , for all  $\mathbf{u} \in U$ ,
- (iii) there exist  $c_U, c_V > 0$  such that

$$c_U \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_U^2 \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + c_V \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2$$

for all  $\mathbf{u}_1, \mathbf{u}_2 \in U$ .

Then, for all  $\kappa \geq \kappa_0 := (c_V + c_U c_B) c_R$ ,  $\mathbf{T}_\kappa$  is bijective with Lipschitz continuous inverse. In particular,

$$(3) \quad c_U \|\mathbf{u}_1 - \mathbf{u}_2\|_U^2 \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + \kappa \langle R^{-1} B(\mathbf{u}_1 - \mathbf{u}_2), B(\mathbf{u}_1 - \mathbf{u}_2) \rangle$$

for all  $\mathbf{u}_1, \mathbf{u}_2 \in U$ .

*Proof.* Since  $R$  is bijective we can consider the Schur complement of  $-\frac{1}{\kappa}R$  to obtain

$$(4) \quad \begin{pmatrix} C + \kappa B^* R^{-1} B & 0 \\ 0 & -\kappa^{-1} R \end{pmatrix} = \mathbf{S}_\kappa^\dagger \mathbf{T}_\kappa \mathbf{S}_\kappa^\dagger.$$

Here,  $\mathbf{S}_\kappa^\dagger : U \times V \rightarrow U \times V$  and  $\mathbf{S}_\kappa^\dagger : U' \times V' \rightarrow U' \times V'$  are the respective involutory operators (i.e., equal to their inverses) defined as

$$\mathbf{S}_\kappa^\dagger = \begin{pmatrix} \text{Id}_U & 0 \\ \kappa R^{-1} B & -\text{Id}_V \end{pmatrix}, \quad \text{and} \quad \mathbf{S}_\kappa^\dagger = \begin{pmatrix} \text{Id}_{U'} & \kappa B^* R^{-1} \\ 0 & -\text{Id}_{V'} \end{pmatrix}.$$

It then follows that  $\mathbf{T}_\kappa$  is invertible if and only if the operator  $\mathbf{D}_\kappa = C + \kappa B^* R^{-1} B : U \rightarrow U'$  is invertible. Since  $B^* R^{-1} B$  defines a linear bounded operator we infer that  $\mathbf{D}_\kappa$  is continuous. Let us now prove that  $\mathbf{D}_\kappa$  is strongly monotone, specifically that (3) holds. Let  $\mathbf{u}_1, \mathbf{u}_2 \in U$  and start observing that

$$\|\mathbf{u}_1 - \mathbf{u}_2\|_U^2 = \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_U^2 + \|\mathbf{u}_1 - \mathbf{u}_2 - P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_U^2,$$

owing to the definition of  $P_B$ . By assumption (i) we can bound

$$\|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2 \leq c_R \langle R^{-1} B(\mathbf{u}_1 - \mathbf{u}_2), B(\mathbf{u}_1 - \mathbf{u}_2) \rangle.$$

Combining this bound with assumptions (ii) and (iii), it follows that

$$c_U \|\mathbf{u}_1 - \mathbf{u}_2\|_U^2 \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + (c_V + c_B c_U) c_R \langle R^{-1} B(\mathbf{u}_1 - \mathbf{u}_2), B(\mathbf{u}_1 - \mathbf{u}_2) \rangle.$$

As a result,  $\mathbf{D}_\kappa$  is continuous and strongly monotone for all  $\kappa \geq \kappa_0 := (c_V + c_B c_U) c_R$ , i.e.,  $\mathbf{D}_\kappa$  is invertible for all  $\kappa \geq \kappa_0$  with Lipschitz continuous inverse, see, e.g., [26, Chap. 2]  $\square$

*Remark 2.2.* (i) In the case of a linear operator  $C$ ,  $-\frac{1}{\kappa}R$  plays the role of a regularizing operator needed for the bijectivity of  $\mathbf{T}_\kappa$ . As in the linear case, here the invertibility of  $\mathbf{T}_\kappa$  does not require the surjectivity of  $B$ , but only that the range  $\mathcal{R}(B)$  is closed in  $V'$  (see, e.g., [22, Remark 4.3]).

(ii) In Proposition 2.1, if we assume in addition that  $C$  is monotone, i.e.,  $\langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle \geq 0$  for all  $\mathbf{u}_1, \mathbf{u}_2 \in U$ , then we infer from (3) that  $\mathbf{T}_\kappa$  is bijective with Lipschitz continuous inverse for all  $\kappa \geq 1$ .

Now, given  $F \in \mathbf{V}'$  and  $G \in \mathbf{U}'$ , we consider two variational formulations, of penalized mixed form

$$(5) \quad (\mathbf{u}, \mathbf{v}) \in \mathbf{U} \times \mathbf{V} : \quad \mathbf{T}_\kappa(\mathbf{u}, \mathbf{v}) = (G, F) \quad \text{in } \mathbf{U}' \times \mathbf{V}',$$

and the reduced Schur variant

$$(6) \quad \mathbf{u} \in \mathbf{U} : \quad \langle \mathbf{D}_\kappa \mathbf{u}, \mathbf{w} \rangle = \kappa \langle R^{-1}F, B\mathbf{w} \rangle + \langle G, \mathbf{w} \rangle \quad \forall \mathbf{w} \in \mathbf{U},$$

with

$$\langle \mathbf{D}_\kappa \mathbf{u}, \mathbf{w} \rangle := \kappa \langle R^{-1}B\mathbf{u}, B\mathbf{w} \rangle + \langle C(\mathbf{u}), \mathbf{w} \rangle \quad (\mathbf{u}, \mathbf{w} \in \mathbf{U}).$$

The following result is a consequence of the Schur factorization (4) and Proposition 2.1.

**THEOREM 2.3.** *Let the assumptions of Proposition 2.1 hold true. Then, for all  $\kappa \geq \kappa_0$  with  $\kappa_0$  as before, problems (5) and (6) are uniquely solvable and equivalent. Specifically, if  $(\mathbf{u}, \mathbf{v}) \in \mathbf{U} \times \mathbf{V}$  solves (5), then  $\mathbf{u}$  solves (6). Conversely, if  $\mathbf{u} \in \mathbf{U}$  solves (6), then  $(\mathbf{u}, \kappa R^{-1}(B\mathbf{u} - F)) \in \mathbf{U} \times \mathbf{V}$  solves (5).*

Under additional assumptions on the operators  $B$  and  $C$ , we have a strong characterization of the solution of (6).

**PROPOSITION 2.4.** *Assume that*

- (i)  $\langle C(\mathbf{u}) - G, \mathbf{v} \rangle = 0$  for all  $\mathbf{v} \in \mathcal{N}(B)$  implies  $C(\mathbf{u}) = G$  in  $\mathbf{U}'$ ,
- (ii)  $B : \mathbf{U} \rightarrow \mathbf{V}'$  is surjective.

*Then,  $\mathbf{u} \in \mathbf{U}$  solves (6) if and only if  $C(\mathbf{u}) = G$  in  $\mathbf{U}'$  and  $B\mathbf{u} = F$  in  $\mathbf{V}'$ .*

*Proof.* Let  $\mathbf{u} \in \mathbf{U}$  be solving (6). It follows that  $\langle C(\mathbf{u}) - G, \mathbf{w} \rangle = 0$  for all  $\mathbf{w} \in \mathcal{N}(B)$ . By assumption (i) we infer that  $C(\mathbf{u}) = G$  in  $\mathbf{U}'$ . Therefore,  $\mathbf{u}$  satisfies  $\langle R^{-1}B\mathbf{u}, B\mathbf{w} \rangle = \langle R^{-1}F, B\mathbf{w} \rangle$  for all  $\mathbf{w} \in \mathbf{U}$ , i.e.,  $B^*R^{-1}(B\mathbf{u} - F) = 0$  in  $\mathbf{U}'$ . Using assumption (ii) and recalling that  $R : \mathbf{V} \rightarrow \mathbf{V}'$  is an isomorphism, we conclude that  $B^*R^{-1} : \mathbf{V}' \rightarrow \mathbf{U}'$  is injective, so that  $B\mathbf{u} = F$  in  $\mathbf{V}'$ . The other direction is immediate.  $\square$

**2.2. Discretization.** We analyze approximations of the continuous problems (5) and (6). At the continuous level, these two problems are equivalent by Theorem 2.3 so that considering  $\mathbf{T}_\kappa$  or  $\mathbf{D}_\kappa$  with their respective right-hand side yields the same problem. However, considering one operator or the other at the discrete level is no longer equivalent.

For an index parameter  $h > 0$ , let  $\mathbf{U}_h$  and  $\mathbf{V}_h$  be two (families of) finite-dimensional spaces such that  $\mathbf{U}_h \subset \mathbf{U}$  and  $\mathbf{V}_h \subset \mathbf{V}$ . Of course, later  $h$  will be a mesh parameter. We denote the canonical injection maps by  $\mathbf{i}_h : \mathbf{U}_h \rightarrow \mathbf{U}$  and  $\mathbf{j}_h : \mathbf{V}_h \rightarrow \mathbf{V}$ , with  $\mathbf{i}_h^* : \mathbf{U}' \rightarrow \mathbf{U}'_h$  and  $\mathbf{j}_h^* : \mathbf{V}' \rightarrow \mathbf{V}'_h$  the respective adjoints. The discrete spaces  $\mathbf{U}_h$  and  $\mathbf{V}_h$  are provided with the induced norms  $\|\cdot\|_{\mathbf{U}_h} := \|\mathbf{i}_h \cdot\|_{\mathbf{U}}$  and  $\|\cdot\|_{\mathbf{V}_h} := \|\mathbf{j}_h \cdot\|_{\mathbf{V}}$ .

**2.2.1. Semi-discrete scheme.** One possibility is to discretize the operator  $\mathbf{D}_\kappa$  in the standard way, that is, considering  $\mathbf{D}_{\kappa,h}^* : \mathbf{U}_h \rightarrow \mathbf{U}'_h$  defined as

$$\mathbf{D}_{\kappa,h}^* = \mathbf{i}_h^* \mathbf{D}_\kappa \mathbf{i}_h, \quad \text{with } \mathbf{D}_\kappa = C + \kappa B^* R^{-1} B.$$

This discretization still requires to calculate  $R^{-1}$ , which is not feasible in practice. In DPG discretizations,  $R$  is the Riesz operator  $R_{\mathbf{V}} : \mathbf{V} \rightarrow \mathbf{V}'$  and such a semi-discrete scheme is sometimes called *ideal* DPG method. It is distinguished from the *practical* variant which includes a discretization of  $R_{\mathbf{V}}^{-1}$ , cf. [23].

The operator  $\mathbf{D}_{\kappa,h}^*$  induces the problem

$$(7) \quad \mathbf{u}_h \in \mathbf{U}_h : \quad \langle \mathbf{D}_{\kappa,h}^* \mathbf{u}_h, \mathbf{w}_h \rangle = \kappa \langle R^{-1} F, B \mathbf{i}_h \mathbf{w}_h \rangle + \langle G, \mathbf{i}_h \mathbf{w}_h \rangle \quad \forall \mathbf{w}_h \in \mathbf{U}_h.$$

**THEOREM 2.5.** *Assume that the assumptions from Proposition 2.1 hold true with constants  $\kappa_0$  and  $c_U$  specified there. Then, for all  $\kappa \geq \kappa_0$ ,  $\mathbf{D}_{\kappa,h}^*$  is invertible with uniformly Lipschitz continuous inverse, and problem (7) is well posed. In addition, assuming that  $C$  is Lipschitz continuous with constant  $c_{\text{Lip}}$ , we have the quasi-optimal error estimate*

$$\|\mathbf{u} - \mathbf{i}_h(\mathbf{u}_h)\|_{\mathbf{U}} \leq (1 + c_U^{-1}(c_{\text{Lip}} + \kappa \|B^* R^{-1} B\|_{\mathcal{L}(\mathbf{U}, \mathbf{U}')})) \inf_{\mathbf{w}_h \in \mathbf{U}_h} \|\mathbf{u} - \mathbf{i}_h(\mathbf{w}_h)\|_{\mathbf{U}}.$$

Here,  $\mathbf{u} \in \mathbf{U}$  and  $\mathbf{u}_h \in \mathbf{U}_h$  are the unique solutions of (6) and (7), respectively.

*Proof.* The discrete operator  $\mathbf{D}_{\kappa,h}^*$  defines a conforming approximation of the continuous problem (6). Therefore, its uniform Lipschitz continuous invertibility follows from the Lipschitz continuous invertibility of  $\mathbf{D}_{\kappa}$ , cf. Theorem 2.3. Again, by the conformity of the approximation, the a priori error estimate follows by standard arguments using the monotonicity (3).  $\square$

**2.2.2. Fully discrete scheme.** In order to avoid the inversion of the operator  $R$ , present in (7), we discretize, instead of the Schur complement  $\mathbf{D}_{\kappa}$ , the full operator  $\mathbf{T}_{\kappa}$  as  $\mathbf{T}_{\kappa,h} : \mathbf{U}_h \times \mathbf{V}_h \rightarrow \mathbf{U}'_h \times \mathbf{V}'_h$  defined by

$$(8) \quad \mathbf{T}_{\kappa,h} = \begin{pmatrix} C_h & B_h^* \\ B_h & -\frac{1}{\kappa} R_h \end{pmatrix} = \begin{pmatrix} \mathbf{i}_h^* & 0 \\ 0 & \mathbf{j}_h^* \end{pmatrix} \mathbf{T}_{\kappa} \begin{pmatrix} \mathbf{i}_h & 0 \\ 0 & \mathbf{j}_h \end{pmatrix}.$$

Applying the Schur factorization to this discrete operator we obtain the fully discrete problem

$$(9) \quad \mathbf{u}_h \in \mathbf{U}_h : \quad \langle \mathbf{D}_{\kappa,h} \mathbf{u}_h, \mathbf{w}_h \rangle = \kappa \langle R_h^{-1} \mathbf{j}_h^* F, B_h \mathbf{w}_h \rangle + \langle \mathbf{i}_h^* G, \mathbf{w}_h \rangle \quad \forall \mathbf{w}_h \in \mathbf{U}_h,$$

with

$$(10) \quad \langle \mathbf{D}_{\kappa,h} \mathbf{u}_h, \mathbf{w}_h \rangle := \langle C_h(\mathbf{u}_h), \mathbf{w}_h \rangle + \kappa \langle R_h^{-1} B_h \mathbf{u}_h, B_h \mathbf{w}_h \rangle.$$

The well-posedness of (9) follows similarly as in Proposition 2.1, by using the existence of a Fortin operator  $\Pi : \mathbf{V} \rightarrow \mathbf{V}_h$  satisfying, uniformly in  $h$ ,

$$(11a) \quad \langle B \mathbf{i}_h \mathbf{u}_h, \mathbf{v} - \mathbf{j}_h \Pi \mathbf{v} \rangle = 0 \quad \forall \mathbf{u}_h \in \mathbf{U}_h, \mathbf{v} \in \mathbf{V},$$

$$(11b) \quad \exists c_{\Pi} > 0, \quad \|\mathbf{j}_h \Pi \mathbf{v}\|_{\mathbf{V}}^2 \leq c_{\Pi} \|\mathbf{v}\|_{\mathbf{V}}^2 \quad \forall \mathbf{v} \in \mathbf{V}.$$

In the context of DPG methods, Gopalakrishnan and Qiu [23] have employed such an operator to analyze the approximation of optimal test functions. This is precisely our motivation. Let us recall this result, cf. [23, Proof of Theorem 2.1].

**LEMMA 2.6.** *Assume that (11) holds. Then,  $\|B \mathbf{i}_h \mathbf{u}_h\|_{\mathbf{V}'}^2 \leq c_{\Pi} \|B_h \mathbf{u}_h\|_{\mathbf{V}'_h}^2$ , for all  $\mathbf{u}_h \in \mathbf{U}_h$ .*

The following proposition extends the statements of Proposition 2.1 to the discrete level.

**PROPOSITION 2.7.** *Assume that,*

(i')  $R_h^{-1}$  is coercive on  $\mathcal{R}(B_h)$  with coercivity constant  $c_{R_h}^{-1} > 0$ ,

(ii) there exists  $c_B > 0$  such that  $\|\mathbf{u} - P_B(\mathbf{u})\|^2 \leq c_B \|B \mathbf{u}\|_{\mathbf{V}'}^2$ , for all  $\mathbf{u} \in \mathbf{U}$ ,

(iii) there exist  $c_U, c_V > 0$  such that

$$c_U \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{U}}^2 \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + c_V \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2$$

for all  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbf{U}$ ,

(iv) there exists a Fortin operator  $\Pi : \mathbf{V} \rightarrow \mathbf{V}_h$  satisfying (11).

Then, for all  $\kappa \geq \kappa_{0,h} := (c_V + c_{BCU})c_{\Pi}c_{R_h}$ ,  $\mathbf{T}_{\kappa,h}$  is invertible with Lipschitz continuous inverse. In particular,

$$(12) \quad c_U \|\mathbf{u}_h - \mathbf{v}_h\|_{\mathbf{U}_h}^2 \leq \langle C_h(\mathbf{u}_h) - C_h(\mathbf{v}_h), \mathbf{u}_h - \mathbf{v}_h \rangle + \kappa \langle R_h^{-1} B_h(\mathbf{u}_h - \mathbf{v}_h), B_h(\mathbf{u}_h - \mathbf{v}_h) \rangle$$

for all  $\mathbf{u}_h, \mathbf{v}_h \in \mathbf{U}_h$ .

*Proof.* We follow the same route as in the proof of Proposition 2.1. Specifically, the discrete operator  $\mathbf{T}_{\kappa,h}$  is invertible if and only if  $\mathbf{D}_{\kappa,h}$  is invertible. We prove that  $\mathbf{D}_{\kappa,h}$  is continuous and strongly monotone.

Let  $\mathbf{u}_h, \mathbf{v}_h \in \mathbf{U}_h$ . By assumptions (ii), (iv) and Lemma 2.6 we have

$$\|\mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h) - P_B \mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{U}}^2 \leq c_B \|\mathbf{B} \mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{V}'}^2 \leq c_B c_{\Pi} \|B_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{V}_h'}^2,$$

so that, using assumption (i'), it follows that

$$\|\mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h) - P_B \mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{U}}^2 \leq c_B c_{\Pi} c_{R_h} \langle R_h^{-1} B_h(\mathbf{u}_h - \mathbf{v}_h), B_h(\mathbf{u}_h - \mathbf{v}_h) \rangle.$$

Owing to assumption (iii) and the definition of  $C_h$ , we also have

$$c_U \|P_B \mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{U}}^2 \leq \langle C_h(\mathbf{u}_h) - C_h(\mathbf{v}_h), \mathbf{u}_h - \mathbf{v}_h \rangle + c_V \|\mathbf{B} \mathbf{i}_h(\mathbf{u}_h - \mathbf{v}_h)\|_{\mathbf{V}'}^2.$$

Again applying Lemma 2.6 and combining the two last inequalities, it follows that

$$c_U \|\mathbf{u}_h - \mathbf{v}_h\|_{\mathbf{U}_h}^2 \leq \langle C_h(\mathbf{u}_h) - C_h(\mathbf{v}_h), \mathbf{u}_h - \mathbf{v}_h \rangle + (c_V + c_{BCU})c_{\Pi}c_{R_h} \langle B_h(\mathbf{u}_h - \mathbf{v}_h), R_h^{-1} B_h(\mathbf{u}_h - \mathbf{v}_h) \rangle.$$

In particular, (12) holds for  $\kappa \geq \kappa_{0,h} := (c_V + c_{BCU})c_{\Pi}c_{R_h}$ . We conclude that  $\mathbf{D}_{\kappa,h}$  is continuous and strongly monotone for all  $\kappa \geq \kappa_{0,h}$  with Lipschitz continuous inverse.  $\square$

**THEOREM 2.8.** *Assume that the assumptions from Propositions 2.1, 2.7 hold true, and that  $\kappa \geq \max(\kappa_0; \kappa_{0,h})$ , with the constants  $\kappa_0, \kappa_{0,h}$  and  $c_U$  from before. Then, problems (6) and (9) are well posed. In addition, assuming that the assumptions from Proposition 2.4 hold true and that  $C$  is Lipschitz continuous with Lipschitz constant  $c_{\text{Lip}}$ , we have the quasi-optimal error estimate*

$$\|\mathbf{u} - \mathbf{i}_h(\mathbf{u}_h)\|_{\mathbf{U}} \leq (1 + c_U^{-1} (c_{\text{Lip}} + \kappa \|B^* \mathbf{j}_h R_h^{-1} \mathbf{j}_h^* B\|_{\mathcal{L}(\mathbf{U}, \mathbf{U}')})) \inf_{\mathbf{w}_h \in \mathbf{U}_h} \|\mathbf{u} - \mathbf{i}_h(\mathbf{w}_h)\|_{\mathbf{U}}.$$

Here,  $\mathbf{u} \in \mathbf{U}$  and  $\mathbf{u}_h \in \mathbf{U}_h$  are the unique solutions of (6) and (9), respectively.

*Proof.* Let  $\mathbf{w}_h \in \mathbf{U}_h$  and denote  $\boldsymbol{\xi}_h = \mathbf{u}_h - \mathbf{w}_h$ . By Proposition 2.7 we have

$$c_U \|\boldsymbol{\xi}_h\|_{\mathbf{U}_h}^2 \leq \langle C_h(\mathbf{u}_h) - C_h(\mathbf{w}_h), \boldsymbol{\xi}_h \rangle + \kappa \langle R_h^{-1} B_h(\mathbf{u}_h - \mathbf{w}_h), B_h \boldsymbol{\xi}_h \rangle$$

for all  $\kappa \geq \kappa_{0,h}$ . In addition, since  $B\mathbf{u} = F$  and  $C(\mathbf{u}) = G$  according to Proposition 2.4, the relations

$$\langle R_h^{-1} \mathbf{j}_h^* B \mathbf{u}, B_h \boldsymbol{\xi}_h \rangle = \langle R_h^{-1} \mathbf{j}_h^* F, B_h \boldsymbol{\xi}_h \rangle \quad \text{and} \quad \langle \mathbf{i}_h^* C(\mathbf{u}), \boldsymbol{\xi}_h \rangle = \langle \mathbf{i}_h^* G, \boldsymbol{\xi}_h \rangle.$$

hold. Hence, we conclude that

$$\kappa \langle R_h^{-1} \mathbf{j}_h^* B(\mathbf{u} - \mathbf{i}_h \mathbf{u}_h), B_h \boldsymbol{\xi}_h \rangle + \langle C(\mathbf{u}) - C(\mathbf{i}_h \mathbf{u}_h), \mathbf{j}_h \boldsymbol{\xi}_h \rangle = 0,$$



yielding

$$c_U \|\boldsymbol{\xi}_h\|_{\mathcal{U}_h}^2 \leq \langle C(\mathbf{u}) - C(\mathbf{i}_h \mathbf{w}_h), \mathbf{i}_h \boldsymbol{\xi}_h \rangle + \kappa \langle B^* \mathbf{j}_h R_h^{-1} \mathbf{j}_h^* B(\mathbf{u} - \mathbf{i}_h \mathbf{w}_h), \mathbf{i}_h \boldsymbol{\xi}_h \rangle.$$

The triangle inequality and the Lipschitz continuity of  $C$  then prove the statement.  $\square$

**3. Nonlinear model problem and functional setting.** In the remainder of this paper we show how our abstract DPG framework applies to an advection-diffusion model problem with nonlinear diffusion. In this section we specify the model problem and consider its continuous formulation.

Given a source term  $f$ , let us consider  $u : \Omega \rightarrow \mathbb{R}$ , with  $\Omega \subset \mathbb{R}^d$  ( $d \geq 2$ ) a connected Lipschitz domain, the solution of

$$(13a) \quad -\nabla \cdot (\boldsymbol{\lambda}(\mathbf{x}, |\nabla u(\mathbf{x})|) \nabla u(\mathbf{x}) + \boldsymbol{\beta}(\mathbf{x}) u(\mathbf{x})) = f(\mathbf{x}) \quad \text{for a.a. } \mathbf{x} \in \Omega,$$

$$(13b) \quad u(\mathbf{x}) = 0 \quad \text{for a.a. } \mathbf{x} \in \partial\Omega.$$

Here,  $\boldsymbol{\beta}$  denotes an  $\mathbb{R}^d$ -valued advection field and  $\boldsymbol{\lambda}$  an  $\mathbb{R}^{d \times d}$ -valued diffusion tensor. By  $\partial\Omega$  we denote the boundary of  $\Omega$ , with outwardly oriented unit normal vector  $\mathbf{n}$ . For simplicity, we write  $\boldsymbol{\lambda}(|\nabla u|)$  for  $\boldsymbol{\lambda}(\mathbf{x}, |\nabla u(\mathbf{x})|)$  for almost all  $\mathbf{x} \in \Omega$ .

Owing to the theory of continuous strongly monotone operators (see, e.g., [26, Chap. 2]), this model problem admits a unique solution  $u \in H_0^1(\Omega)$  for any source term  $f \in L^2(\Omega)$  if the physical parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\lambda}$  satisfy the following assumptions:

- $\boldsymbol{\beta}$  is Lipschitz continuous on  $\Omega$  and satisfies

$$(14a) \quad \text{ess inf}_{\Omega} (-\nabla \cdot \boldsymbol{\beta}) \geq 0.$$

- There exist constants  $0 < \lambda_b \leq \lambda_{\sharp}$  such that, for all  $\boldsymbol{\sigma}, \boldsymbol{\theta} \in \mathbf{L}^2(\Omega)$ ,

$$(14b) \quad \lambda_b \|\boldsymbol{\sigma} - \boldsymbol{\theta}\|_{\mathbf{L}^2(\Omega)}^2 \leq (\boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} - \boldsymbol{\lambda}(|\boldsymbol{\theta}|) \boldsymbol{\theta}, \boldsymbol{\sigma} - \boldsymbol{\theta})_{\mathbf{L}^2(\Omega)},$$

$$(14c) \quad \lambda_{\sharp} \|\boldsymbol{\sigma} - \boldsymbol{\theta}\|_{\mathbf{L}^2(\Omega)} \geq \|\boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} - \boldsymbol{\lambda}(|\boldsymbol{\theta}|) \boldsymbol{\theta}\|_{\mathbf{L}^2(\Omega)}.$$

Throughout the remainder of this paper we assume that all these conditions are satisfied, specifically that  $f \in L^2(\Omega)$ .

**3.1. Standard Sobolev spaces and Péclet number.** Let  $L^2(\Omega)$  and  $\mathbf{L}^2(\Omega)$  be the standard Lebesgue spaces collecting  $\mathbb{R}$ -valued and  $\mathbb{R}^d$ -valued functions, respectively, satisfying

$$(v, v)_{\Omega} = \|v\|_{L^2(\Omega)}^2 = \int_{\Omega} |v|^2 < +\infty \quad \text{and} \quad (\mathbf{v}, \mathbf{v})_{\Omega} = \|\mathbf{v}\|_{\mathbf{L}^2(\Omega)}^2 = \int_{\Omega} |\mathbf{v}|^2 < +\infty.$$

We denote by  $H^1(\Omega)$  and  $\mathbf{H}(\text{div}; \Omega)$  the classical Sobolev spaces equipped with the scaled inner products

$$(v, w)_{H^1(\Omega)} := (v, w)_{\Omega} + \ell_{\Omega}^2 (\nabla v, \nabla w)_{\Omega} \quad \forall v, w \in H^1(\Omega),$$

$$(\boldsymbol{\tau}, \boldsymbol{\eta})_{\mathbf{H}(\text{div}; \Omega)} := (\boldsymbol{\tau}, \boldsymbol{\eta})_{\Omega} + \ell_{\Omega}^2 (\nabla \cdot \boldsymbol{\tau}, \nabla \cdot \boldsymbol{\eta})_{\Omega} \quad \forall \boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbf{H}(\text{div}; \Omega).$$

We also denote by  $H_0^1(\Omega)$  the closure of the space collecting infinitely differentiable functions with compact support in  $\Omega$  with the norm  $\|\cdot\|_{H^1(\Omega)}$ .

The characteristic length  $\ell_{\Omega} > 0$  is introduced so that the above inner products are dimensionally coherent. Its definition is arbitrary, but fixed once and for all. To avoid the proliferation of constants, the reference length  $\ell_{\Omega}$  is chosen such that

$$(15) \quad \|\boldsymbol{\beta}\|_{\mathbf{L}^{\infty}(\Omega)} \ell_{\Omega} \lambda_{\sharp}^{-1} = 1.$$

This means that the global Péclet number is of order 1, i.e., the magnitude of the advective and diffusive effects are comparable.

**3.2. Mesh partition and product spaces.** Let  $\Omega_h$  be a non-overlapping partition of  $\Omega$  composed open elements  $T \in \Omega_h$  with Lipschitz boundary  $\partial T$  and outwardly oriented by  $\mathbf{n}_T$ . Let  $H^1(\Omega_h)$  and  $\mathbf{H}(\text{div}; \Omega_h)$  be the product or “broken” Sobolev spaces equipped with inner products

$$\begin{aligned} (v, w)_{H^1(\Omega_h)} &:= (v, w)_{\Omega_h} + \ell_\Omega^2 (\nabla v, \nabla w)_{\Omega_h} \quad \forall v, w \in H^1(\Omega_h), \\ (\boldsymbol{\tau}, \boldsymbol{\eta})_{\mathbf{H}(\text{div}; \Omega_h)} &:= (\boldsymbol{\tau}, \boldsymbol{\eta})_{\Omega_h} + \ell_\Omega^2 (\nabla \cdot \boldsymbol{\tau}, \nabla \cdot \boldsymbol{\eta})_{\Omega_h} \quad \forall \boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbf{H}(\text{div}; \Omega_h). \end{aligned}$$

Here,  $(\cdot, \cdot)_{\Omega_h} = \sum_{T \in \Omega_h} (\cdot, \cdot)_T$  denotes the element-wise  $L^2$ -inner product, that is, appearing differential operators are taken in a piecewise form. For all  $T \in \Omega_h$  we denote by  $H^{1/2}(\partial T)$  and  $H^{-1/2}(\partial T)$  the trace spaces of  $H^1(T)$  and  $\mathbf{H}(\text{div}; T)$ , respectively. They are dual to each other. Traces on the mesh skeleton  $\partial\Omega_h$  are defined with the trace maps  $\gamma : H^1(\Omega_h) \rightarrow \times_{T \in \Omega_h} H^{1/2}(\partial T)$  and  $\gamma_n : \mathbf{H}(\text{div}; \Omega_h) \rightarrow \times_{T \in \Omega_h} H^{-1/2}(\partial T)$ , defined as

$$\begin{aligned} \gamma(u)|_{\partial T} &= u|_{\partial T} \quad \forall T \in \Omega_h \quad \forall u \in H^1(\Omega_h), \\ \gamma_n(\boldsymbol{\rho})|_{\partial T} &= (\boldsymbol{\rho} \cdot \mathbf{n}_T)|_{\partial T} \quad \forall T \in \Omega_h \quad \forall \boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega_h). \end{aligned}$$

The duality product between  $\times_{T \in \Omega_h} H^{-1/2}(\partial T)$  and  $\times_{T \in \Omega_h} H^{1/2}(\partial T)$  is denoted by  $\langle \cdot, \cdot \rangle_{\partial\Omega_h} = \sum_{T \in \Omega_h} \langle \cdot, \cdot \rangle_{\partial T}$  with duality  $\langle \cdot, \cdot \rangle_{\partial T}$  between  $H^{-1/2}(\partial T)$  and  $H^{1/2}(\partial T)$  ( $T \in \Omega_h$ ).

We also introduce the trace spaces

$$\begin{aligned} H_{00}^{1/2}(\partial\Omega_h) &:= \left\{ \hat{u} \in \times_{T \in \Omega_h} H^{1/2}(\partial T) \mid \exists w \in H_0^1(\Omega), \hat{u} = \gamma(w) \right\}, \\ H^{-1/2}(\partial\Omega_h) &:= \left\{ \hat{\rho} \in \times_{T \in \Omega_h} H^{-1/2}(\partial T) \mid \exists \boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega), \hat{\rho} = \gamma_n(\boldsymbol{\rho}) \right\}, \end{aligned}$$

equipped with their respective quotient norms,

$$(16a) \quad \|\hat{v}\|_{H_{00}^{1/2}(\partial\Omega_h)} := \inf_{w \in H_0^1(\Omega)} \{ \|w\|_{H^1(\Omega)}; \gamma(w) = \hat{v} \} \quad \forall \hat{v} \in H_{00}^{1/2}(\partial\Omega_h),$$

$$(16b) \quad \|\hat{\rho}\|_{H^{-1/2}(\partial\Omega_h)} := \inf_{\boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega)} \{ \|\boldsymbol{\rho}\|_{\mathbf{H}(\text{div}; \Omega)}; \gamma_n(\boldsymbol{\rho}) = \hat{\rho} \} \quad \forall \hat{\rho} \in H^{-1/2}(\partial\Omega_h).$$

Finally, we close this section by recalling the Poincaré-Steklov inequality in the product space  $H^1(\Omega_h)$ , cf. [2]. The proof is given for completeness, here including the length scale parameter  $\ell_\Omega$ .

LEMMA 3.1. *We have*

$$c_{PS} \ell_\Omega^{-2} \|v\|_{H^1(\Omega_h)}^2 \leq \|\nabla v\|_{L^2(\Omega_h)}^2 + \left( \sup_{\hat{\tau} \in H^{-1/2}(\partial\Omega_h)} \frac{\langle \hat{\tau}, \gamma(v) \rangle_{\partial\Omega_h}}{\|\hat{\tau}\|_{H^{-1/2}(\partial\Omega_h)}} \right)^2 \quad \forall v \in H^1(\Omega_h),$$

with  $c_{PS}^{-1} = 2(1 + c_{PS,0})$  and  $c_{PS,0} > 0$  the Poincaré-Steklov constant in  $H_0^1(\Omega)$  satisfying  $\|\xi\|_{L^2(\Omega)}^2 \leq c_{PS,0} \ell_\Omega^2 \|\nabla \xi\|_{L^2(\Omega)}^2$  for all  $\xi \in H_0^1(\Omega)$ .

*Proof.* Let  $v \in H^1(\Omega_h)$  and let  $\xi \in H_0^1(\Omega)$  such that  $-\ell_\Omega^2 \Delta \xi = v$ . It holds

$$\begin{aligned} \|v\|_{L^2(\Omega)}^2 &= \ell_\Omega^2 (\nabla v, \nabla \xi)_{\Omega_h} - \ell_\Omega^2 \langle \gamma_n(\nabla \xi), \gamma(v) \rangle_{\partial\Omega_h} \\ &\leq \ell_\Omega^2 \|\nabla v\|_{L^2(\Omega_h)} \|\nabla \xi\|_{L^2(\Omega)} + \ell_\Omega^2 \|\gamma_n(\nabla \xi)\|_{H^{-1/2}(\partial\Omega_h)} \sup_{\hat{\tau} \in H^{-1/2}(\partial\Omega_h)} \frac{\langle \hat{\tau}, \gamma(v) \rangle_{\partial\Omega_h}}{\|\hat{\tau}\|_{H^{-1/2}(\partial\Omega_h)}}. \end{aligned}$$

Furthermore,  $\ell_\Omega^2 \|\nabla \xi\|_{\mathbf{L}^2(\Omega)}^2 \leq c_{PS,0} \|v\|_{L^2(\Omega)}^2$  by the standard Poincaré-Steklov inequality in  $H_0^1(\Omega)$ . Hence, observing that  $\|\gamma_n(\nabla \xi)\|_{H^{-1/2}(\partial\Omega_h)}^2 \leq (1 + c_{PS,0}) \ell_\Omega^{-2} \|v\|_{L^2(\Omega)}^2$ , the statement follows.  $\square$

**4. Penalized variational formulation of the model problem.** In this section we apply the results of Section 2 to devise and approximate a penalized formulation of (13). The objective of our non-standard formulation is to separate the linear and the nonlinear parts of our problem, namely rewriting formally (13a) as

$$-\nabla \cdot (\boldsymbol{\rho} + \boldsymbol{\beta}u) = f, \quad \boldsymbol{\sigma} = \nabla u \quad \text{and} \quad \boldsymbol{\rho} = \boldsymbol{\lambda}(|\boldsymbol{\sigma}|)\boldsymbol{\sigma}.$$

In the following, the operator  $B$  stands for the representation of the ultra-weak formulation of the first two linear equations, and the nonlinear operator  $C$  is used to enforce the nonlinear closure relation  $\boldsymbol{\rho} = \boldsymbol{\lambda}(|\boldsymbol{\sigma}|)\boldsymbol{\sigma}$ .

**4.1. Ultra-weak formulation of the linear part.** We start by specifying a variational formulation of the linear part of the model problem. In this case we select an ultra-weak variant. In some cases like singularly perturbed problems the ultra-weak form has its advantages but for our model problem this selection is not essential.

We consider the following linear problem. Find  $u \in L^2(\Omega)$ ,  $\boldsymbol{\sigma} \in \mathbf{L}^2(\Omega)$ ,  $\boldsymbol{\rho} \in \mathbf{L}^2(\Omega)$ ,  $\hat{u} \in H_{00}^{1/2}(\partial\Omega_h)$ ,  $\hat{\rho} \in H^{-1/2}(\partial\Omega_h)$  such that

$$(17a) \quad (u, \boldsymbol{\beta} \cdot \nabla v)_{\Omega_h} + (\boldsymbol{\rho}, \nabla v)_{\Omega_h} - \langle \hat{\rho}, \gamma(v) \rangle_{\partial\Omega_h} = (f, v)_\Omega \quad \forall v \in H^1(\Omega_h),$$

$$(17b) \quad (\boldsymbol{\sigma}, \boldsymbol{\tau})_\Omega + (u, \nabla \cdot \boldsymbol{\tau})_{\Omega_h} - \langle \gamma_n(\boldsymbol{\tau}), \hat{u} \rangle_{\partial\Omega_h} = 0 \quad \forall \boldsymbol{\tau} \in \mathbf{H}(\text{div}; \Omega_h),$$

and denote the spaces

$$\mathbf{U} := L^2(\Omega) \times \mathbf{L}^2(\Omega) \times \mathbf{L}^2(\Omega) \times H_{00}^{1/2}(\partial\Omega_h) \times H^{-1/2}(\partial\Omega_h),$$

$$\mathbf{V} := H^1(\Omega_h) \times \mathbf{H}(\text{div}; \Omega_h).$$

Define the operator  $B : \mathbf{U} \rightarrow \mathbf{V}'$  as

$$(18) \quad \langle B\mathbf{u}, \mathbf{v} \rangle := (u, \boldsymbol{\beta} \cdot \nabla v + \nabla \cdot \boldsymbol{\tau})_{\Omega_h} + (\boldsymbol{\sigma}, \boldsymbol{\tau})_\Omega + (\boldsymbol{\rho}, \nabla v)_{\Omega_h} \\ - \langle \hat{\rho}, \gamma(v) \rangle_{\partial\Omega_h} - \langle \gamma_n(\boldsymbol{\tau}), \hat{u} \rangle_{\partial\Omega_h}$$

for  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathbf{U}$  and  $\mathbf{v} = (v, \boldsymbol{\tau}) \in \mathbf{V}$ . Problem (17) is then reformulated as

$$(19) \quad \mathbf{u} \in \mathbf{U} : \quad B\mathbf{u} = F \in \mathbf{V}',$$

with  $F \in \mathbf{V}'$  such that  $F : \mathbf{v} \mapsto (f, v)_\Omega$  for all  $\mathbf{v} = (v, \boldsymbol{\tau}) \in \mathbf{V}$ . The following lemma gives a strong characterization of  $\mathbf{u} \in \mathbf{U}$  solving (19).

**LEMMA 4.1.** *Let  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathbf{U}$  be a solution of (19). Then,  $u \in H_0^1(\Omega)$ ,  $\boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega)$ , and  $\boldsymbol{\sigma} = \nabla u$ ,  $-\nabla \cdot (\boldsymbol{\rho} + \boldsymbol{\beta}u) = f$ ,  $\hat{u} = \gamma(u)$  and  $\hat{\rho} = \gamma_n(\boldsymbol{\rho} + \boldsymbol{\beta}u)$ .*

In the following analysis, the Cartesian space  $\mathbf{U}$  is equipped with the scaled norm  $\|\cdot\|_{\mathbf{U}}$  defined for all  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathbf{U}$  as

$$\|\mathbf{u}\|_{\mathbf{U}}^2 := \|u\|_{L^2(\Omega)}^2 + \ell_\Omega^2 \|\boldsymbol{\sigma}\|_{\mathbf{L}^2(\Omega)}^2 + \ell_\Omega^2 \lambda_\#^{-2} \|\boldsymbol{\rho}\|_{\mathbf{L}^2(\Omega)}^2 + \|\hat{u}\|_{H_{00}^{1/2}(\partial\Omega_h)}^2 + \ell_\Omega^2 \lambda_\#^{-2} \|\hat{\rho}\|_{H^{-1/2}(\partial\Omega_h)}^2.$$

Similarly, the space  $\mathbf{V}$  is equipped with the inner product  $(\cdot, \cdot)_{\mathbf{V}}$  defined for  $\mathbf{v} = (v, \boldsymbol{\tau}), \mathbf{w} = (w, \boldsymbol{\eta}) \in \mathbf{V}$  as

$$(\mathbf{v}, \mathbf{w})_{\mathbf{V}} := \lambda_{\sharp}^2 \ell_{\Omega}^{-4} (v, w)_{H^1(\Omega_h)} + \ell_{\Omega}^{-2} (\boldsymbol{\tau}, \boldsymbol{\eta})_{\mathbf{H}(\text{div}; \Omega_h)}.$$

With these norms in  $\mathbf{U}$  and  $\mathbf{V}$ ,  $B$  is uniformly bounded. Owing to the following lemma,  $B$  is surjective or, equivalently,  $B^*$  is injective with closed range in  $\mathbf{U}'$ .

LEMMA 4.2.  $\|B^* \mathbf{v}\|_{\mathbf{U}'}^2 \geq c_B \|\mathbf{v}\|_{\mathbf{V}}^2$  for all  $\mathbf{v} \in \mathbf{V}$  with  $c_B = c_{PS}/2$  and  $c_{PS}$  the Poincaré-Steklov constant from Lemma 3.1.

*Proof.* Let  $\mathbf{v} \in \mathbf{V}$  with  $\mathbf{v} = (v, \boldsymbol{\tau})$  and note that

$$\begin{aligned} \|B^* \mathbf{v}\|_{\mathbf{U}'}^2 &\geq \|\boldsymbol{\beta} \cdot \nabla v + \nabla \cdot \boldsymbol{\tau}\|_{L^2(\Omega_h)}^2 + \ell_{\Omega}^{-2} \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \ell_{\Omega}^{-2} \lambda_{\sharp}^2 \|\nabla v\|_{L^2(\Omega_h)}^2 \\ &\quad + \ell_{\Omega}^{-2} \lambda_{\sharp}^2 \left( \sup_{\hat{\tau} \in H^{-1/2}(\partial\Omega_h)} \frac{|\langle \hat{\tau}, \gamma(v) \rangle_{\partial\Omega_h}|}{\|\hat{\tau}\|_{H^{-1/2}(\partial\Omega_h)}} \right)^2. \end{aligned}$$

Then, using the Poincaré-Steklov inequality in  $H^1(\Omega_h)$  from Lemma 3.1, it follows

$$\begin{aligned} \|B^* \mathbf{v}\|_{\mathbf{U}'}^2 &\geq \|\boldsymbol{\beta} \cdot \nabla v + \nabla \cdot \boldsymbol{\tau}\|_{L^2(\Omega_h)}^2 + \ell_{\Omega}^{-2} \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \frac{1}{2} \ell_{\Omega}^{-2} \lambda_{\sharp}^2 \|\nabla v\|_{L^2(\Omega_h)}^2 \\ &\quad + \frac{1}{2} c_{PS} \ell_{\Omega}^{-4} \lambda_{\sharp}^2 \|v\|_{H^1(\Omega_h)}^2. \end{aligned}$$

Hence, observing that  $\|\boldsymbol{\tau}\|_{\mathbf{H}(\text{div}; \Omega_h)}^2 \leq \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \ell_{\Omega}^2 \|\nabla \cdot \boldsymbol{\tau} + \boldsymbol{\beta} \cdot \nabla v\|_{L^2(\Omega_h)}^2 + \lambda_{\sharp}^2 \|\nabla v\|_{L^2(\Omega_h)}^2$  by the triangle inequality and assumption (15), it follows that

$$\|B^* \mathbf{v}\|_{\mathbf{U}'}^2 \geq \frac{1}{2} \ell_{\Omega}^{-2} \|\boldsymbol{\tau}\|_{\mathbf{H}(\text{div}; \Omega_h)}^2 + \frac{1}{2} c_{PS} \ell_{\Omega}^{-4} \lambda_{\sharp}^2 \|v\|_{H^1(\Omega_h)}^2.$$

Observing that  $c_{PS} \leq 1$  in Lemma 3.1, the statement follows by definition of the norm in  $\mathbf{V}$ .  $\square$

Denoting by  $P_B : \mathbf{U} \rightarrow \mathcal{N}(B)$  the projector of  $\mathbf{U}$  on  $\mathcal{N}(B)$ , the following lemma is consequence of the boundedness of  $B$  and Lemma 4.2. In particular, assumption (ii) of Proposition 2.1 is satisfied.

LEMMA 4.3.  $\|\mathbf{u} - P_B \mathbf{u}\|_{\mathbf{U}}^2 \leq c_B \|B \mathbf{u}\|_{\mathbf{V}'}$ , for all  $\mathbf{u} \in \mathbf{U}$ , with the constant  $c_B$  defined in Lemma 4.2.

*Proof.* By standard arguments the statement is equivalent to the statement of Lemma 4.2. Indeed, by the continuity of  $B : \mathbf{U} \rightarrow \mathbf{V}'$  and the boundedness below of  $B^*$ ,  $B : \mathcal{N}(B)^{\perp} \rightarrow \mathbf{V}'$  is injective with closed range, i.e.,  $B^* : \mathbf{V} \rightarrow (\mathcal{N}(B)^{\perp})'$  is surjective, so that

$$\|B \mathbf{u}\|_{\mathbf{V}'} = \sup_{\mathbf{v} \in \mathbf{V}} \frac{\langle B^* \mathbf{v}, \mathbf{u} - P_B(\mathbf{u}) \rangle}{\|\mathbf{v}\|_{\mathbf{V}}} \geq c_B^{1/2} \sup_{\mathbf{v} \in \mathbf{V}} \frac{\langle B^* \mathbf{v}, \mathbf{u} - P_B(\mathbf{u}) \rangle}{\|B^* \mathbf{v}\|_{\mathbf{U}'}} = c_B^{1/2} \|\mathbf{u} - P_B(\mathbf{u})\|_{\mathbf{U}}$$

for any  $\mathbf{u} \in \mathbf{U}$ .  $\square$

The non-trivial kernel  $\mathcal{N}(B)$  can be represented as the image of the map  $E : H_0^1(\Omega) \times \mathbf{H} \rightarrow \mathbf{U}$  defined, for all  $\psi \in H_0^1(\Omega)$  and  $\boldsymbol{\eta} \in \mathbf{H} = \{\mathbf{v} \in \mathbf{H}(\text{div}; \Omega) \mid \nabla \cdot \mathbf{v} = 0\}$ , as

$$(20) \quad E(\psi, \boldsymbol{\eta}) = (\psi, \nabla \psi, \boldsymbol{\eta} - \boldsymbol{\beta} \psi, \gamma(\psi), \gamma_n(\boldsymbol{\eta})).$$

LEMMA 4.4. The map  $E : H_0^1(\Omega) \times \mathbf{H} \rightarrow \mathbf{U}$  satisfies

$$(21) \quad c_E \|E(\psi, \boldsymbol{\eta})\|_{\mathbf{U}}^2 \leq \ell_{\Omega}^2 \|\nabla \psi\|_{L^2(\Omega)}^2 + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\boldsymbol{\eta} - \boldsymbol{\beta} \psi\|_{L^2(\Omega)}^2 \quad \forall \psi \in H_0^1(\Omega), \boldsymbol{\eta} \in \mathbf{H},$$

with  $c_E^{-1} = 4c_{PS,0} + 3$  and  $c_{PS,0}$  the Poincaré-Steklov constant defined in Lemma 3.1. In addition,  $E(H_0^1(\Omega) \times \mathbf{H}) = \mathcal{N}(B)$ .

*Proof.* Estimate (21) is consequence of the Poincaré-Steklov inequality in  $H_0^1(\Omega)$  and the triangle inequality. Let us prove that  $E(H_0^1(\Omega) \times \mathbf{H}) = \mathcal{N}(B)$ . Consider  $\psi \in H_0^1(\Omega)$ ,  $\boldsymbol{\eta} \in \mathbf{H}$  and let  $\mathbf{w} = E(\psi, \boldsymbol{\eta})$ . For any  $\mathbf{v} = (v, \boldsymbol{\tau}) \in \mathbf{V}$  it follows that

$$\begin{aligned} \langle B\mathbf{w}, \mathbf{v} \rangle &= (\psi, \nabla \cdot \boldsymbol{\tau})_{\Omega_h} + (\nabla \psi, \boldsymbol{\tau})_{\Omega} + (\boldsymbol{\eta}, \nabla v)_{\Omega_h} \\ &\quad - \langle \gamma_n(\boldsymbol{\eta}), \gamma(v) \rangle_{\partial\Omega_h} - \langle \gamma_n(\boldsymbol{\tau}), \gamma(\psi) \rangle_{\partial\Omega_h}. \end{aligned}$$

Hence, integrating by parts we obtain  $\langle B\mathbf{w}, \mathbf{v} \rangle = 0$  so that  $\mathbf{w} \in \mathcal{N}(B)$ . Conversely, assume that  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathcal{N}(B)$ . Then, computing  $\langle B\mathbf{u}, (0, \boldsymbol{\tau}) \rangle$  we deduce that  $u \in H^1(\Omega_h)$  and that  $\boldsymbol{\sigma}|_T = (\nabla u)|_T$  for all  $T \in \Omega_h$ . We conclude that

$$0 = (u, \nabla \cdot \boldsymbol{\tau})_{\Omega_h} + (\boldsymbol{\sigma}, \boldsymbol{\tau})_{\Omega} - \langle \gamma_n(\boldsymbol{\tau}), \hat{u} \rangle_{\partial\Omega_h} = \langle \gamma_n(\boldsymbol{\tau}), \gamma(u) \rangle_{\partial\Omega_h} - \langle \gamma_n(\boldsymbol{\tau}), \hat{u} \rangle_{\partial\Omega_h},$$

so that  $\gamma(u) = \hat{u} \in H_{00}^{1/2}(\partial\Omega_h)$ . Then, by definition of  $H_{00}^{1/2}(\partial\Omega_h)$ , we infer that  $u \in H_0^1(\Omega)$ . Proceeding similarly by testing with  $\mathbf{v} = (v, \mathbf{0})$ , we infer that  $\boldsymbol{\beta}u + \boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega_h)$  with  $(\nabla \cdot (\boldsymbol{\beta}u + \boldsymbol{\rho}))|_T = 0$  for all  $T \in \Omega_h$ . It follows that  $\hat{\rho} = \gamma_n(\boldsymbol{\beta}u + \boldsymbol{\rho}) \in H^{-1/2}(\partial\Omega_h)$  so that  $\boldsymbol{\beta}u + \boldsymbol{\rho} \in \mathbf{H}(\text{div}; \Omega)$ , and then  $\boldsymbol{\beta}u + \boldsymbol{\rho} \in \mathbf{H}$ . As a result, there exists  $\boldsymbol{\eta} \in \mathbf{H}$  such that  $\boldsymbol{\rho} = -\boldsymbol{\beta}u + \boldsymbol{\eta}$ . To conclude, if  $\mathbf{u} \in \mathcal{N}(B)$ , then  $u \in H_0^1(\Omega)$ ,  $\boldsymbol{\rho} + \boldsymbol{\beta}u \in \mathbf{H}$  and  $\mathbf{u} = E(u, \boldsymbol{\rho} + \boldsymbol{\beta}u)$ .  $\square$

**4.2. Nonlinear penalty term.** In this section we devise a nonlinear penalty form to enforce the closure relation  $\boldsymbol{\rho} = \boldsymbol{\lambda}(|\boldsymbol{\sigma}|)\boldsymbol{\sigma}$  and to control  $\mathcal{N}(B)$ . To simplify the presentation let us introduce  $\pi : \mathbf{U} \rightarrow \mathbf{L}^2(\Omega) \times \mathbf{L}^2(\Omega)$  by defining  $\pi(\mathbf{u}) := (\boldsymbol{\sigma}, \boldsymbol{\rho})$  for  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathbf{U}$ . Then we define a nonlinear operator  $C : \mathbf{U} \rightarrow \mathbf{U}'$  by

$$(22) \quad \langle C(\mathbf{u}), \mathbf{v} \rangle := \ell_{\Omega}^2 \lambda_{\sharp}^{-2} (\boldsymbol{\lambda}(|\boldsymbol{\sigma}|)\boldsymbol{\sigma} - \boldsymbol{\rho}, \alpha \lambda_b \boldsymbol{\theta} - \boldsymbol{\eta})_{\Omega} \quad \forall \mathbf{u}, \mathbf{v} \in \mathbf{U},$$

with  $\pi(\mathbf{u}) = (\boldsymbol{\sigma}, \boldsymbol{\rho})$  and  $\pi(\mathbf{v}) = (\boldsymbol{\theta}, \boldsymbol{\eta})$ . Here,  $\alpha > 0$  denotes a stability parameter that will be chosen greater than  $\lambda_{\sharp/b}^2$ , with  $\lambda_{\sharp/b} := \lambda_{\sharp}/\lambda_b$  the diffusive anisotropy ratio.

We start by establishing the Lipschitz continuity of  $C$ .

LEMMA 4.5.  $\|C(\mathbf{u}_1) - C(\mathbf{u}_2)\|_{\mathbf{U}'} \leq c_{\text{Lip}} \|\mathbf{u}_1 - \mathbf{u}_2\|_{\mathbf{U}}$  for all  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbf{U}$ , with  $c_{\text{Lip}} = 2\max(\alpha \lambda_{\sharp/b}^{-1}, 1)$ .

*Proof.* Let  $\mathbf{u}_1, \mathbf{u}_2, \mathbf{v} \in \mathbf{U}$  such that  $\pi(\mathbf{u}_1) = (\boldsymbol{\sigma}_1, \boldsymbol{\rho}_1)$ ,  $\pi(\mathbf{u}_2) = (\boldsymbol{\sigma}_2, \boldsymbol{\rho}_2)$  and  $\pi(\mathbf{v}) = (\boldsymbol{\theta}, \boldsymbol{\eta})$ . Owing to (22), it follows that

$$\begin{aligned} |\langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{v} \rangle| &\leq \ell_{\Omega}^2 \lambda_{\sharp}^{-2} |(\boldsymbol{\lambda}(|\boldsymbol{\sigma}_1|)\boldsymbol{\sigma}_1 - \boldsymbol{\lambda}(|\boldsymbol{\sigma}_2|)\boldsymbol{\sigma}_2, \alpha \lambda_b \boldsymbol{\theta} - \boldsymbol{\eta})_{\Omega}| \\ &\quad + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} |(\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \alpha \lambda_b \boldsymbol{\theta} - \boldsymbol{\eta})_{\Omega}|. \end{aligned}$$

The Cauchy-Schwarz inequality and assumption (14c) yield

$$\begin{aligned} |\langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{v} \rangle| &\leq \ell_{\Omega}^2 \lambda_{\sharp}^{-1} \left( \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)} + \lambda_{\sharp}^{-1} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)} \right) \|\alpha \lambda_b \boldsymbol{\theta} - \boldsymbol{\eta}\|_{\mathbf{L}^2(\Omega)}. \end{aligned}$$

Observing that  $\|\alpha \lambda_b \boldsymbol{\theta} - \boldsymbol{\eta}\|_{\mathbf{L}^2(\Omega)} \leq 2\lambda_{\sharp} \ell_{\Omega}^{-1} \max(\alpha \lambda_{\sharp/b}^{-1}, 1) \|\mathbf{v}\|_{\mathbf{U}}$ , the desired estimate follows.  $\square$

The next lemma verifies assumption (iii) of Proposition 2.1 for the model problem.

LEMMA 4.6. *Assume that  $\alpha > \lambda_{\sharp}^2/b$ . Then*

$$c_U \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{U}}^2 \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + c_V \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2, \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathbf{U},$$

where  $c_U$  and  $c_V$  are given by

$$c_U = \frac{c_E}{2} \min\left(\frac{1}{2}, \alpha \lambda_{\sharp}^{-2/b} - 1\right), \quad c_V = \frac{c_B}{2} \min\left(\frac{1}{2}, \alpha \lambda_{\sharp}^{-2/b} - 1\right) + \frac{5}{2} c_B \lambda_{\sharp}^{-2} \alpha^2,$$

with  $c_B$  and  $c_E$  the constants defined in Lemmata 4.2 and 4.4, respectively.

*Proof.* Let  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbf{U}$  be such that  $\pi(\mathbf{u}_1) = (\boldsymbol{\sigma}_1, \boldsymbol{\rho}_1)$  and  $\pi(\mathbf{u}_2) = (\boldsymbol{\sigma}_2, \boldsymbol{\rho}_2)$ . By Lemma 4.4 there exist  $\psi_1, \psi_2 \in H_0^1(\Omega)$  and  $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2 \in \mathbf{H}$  such that

$$(23a) \quad P_B(\mathbf{u}_1) = E(\psi_1, \boldsymbol{\eta}_1) = (\psi_1, \nabla \psi_1, \boldsymbol{\eta}_1 - \boldsymbol{\beta} \psi_1, \gamma(\psi_1), \gamma_n(\boldsymbol{\eta}_1)),$$

$$(23b) \quad P_B(\mathbf{u}_2) = E(\psi_2, \boldsymbol{\eta}_2) = (\psi_2, \nabla \psi_2, \boldsymbol{\eta}_2 - \boldsymbol{\beta} \psi_2, \gamma(\psi_2), \gamma_n(\boldsymbol{\eta}_2)).$$

Owing to the continuity estimate (21) from Lemma 4.4, it follows that

$$c_E \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{U}}^2 \leq \ell_{\Omega}^2 \|\nabla \psi_1 - \nabla \psi_2\|_{\mathbf{L}^2(\Omega)}^2 + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\bar{\boldsymbol{\eta}}_1 - \bar{\boldsymbol{\eta}}_2\|_{\mathbf{L}^2(\Omega)}^2,$$

where we have denoted  $\bar{\boldsymbol{\eta}}_i = \boldsymbol{\eta}_i - \boldsymbol{\beta} \psi_i$  for  $i = 1, 2$ . By the triangle inequality we have

$$\begin{aligned} c_E \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{U}}^2 &\leq \ell_{\Omega}^2 \|\nabla \psi_1 - \nabla \psi_2 - (\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)\|_{\mathbf{L}^2(\Omega)}^2 + \ell_{\Omega}^2 \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 \\ &\quad + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\bar{\boldsymbol{\eta}}_1 - \bar{\boldsymbol{\eta}}_2 - (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2)\|_{\mathbf{L}^2(\Omega)}^2 + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2. \end{aligned}$$

Observing that  $\|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2 - (\nabla \psi_1 - \nabla \psi_2)\|_{\mathbf{L}^2(\Omega)}^2 \leq c_B \ell_{\Omega}^{-2} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2$ , and  $\|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2 - (\bar{\boldsymbol{\eta}}_1 - \bar{\boldsymbol{\eta}}_2)\|_{\mathbf{L}^2(\Omega)}^2 \leq c_B \lambda_{\sharp}^2 \ell_{\Omega}^{-2} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2$ , owing to Lemma 4.3, we obtain

$$(24) \quad c_E \|P_B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{U}}^2 \leq c_B \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2 + \ell_{\Omega}^2 \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2.$$

It remains to prove that there exist constants  $c, c' > 0$  such that

$$(25) \quad c \left( \ell_{\Omega}^2 \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \ell_{\Omega}^2 \lambda_{\sharp}^{-2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2 \right) \leq \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle + c' \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2.$$

Let  $F(\mathbf{u}_1, \mathbf{u}_2) = \ell_{\Omega}^{-2} \lambda_{\sharp}^2 \langle C(\mathbf{u}_1) - C(\mathbf{u}_2), \mathbf{u}_1 - \mathbf{u}_2 \rangle$  and observe that

$$\begin{aligned} F(\mathbf{u}_1, \mathbf{u}_2) &= \alpha \lambda_b (\boldsymbol{\lambda}(|\boldsymbol{\sigma}_1|) \boldsymbol{\sigma}_1 - \boldsymbol{\lambda}(|\boldsymbol{\sigma}_2|) \boldsymbol{\sigma}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)_{\Omega} + (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\rho}_1 - \boldsymbol{\rho}_2)_{\Omega} \\ &\quad - (\boldsymbol{\lambda}(|\boldsymbol{\sigma}_1|) \boldsymbol{\sigma}_1 - \boldsymbol{\lambda}(|\boldsymbol{\sigma}_2|) \boldsymbol{\sigma}_2, \boldsymbol{\rho}_1 - \boldsymbol{\rho}_2)_{\Omega} - \alpha \lambda_b (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)_{\Omega}. \end{aligned}$$

Owing to assumptions (14b) and (14c), it follows that

$$\begin{aligned} F(\mathbf{u}_1, \mathbf{u}_2) &\geq \alpha \lambda_b^2 \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2 \\ &\quad - \lambda_{\sharp} \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)} - \alpha \lambda_b (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)_{\Omega}. \end{aligned}$$

Hence, applying Young's inequality, we conclude that

$$F(\mathbf{u}_1, \mathbf{u}_2) \geq \left( \alpha \lambda_b^2 - \frac{\lambda_\#^2}{2} \right) \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \frac{1}{2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2 - \alpha \lambda_b (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)_\Omega.$$

The last term above can be bounded from below by using the representation (23) of  $\mathcal{N}(B)$ . First, we observe that

$$\begin{aligned} (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2)_\Omega &= (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2, \boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2 - (\nabla\psi_1 - \nabla\psi_2))_\Omega + (\bar{\boldsymbol{\eta}}_1 - \bar{\boldsymbol{\eta}}_2, \nabla\psi_1 - \nabla\psi_2)_\Omega \\ &\quad + (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2 - (\bar{\boldsymbol{\eta}}_1 - \bar{\boldsymbol{\eta}}_2), \nabla\psi_1 - \nabla\psi_2)_\Omega =: T_1 + T_2 + T_3. \end{aligned}$$

Applying the Cauchy-Schwarz and triangle inequalities, and Lemma 4.3 for  $T_1$  and  $T_3$ , it follows

$$|T_1| \leq c_B^{1/2} \ell_\Omega^{-1} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)},$$

and

$$\begin{aligned} |T_3| &\leq c_B^{1/2} \lambda_\# \ell_\Omega^{-1} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'} \|\nabla\psi_1 - \nabla\psi_2\|_{\mathbf{L}^2(\Omega)} \\ &\leq c_B^{1/2} \lambda_\# \ell_\Omega^{-1} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'} \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)} + c_B \lambda_\# \ell_\Omega^{-2} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2. \end{aligned}$$

Recalling the definition  $\bar{\boldsymbol{\eta}}_i = -\boldsymbol{\beta}\psi_i + \boldsymbol{\eta}_i$ ,  $i = 1, 2$ , with  $\boldsymbol{\eta}_i \in \mathbf{H} = (\nabla H_0^1(\Omega))^\perp$ , we infer that

$$T_2 = -(\boldsymbol{\beta}\psi_1 - \boldsymbol{\beta}\psi_2, \nabla\psi_1 - \nabla\psi_2)_\Omega \leq 0$$

owing to assumption (14a). Collecting the previous estimates we obtain

$$\begin{aligned} \left( \alpha \lambda_b^2 - \frac{\lambda_\#^2}{2} \right) \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \frac{1}{2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|^2 &\leq F(\mathbf{u}_1, \mathbf{u}_2) \\ &\quad + c_B \alpha \lambda_b^2 \lambda_{\# / b} \ell_\Omega^{-2} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2 \\ &\quad + c_B^{1/2} \alpha \lambda_b \ell_\Omega^{-1} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)} \\ &\quad + c_B^{1/2} \alpha \lambda_b^2 \lambda_{\# / b} \ell_\Omega^{-1} \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'} \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}. \end{aligned}$$

Young's inequality and assumption  $\alpha > \lambda_{\# / b}^2$  yield

$$\frac{1}{2} (\alpha \lambda_b^2 - \lambda_\#^2) \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \frac{1}{4} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|^2 \leq F(\mathbf{u}_1, \mathbf{u}_2) + \frac{5}{2} c_B \lambda_b^2 \ell_\Omega^{-2} \alpha^2 \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2.$$

As a result, multiplying by  $\lambda_\#^{-2}$ , we obtain

$$\begin{aligned} \frac{1}{2} \min \left( \frac{1}{2}, \alpha \lambda_{\# / b}^{-2} - 1 \right) \left( \|\boldsymbol{\sigma}_1 - \boldsymbol{\sigma}_2\|_{\mathbf{L}^2(\Omega)}^2 + \lambda_\#^{-2} \|\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2\|_{\mathbf{L}^2(\Omega)}^2 \right) &\leq \lambda_\#^{-2} F(\mathbf{u}_1, \mathbf{u}_2) \\ &\quad + \frac{5}{2} c_B \ell_\Omega^{-2} \lambda_{\# / b}^{-2} \alpha^2 \|B(\mathbf{u}_1 - \mathbf{u}_2)\|_{\mathbf{V}'}^2. \end{aligned}$$

The statement follows multiplying (24) by  $\frac{1}{2} \min \left( \frac{1}{2}, \alpha \lambda_{\# / b}^{-2} - 1 \right)$ .  $\square$

**4.3. Penalized variational formulation and well-posedness.** We are now in a position to present our penalized variational formulation of the model problem and to prove its well-posedness.

Let  $R_V : \mathbf{V} \rightarrow \mathbf{V}'$  be the Riesz operator, i.e.,  $\langle R_V \mathbf{v}, \mathbf{w} \rangle = (\mathbf{v}, \mathbf{w})_V$  for  $\mathbf{v}, \mathbf{w} \in \mathbf{V}$ , and, for  $\kappa \geq 0$ , let  $\mathbf{T}_\kappa : \mathbf{U} \times \mathbf{V} \rightarrow \mathbf{U}' \times \mathbf{V}'$  be the nonlinear operator

$$(26) \quad \mathbf{T}_\kappa := \begin{pmatrix} C & B^* \\ B & -\frac{1}{\kappa} R_V \end{pmatrix}.$$

Here,  $B : \mathbf{U} \rightarrow \mathbf{V}'$  and  $C : \mathbf{U} \rightarrow \mathbf{U}'$  are the linear and nonlinear operators defined by (18) and (22), respectively.

Following Section 2, problem (13) has the variational formulations

$$(27) \quad (\mathbf{u}, \mathbf{v}) \in \mathbf{U} \times \mathbf{V} : \quad \mathbf{T}_\kappa(\mathbf{u}, \mathbf{v}) = (0, F) \quad \text{in } \mathbf{U}' \times \mathbf{V}'$$

and

$$(28) \quad \mathbf{u} \in \mathbf{U} : \quad \kappa \langle R_V^{-1} B \mathbf{u}, B \mathbf{v} \rangle + \langle C(\mathbf{u}), \mathbf{v} \rangle = \kappa \langle R_V^{-1} F, B \mathbf{v} \rangle \quad \forall \mathbf{v} \in \mathbf{U}.$$

In addition, since  $R_V^{-1}$  is self-adjoint, (28) can be reformulated as

$$(29) \quad \mathbf{u} \in \mathbf{U} : \quad \kappa \langle B \mathbf{u}, \Theta \mathbf{v} \rangle + \langle C(\mathbf{u}), \mathbf{v} \rangle = \kappa \langle F, \Theta \mathbf{v} \rangle_\Omega \quad \forall \mathbf{v} \in \mathbf{U}.$$

Here, we have used the so-called *trial-to-test operator*  $\Theta := R_V^{-1} B : \mathbf{U} \rightarrow \mathbf{V}$ , cf., e.g., [12].

Selecting  $\kappa$  sufficiently large, both problems are well posed and equivalent.

**THEOREM 4.7.** *Assume that  $\alpha > \lambda_{\sharp}^2/b$ , and  $\kappa \geq c_V + c_{BCU}$  with  $c_B, c_U, c_V$  defined in Lemmata 4.2 and 4.6. Then*

- (a) *problems (27) and (29) are equivalent,*
- (b)  *$\mathbf{T}_\kappa : \mathbf{U} \times \mathbf{V} \rightarrow \mathbf{U}' \times \mathbf{V}'$  defined by (26) is invertible with Lipschitz continuous inverse,*
- (c) *the solution  $\mathbf{u} \in \mathbf{U}$  of (29) satisfies  $B \mathbf{u} = F$  in  $\mathbf{V}'$  and  $C(\mathbf{u}) = 0$  in  $\mathbf{U}'$ .*

*Proof.* Statement (a) is a direct consequence of Theorem 2.3 and the definition of  $\Theta$ . The second statement (b) follows from Proposition 2.1. Indeed,  $\langle B \mathbf{v}, R_V^{-1} B \mathbf{v} \rangle = \|B \mathbf{v}\|_{V'}^2$  for all  $\mathbf{v} \in \mathbf{V}$ , so that assumption (i) of Proposition 2.1 is satisfied with coercivity constant  $c_R = 1$ . Since assumptions (ii) and (iii) of Proposition 2.1 hold by Lemmata 4.3 and 4.6, we infer statement (b). Finally, statement (c) is a consequence of Proposition 2.4. The operator  $B$  being surjective by Lemma 4.2, it remains to prove that

$$\langle C(\mathbf{u}), \mathbf{v} \rangle = 0 \quad \forall \mathbf{v} \in \mathcal{N}(B) \quad \implies \quad C(\mathbf{u}) = 0 \text{ in } \mathbf{U}'.$$

In (28) we select  $\mathbf{v} = (0, \mathbf{0}, \boldsymbol{\eta}, 0, \gamma_n(\boldsymbol{\eta})) \in \mathcal{N}(B)$  with  $\boldsymbol{\eta} \in \mathbf{H}$ , giving

$$\langle C(\mathbf{u}), \mathbf{v} \rangle_{\mathbf{U}', \mathbf{U}} = -\ell_\Omega^2 \lambda_{\sharp}^{-2} (\boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} - \boldsymbol{\rho}, \boldsymbol{\eta})_\Omega = 0,$$

since  $\Theta \mathbf{v} = R_V^{-1} B \mathbf{v} = 0$ . Therefore,  $\boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} - \boldsymbol{\rho} \in \mathbf{H}^\perp = \nabla H_0^1(\Omega)$  since  $\nabla H_0^1(\Omega)$  is closed in  $\mathbf{L}^2(\Omega)$ . That is, there exists  $\psi \in H_0^1(\Omega)$  such that  $\nabla \psi = \boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} - \boldsymbol{\rho}$ . Next, choosing  $\mathbf{v} = (\varphi, \nabla \varphi, -\beta \varphi, \gamma(\varphi), 0) \in \mathcal{N}(B)$  with  $\varphi \in H_0^1(\Omega)$ , it follows that  $(\nabla \psi, \alpha \lambda_b \nabla \varphi + \beta \varphi)_\Omega = 0$  for all  $\varphi \in H_0^1(\Omega)$ . As a result, owing to assumption (14a), we infer that  $\psi = 0$ , and then  $\boldsymbol{\lambda}(|\boldsymbol{\sigma}|) \boldsymbol{\sigma} = \boldsymbol{\rho}$ , so that  $C(\mathbf{u}) = 0$  in  $\mathbf{U}'$ .  $\square$



**5. Relaxed DPG scheme.** With all the preparations at hand, the formulation of our DPG scheme for the model problem is immediate and the proof of its quasi-optimal convergence is straightforward.

Considering the continuous problem (27), which is equivalent to (29) in the sense of Theorem 2.3, we follow the presentation of Section 2 and consider two finite-dimensional spaces  $\mathbf{U}_h$  and  $\mathbf{V}_h$  with  $\mathbf{U}_h \subset \mathbf{U}$  and  $\mathbf{V}_h \subset \mathbf{V}$ . Our discrete problem is

$$(30) \quad \mathbf{u}_h \in \mathbf{U}_h : \quad \langle C_h(\mathbf{u}_h), \mathbf{w}_h \rangle + \kappa \langle B_h \mathbf{u}_h, \Theta_h \mathbf{w}_h \rangle = \kappa \langle \mathbf{j}_h^* F, \Theta_h \mathbf{w}_h \rangle \quad \forall \mathbf{w}_h \in \mathbf{V}_h,$$

with

$$C_h := \mathbf{i}_h^* C \mathbf{i}_h, \quad B_h := \mathbf{j}_h^* B \mathbf{i}_h, \quad R_{\mathbf{V}_h} = \mathbf{j}_h^* R_{\mathbf{V}} \mathbf{j}_h, \quad \text{and} \quad \Theta_h = R_{\mathbf{V}_h}^{-1} B_h.$$

The well-posedness of this problem follows from Proposition 2.7.

**THEOREM 5.1.** *Assume that there exists a Fortin operator  $\Pi : \mathbf{V} \rightarrow \mathbf{V}_h$  satisfying (11). Then, for all  $\kappa \geq (c_{\mathbf{V}} + c_B c_{\mathbf{U}}) c_{\Pi}$ , with  $c_B, c_{\mathbf{U}}, c_{\mathbf{V}}$  defined in Lemmata 4.2 and 4.6, (30) is well-posed and converges quasi-optimally,*

$$\|\mathbf{u} - \mathbf{i}_h(\mathbf{u}_h)\|_{\mathbf{U}} \leq (1 + c_{\mathbf{U}}^{-1} (c_{\text{Lip}} + \kappa \|B^* \mathbf{j}_h R_{\mathbf{V}_h}^{-1} \mathbf{j}_h^* B\|_{\mathcal{L}(\mathbf{U}, \mathbf{U}')})) \inf_{\mathbf{w}_h \in \mathbf{U}_h} \|\mathbf{u} - \mathbf{i}_h(\mathbf{w}_h)\|_{\mathbf{U}},$$

with  $c_{\text{Lip}} > 0$  defined in Lemma 4.5. Here,  $\mathbf{u} \in \mathbf{U}$  and  $\mathbf{u}_h \in \mathbf{U}_h$  are the unique solution of (28) and (30), respectively.

*Proof.* The well-posedness of (30) is a consequence of Proposition 2.7. Indeed, assumptions (ii) and (iii) hold by Lemmata 4.3 and 4.6, respectively, and (i') holds with  $c_{R_h} = 1$  since  $R_{\mathbf{V}_h}$  satisfies  $\|B_h \mathbf{v}_h\|_{\mathbf{V}_h}^2 = \langle B_h \mathbf{v}_h, R_{\mathbf{V}_h}^{-1} B_h \mathbf{v}_h \rangle$  for all  $\mathbf{v}_h \in \mathbf{U}_h$ . The error estimate is finally a consequence of Theorem 2.8 using the statement (c) from Theorem 4.7 and the Lipschitz continuity of  $C$  from Lemma 4.5.  $\square$

**6. Numerical example.** In this section, we present some numerical results of a lowest-order implementation of our nonlinear DPG scheme (30) for a model problem with and without advective field  $\beta$ . The specific discretization including Fortin operator is presented in the first subsection, and numerical results are reported afterwards.

**6.1. Discrete setting and Fortin operator.** We use lowest-order test and trial spaces  $\mathbf{U}_h, \mathbf{V}_h$  defined by

$$(31a) \quad \mathbf{U}_h = \mathbb{P}_0(\Omega_h; \mathbb{R}) \times \mathbb{P}_0(\Omega_h; \mathbb{R}^2) \times \mathbb{P}_0(\Omega_h; \mathbb{R}^2) \times \mathbb{P}_1^0(\partial\Omega_h; \mathbb{R}) \times \mathbb{P}_0(\partial\Omega_h; \mathbb{R}),$$

$$(31b) \quad \mathbf{V}_h = \mathbb{P}_2(\Omega_h; \mathbb{R}) \times \mathbb{P}_2(\Omega_h; \mathbb{R}^2).$$

Here,  $\mathbb{P}_k(\Omega_h; \mathbb{R}^d)$  denotes the spaces of  $\Omega_h$ -piecewise  $d$ -variate polynomials of degree  $k$  (meshes are defined below) and  $\mathbb{P}_1^0(\partial\Omega_h; \mathbb{R}) \subset \mathbb{P}_1(\partial\Omega_h; \mathbb{R})$  is the largest subspace of ‘‘continuous’’ functions satisfying the homogeneous Dirichlet condition, that is,  $\gamma(\mathbb{P}_1^0(\partial\Omega_h; \mathbb{R})) \subset H_{00}^{1/2}(\partial\Omega_h)$ . The test space  $\mathbf{V}_h$  is chosen such that there exists a Fortin operator  $\Pi : \mathbf{V} \rightarrow \mathbf{V}_h$  satisfying (11), ensuring that (30) is well posed by Theorem 5.1.

**LEMMA 6.1.** *Let  $\mathbf{U}_h$  and  $\mathbf{V}_h$  be defined by (31), and consider  $\beta \in \mathbb{P}_0(\Omega_h; \mathbb{R}^2)$ . Then, there exists a Fortin operator  $\Pi : \mathbf{V} \rightarrow \mathbf{V}_h$  satisfying (11).*

*Proof.* The statement is consequence of [23, Lemmata 3.2, 3.3]. Indeed, as proved in [23], there are continuous operators  $\Pi^g : H^1(\Omega_h) \rightarrow \mathbb{P}_2(\Omega_h; \mathbb{R})$  and  $\Pi^d : \mathbf{H}(\text{div}; \Omega_h) \rightarrow \mathbb{P}_2(\Omega_h; \mathbb{R}^2)$  such that

$$(32a) \quad \langle \gamma(\Pi^g v - v)|_{\partial T}, q \rangle_{\partial T} = 0 \quad \forall q \in \mathbb{P}_0(\{F\}_{F \in \mathcal{F}_T}; \mathbb{R}),$$

$$(32b) \quad ((\Pi^d \boldsymbol{\tau} - \boldsymbol{\tau})|_T, \mathbf{q})_{\mathbf{L}^2(T)} = 0 \quad \forall \mathbf{q} \in \mathbb{P}_0(T; \mathbb{R}^2),$$

$$(32c) \quad \langle \gamma_n(\Pi_T^d \boldsymbol{\tau} - \boldsymbol{\tau})|_{\partial T}, \hat{q} \rangle_{\partial T} = 0 \quad \forall \hat{q} \in \mathbb{P}_1(\{F\}_{F \in \mathcal{F}_T}; \mathbb{R}) \cap C^0(\partial T)$$

for all  $v \in H^1(\Omega_h)$ ,  $\boldsymbol{\tau} \in \mathbf{H}(\text{div}; \Omega_h)$ , and  $T \in \Omega_h$ . Here,  $\mathcal{F}_T$  denotes the set of faces of an element  $T \in \Omega_h$ . The operator  $\Pi = (\Pi^g, \Pi^d) : \mathbf{V} \rightarrow \mathbf{V}_h$  is continuous and satisfies assumption (11a). In fact, considering  $\mathbf{u}_h = (u_h, \boldsymbol{\sigma}_h, \boldsymbol{\rho}_h, \hat{u}_h, \hat{\rho}_h) \in \mathbf{U}_h$ , we have

$$\begin{aligned} \langle B\mathbf{i}_h \mathbf{u}_h, \mathbf{j}_h \Pi \mathbf{v} - \mathbf{v} \rangle &= (u_h, \boldsymbol{\beta} \cdot \nabla(\Pi^g(v) - v) + \nabla \cdot (\Pi^d(\boldsymbol{\tau}) - \boldsymbol{\tau}))_{\Omega_h} + (\boldsymbol{\sigma}_h, \Pi^d \boldsymbol{\tau} - \boldsymbol{\tau})_{\Omega_h} \\ &\quad + (\boldsymbol{\rho}_h, \nabla(\Pi^g v - v))_{\Omega_h} - \langle \hat{\rho}_h, \gamma(\Pi^g(v) - v) \rangle_{\partial \Omega_h} - \langle \gamma_n(\Pi^d(\boldsymbol{\tau}) - \boldsymbol{\tau}), \hat{u}_h \rangle_{\partial \Omega_h}. \end{aligned}$$

Since  $\boldsymbol{\beta} \in \mathbb{P}_0(\Omega_h; \mathbb{R}^2)$ , it follows by integration by parts that

$$\begin{aligned} (u_h, \boldsymbol{\beta} \cdot \nabla(\Pi^g(v) - v))_{\Omega_h} &= -(\nabla \cdot (\boldsymbol{\beta} u_h), \Pi^g(v) - v)_{\Omega_h} + \langle \gamma_n(\boldsymbol{\beta} u_h), \gamma(\Pi^g(v) - v) \rangle_{\partial \Omega_h} \\ &= 0, \end{aligned}$$

owing to (32a). Proceeding similarly for the other terms, we infer (11a).  $\square$

The nonlinear discrete problem (30) is solved by using the standard Newton method. The initial iterate is defined as the solution of (13) with  $\boldsymbol{\lambda}(|\nabla \mathbf{u}|) = \boldsymbol{\lambda}(0)$  (it is different from 0 in our experiment). Numerical experiments are performed on a uniform refinement of a two-dimensional triangular mesh sequence  $\{\Omega_h\}_h$ , indexed by the level of refinement  $h := \max_{T \in \Omega_h} |T|$ , and satisfying  $\mathcal{O}(h) = \mathcal{O}(\dim(\mathbf{U}_h)^{-\frac{1}{2}})$ . The numerical parameter  $\kappa$  is chosen equal to 1 and the stability parameter  $\alpha := \lambda_{\sharp/b}^2$ . Denoting by  $\mathbf{u} = (u, \boldsymbol{\sigma}, \boldsymbol{\rho}, \hat{u}, \hat{\rho}) \in \mathbf{U}$  and  $\mathbf{u}_h = (u_h, \boldsymbol{\sigma}_h, \boldsymbol{\rho}_h, \hat{u}_h, \hat{\rho}_h) \in \mathbf{U}_h$  the exact and discrete solutions, respectively, we depict the three errors  $\|u - u_h\|_{L^2(\Omega)}$ ,  $\|\boldsymbol{\sigma} - \boldsymbol{\sigma}_h\|_{\mathbf{L}^2(\Omega)}$  and  $\|\boldsymbol{\rho} - \boldsymbol{\rho}_h\|_{\mathbf{L}^2(\Omega)}$ . In addition, from definitions (16) of the trace norms and Lemma 4.1, we depict, instead, the corresponding upper bounds

$$\|\hat{u} - \hat{u}_h\|_{H_0^{1/2}(\partial \Omega_h)} \leq \|u - \mathcal{I}_h^g(\hat{u}_h)\|_{H^1(\Omega)}, \quad \|\hat{\rho} - \hat{\rho}_h\|_{H^{-1/2}(\partial \Omega_h)} \leq \|\boldsymbol{\rho} + \boldsymbol{\beta} u - \mathcal{I}_h^d(\hat{\rho}_h)\|_{\mathbf{H}(\text{div}; \Omega)}.$$

Here,  $\mathcal{I}_h^g : \mathbb{P}_1^0(\partial \Omega_h; \mathbb{R}) \rightarrow \mathbb{P}_1(\Omega_h; \mathbb{R}) \cap C^0(\Omega)$  denotes the  $\mathbb{P}_1$ -Lagrange interpolation operator, and  $\mathcal{I}_h^d : \mathbb{P}_0(\partial \Omega_h; \mathbb{R}^2) \rightarrow \mathbb{RT}_0(\Omega_h)$  is the lowest-order Raviart–Thomas interpolation operator.

**6.2. Example with and without advection.** We consider the nonlinear example without advection from [6] and a corresponding example with non-zero advection. The exact solution is given by  $u(x, y) = \cos(\pi x/2) \cos(\pi y/2)$  on the domain  $(-1, 1)^2$ . The nonlinear diffusive tensor is  $\mathbb{R}$ -valued and defined by  $\lambda(s) = 2 - (1 + s)^{-2}$  for  $s \geq 0$ . Assumptions (14b), (14c) are satisfied with  $\lambda_b = 1$  and  $\lambda_{\sharp} = 3$ . We also consider the case when the advective field  $\boldsymbol{\beta}$  is different from zero, and given by  $\boldsymbol{\beta}(x, y) = (y, -x)$ , satisfying assumption (14a). The approximation errors are depicted in Figure 1.

As expected, all the errors behave like  $\mathcal{O}(h)$ . Comparing the two plots of Figure 1, the presence of the advective field does not impact the accuracy of our method for this test case. In addition, our method delivers very similar result to those obtained with the approach proposed in [6].

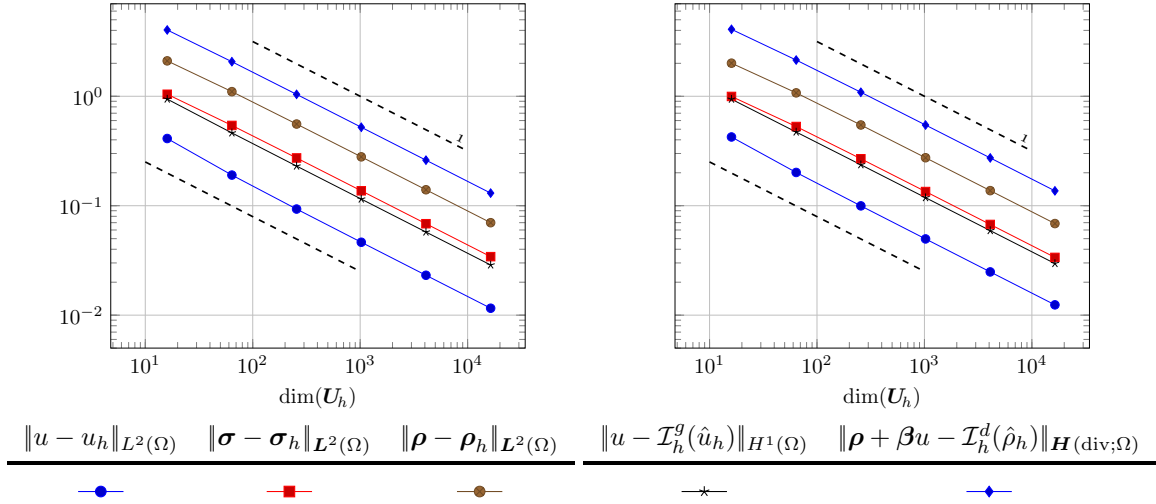


FIG. 1. Left panel: numerical errors with zero advective field. Right panel: numerical error with the advective field  $\beta(x, y) = (y, -x)$ .

- [1] J. BRAMWELL, L. DEMKOWICZ, J. GOPALAKRISHNAN, AND W. QIU, *A locking-free hp DPG method for linear elasticity with symmetric stresses*, *Numerische Mathematik*, 122 (2012), pp. 671–707.
- [2] S. C. BRENNER, *Poincaré–Friedrichs inequalities for piecewise  $H^1$  functions*, *SIAM Journal on Numerical Analysis*, 41 (2003), pp. 306–324.
- [3] D. BROERSEN AND R. STEVENSON, *A robust Petrov-Galerkin discretisation of convection-diffusion equations*, *Computers & Mathematics with Applications*, 68 (2014), pp. 1605–1618.
- [4] D. BROERSEN AND R. STEVENSON, *A Petrov-Galerkin discretization with optimal test space of a mild-weak formulation of convection-diffusion equations in mixed form*, *IMA Journal of Numerical Analysis*, 35 (2015), pp. 39–73.
- [5] T. BUI-THANH AND O. GHATTAS, *A PDE-constrained optimization approach to the discontinuous Petrov-Galerkin method with a trust region inexact Newton-CG solver*, *Computer Methods in Applied Mechanics and Engineering*, 278 (2014), pp. 20–40.
- [6] C. CARSTENSEN, P. BRINGMANN, F. HELLMIG, AND P. WRIGGERS, *Nonlinear discontinuous Petrov–Galerkin methods*, arXiv: 1710.00529, 2017.
- [7] C. CARSTENSEN, L. DEMKOWICZ, AND J. GOPALAKRISHNAN, *Breaking spaces and forms for the DPG method and applications including Maxwell equations*, *Computers & Mathematics with Applications*, 72 (2016), pp. 494–522.
- [8] J. CHAN, L. DEMKOWICZ, AND R. MOSER, *A DPG method for steady viscous compressible flow*, *Computers & Fluids*, 98 (2014), pp. 69–90.
- [9] J. CHAN, N. HEUER, T. BUI-THANH, AND L. DEMKOWICZ, *Robust DPG method for convection-dominated diffusion problems II: Adjoint boundary conditions and mesh-dependent test norms*, *Computers & Mathematics with Applications*, 67 (2014), pp. 771–795.
- [10] P. G. CIARLET, M. H. SCHULTZ, AND R. S. VARGA, *Numerical methods of high-order accuracy for nonlinear boundary value problems*, *Numerische Mathematik*, 13 (1969), pp. 51–77.
- [11] A. COHEN, W. DAHMEN, AND G. WELPER, *Adaptivity and variational stabilization for convection-diffusion equations*, *ESAIM. Mathematical Modelling and Numerical Analysis*, 46 (2012), pp. 1247–1273.
- [12] L. DEMKOWICZ AND J. GOPALAKRISHNAN, *A class of discontinuous Petrov-Galerkin methods. Part I: the transport equation*, *Computer Methods in Applied Mechanics and Engineering*, 199 (2010), pp. 1558–1572.
- [13] L. DEMKOWICZ AND J. GOPALAKRISHNAN, *Analysis of the DPG method for the Poisson problem*, *SIAM Journal on Numerical Analysis*, 49 (2011), pp. 1788–1809.
- [14] L. DEMKOWICZ AND J. GOPALAKRISHNAN, *A class of discontinuous Petrov-Galerkin methods. II. Optimal test functions*, *Numerical Methods for Partial Differential Equations*, 27 (2011), pp. 70–105.
- [15] L. DEMKOWICZ, J. GOPALAKRISHNAN, S. NAGARAJ, AND P. SEPÚLVEDA, *A spacetime DPG method for the Schrödinger equation*, arXiv: 1610.04678, 2016.

- [16] L. DEMKOWICZ AND N. HEUER, *Robust DPG method for convection-dominated diffusion problems*, SIAM Journal on Numerical Analysis, 51 (2013), pp. 2514–2537.
- [17] V. J. ERVIN, T. FÜHRER, N. HEUER, AND M. KARKULIK, *DPG method with optimal test functions for a fractional advection diffusion equation*, Journal of Scientific Computing, 72 (2017), pp. 568–585.
- [18] F. FUENTES, B. KEITH, L. DEMKOWICZ, AND P. LE TALLEC, *Coupled variational formulations of linear elasticity and the DPG methodology*, arXiv: 1609.08160, 2016.
- [19] T. FÜHRER, N. HEUER, AND M. KARKULIK, *On the coupling of DPG and BEM*, Mathematics of Computation, 86 (2017), pp. 2261–2284.
- [20] T. FÜHRER, N. HEUER, M. KARKULIK, AND R. RODRÍGUEZ, *Combining the DPG method with finite elements*, Computational Methods in Applied Mathematics. (in press, DOI:10.1515/cmam-2017-0041).
- [21] T. FÜHRER, N. HEUER, AND E. P. STEPHAN, *On the DPG method for Signorini problems*, IMA Journal of Numerical Analysis. (in press, DOI:10.1093/imanum/drx048).
- [22] V. GIRAULT AND P.-A. RAVIART, *Finite element methods for Navier-Stokes equations*, vol. 5 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 1986. Theory and algorithms.
- [23] J. GOPALAKRISHNAN AND W. QIU, *An analysis of the practical DPG method*, Mathematics of Computation, 83 (2014), pp. 537–552.
- [24] N. HEUER AND M. KARKULIK, *Discontinuous Petrov-Galerkin boundary elements*, Numerische Mathematik, 135 (2017), pp. 1011–1043.
- [25] N. HEUER AND M. KARKULIK, *A robust DPG method for singularly perturbed reaction-diffusion problems*, SIAM Journal on Numerical Analysis, 55 (2017), pp. 1218–1242.
- [26] J.-L. LIONS, *Quelques méthodes de résolution des problèmes aux limites non linéaires*, Études mathématiques, Dunod; Gauthier-Villars, Paris, 1969.
- [27] I. MUGA AND K. G. VAN DER ZEE, *Discretization of linear problems in Banach spaces: Residual minimization, nonlinear Petrov-Galerkin, and monotone mixed methods*, arXiv: 1511.04400, 2015.
- [28] N. V. ROBERTS, T. BUI-THANH, AND L. DEMKOWICZ, *The DPG method for the Stokes problem*, Computers & Mathematics with Applications, 67 (2014), pp. 966–995.
- [29] N. V. ROBERTS, L. DEMKOWICZ, AND R. MOSER, *A discontinuous Petrov-Galerkin methodology for adaptive solutions to the incompressible Navier-Stokes equations*, Journal of Computational Physics, 301 (2015), pp. 456–483.