

# A low-order discontinuous Petrov–Galerkin method for the Stokes equations

Carsten Carstensen<sup>1</sup> · Sophie Puttkammer<sup>1</sup>

Received: 4 April 2016 / Revised: 27 February 2018 / Published online: 9 April 2018  
© Springer-Verlag GmbH Germany, part of Springer Nature 2018

**Abstract** This paper introduces a low-order discontinuous Petrov-Galerkin (dPG) finite element method (FEM) for the Stokes equations. The ultra-weak formulation utilizes piecewise constant and affine ansatz functions and piecewise affine and discontinuous lowest-order Raviart–Thomas test search functions. This low-order discretization for the Stokes equations allows for a direct proof of the discrete inf-sup condition with explicit constants. The general framework of Carstensen et al. (SIAM J Numer Anal 52(3):1335–1353, 2014) then implies a complete a priori and a posteriori error analysis of the dPG FEM in the natural norms. Numerical experiments investigate the performance of the method and underline its quasi-optimal convergence.

**Keywords** Stokes · Discontinuous Petrov Galerkin · Low-order discretization · A priori · A posteriori · Adaptive mesh refinement

**Mathematics Subject Classification** 65N12 · 65N15 · 65N30 · 65Y05 · 65Y20

## 1 Introduction

The initial motivation for the discontinuous Petrov-Galerkin (dPG) methodology in [18–20] was the design of the optimal test function space in applications of fluid mechanics, when a stabilization appears obligatory for many standard finite element methods (FEMs). Since important examples of this class follow as linearizations of

---

✉ Carsten Carstensen  
cc@math.hu-berlin.de

Sophie Puttkammer  
puttkams@math.hu-berlin.de

<sup>1</sup> Department of Mathematics, Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany

the Navier–Stokes equations, the understanding of simple low-order dPG schemes for the Stokes equations appears to be a necessary step. The first dPG FEMs for the Stokes equations in [25] and [10] utilize polynomials of much higher degrees in the trial and test search space. This paper introduces a much simpler lowest-order dPG FEM with an emphasis on a direct estimation of the discrete inf-sup constant and the discussion of the associated Fortin-type operators for a reliable a posteriori analysis to generalize [10] for low-order test functions. Popular alternative simulation tools for the Stokes equations with optimal convergence rates for adaptive mesh-refining algorithm are the nonconforming (restricted to first-order in  $3D$ ) and the pseudostress FEM with a more complicated a posteriori error analysis [8, 16].

The dGP methodology is roughly described as a minimum residual method with discontinuous ansatz and test functions. This leads to piecewise (also called broken) Sobolev spaces with related trace spaces on element boundaries and so requires a careful definition and analysis on the independence of the underlying partition. In return, this results in a local and parallel computation of the underlying dual norms and a simple implementation and allows rather general geometries of the element domains; both regarded as obligatory in particular in higher space dimensions. The detailed description of the ultraweak formulation with piecewise smooth functions and several flux variables on the boundaries of the element domain is cumbersome and follows in Sect. 3. For the sake of this introduction it may suffice to acknowledge that this leads to a continuous formulation in Lebesgue and broken Sobolev spaces  $X$  and  $Y$  such that the continuous problem of the standard Stokes equation leads to a right-hand side  $F \in Y^*$  and an exact solution  $x \in X$  of the (well-posed) equation

$$b(x, y) = F(y) \quad \text{for all } y \in Y. \quad (1)$$

The bounded bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  models the equivalent ultraweak formulation of the Stokes equation as in [10]. This equation is well posed if that  $b$  satisfies an inf-sup condition on the continuous level,

$$0 < \beta := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y}. \quad (2)$$

With a few and well-spotted exceptions, the least-squares FEMs start with the minimization of the residual  $F - b(x_h, \bullet)$  in subspaces of  $L^2$ . For the bilinear form  $b$  at hand, this is impossible as  $Y$  does not solely contain Lebesgue functions. The dPG schemes first approximate the dual norm in  $Y^*$  of the residual by the dual norm  $Y_h^*$  over a finite-dimensional subspace  $Y_h \subset Y$  of  $Y$  and second minimize the residual  $F - b(x_h, \bullet)$  for ansatz functions  $x_h$  in a finite-dimensional subspace  $X_h \subset X$  of  $X$ . In other words, the dPG approximation is the minimizer  $x_h$  in

$$x_h = \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*} = \min_{\xi_h \in X_h} \max_{y_h \in Y_h \setminus \{0\}} (F - b(\xi_h, y_h)) / \|y_h\|_Y. \quad (3)$$

The computational costs are related to the total number  $N + M$  of unknowns with the dimensions  $N := \dim(X_h)$  for the ansatz function space and  $M := \dim(Y_h)$  for

the test search space. The minimal residual method is *not* a mixed finite element scheme in that it allows  $N \leq M$  with significantly larger  $M$ . The benefit is that the *test search space*  $Y_h$  can be much richer and approximate  $Y$  well so that the crucial discrete inf-sup condition

$$0 < \beta_h := \inf_{x_h \in X_h \setminus \{0\}} \sup_{y_h \in Y_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y} \tag{4}$$

can be made larger to approximate the idealized inf-sup constant  $\beta_h^*$ ,

$$\beta_h \leq \beta_h^* := \inf_{x_h \in X_h \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{b(x_h, y)}{\|x_h\|_X \|y\|_Y},$$

which may even be larger than the global inf-sup constant  $\beta$ . Hence a sufficiently large test search space  $Y_h$  may stabilize a situation, when a stable pairing does not exist or is at least unknown and a mixed finite element scheme is not available with  $M = N$ . It is known that the dPG scheme is equivalent to a mixed scheme with an extended bilinear form  $\mathcal{B} : (X \times Y) \times (X \times Y) \rightarrow \mathbb{R}$ , when  $Y$  is a Hilbert space [from L. Demkovicz in personal communication]. Moreover, it can even be reduced to the computation of some subspace  $M_h \subset Y_h$  with  $\dim(M_h) = N$  such that  $x_h$  is a solution to a quadratic mixed FEM with  $b$  reduced to  $X_h \times Y_h$  [20]. This is all related to the numerical linear algebra of the dPG schemes and the computational costs grow with  $N + M$ . It is therefore practically relevant to minimize the test search space  $Y_h$  and so  $M \geq N$ , while  $\beta_h > 0$  is still uniformly bounded away from zero as the underlying partitions become finer and finer. The first proofs of a stability result of this type [22] involve some linear and bounded Fortin operator  $\Pi : Y \rightarrow Y_h$  with operator norm  $\|\Pi\|$  and the annulation property

$$b(x_h, y - \Pi y) = 0 \quad \text{for all } x_h \in X_h \quad \text{and} \quad y \in Y. \tag{5}$$

Given such an operator  $\Pi$ , the analysis in [22] leads to  $\beta/\|\Pi\| \leq \beta_h$  [6, Proposition 5.4.2] and so is a sufficient condition for stability. Conversely, the stability leads to the existence of some Fortin interpolation operator  $\Pi$  with  $\|\Pi\| \leq \|b\|/\beta_h$  [14, Lemma 2.10].

The examples in [10,22] typically involve piecewise polynomials of degree  $k$  (and one variable with  $k + 1$ ) in  $X_h$  and piecewise polynomials of degree  $k + n$  in  $Y_h$  of the underlying partition with  $J$  element domains in  $\mathbb{R}^n$  with  $n$  space dimensions. This leads to  $N = \mathcal{O}(J(k + 1))$  and  $M = \mathcal{O}(J(k + n + 1))$ , which results in overall computational costs which grow with  $J(2k + n + 2)$ . The subsequent discussion concerns the same fixed ansatz space  $X_h$  and so  $N$  is fixed. The overall costs are then expected to be of a monoton function in  $M$  and the precise dependence is less clear for an optimized numerical linear algebra with parallel computation. This paper is motivated in the extreme case  $k = 0$  because then the current dPG schemes require  $M = \mathcal{O}(J(1 + n))$  which is  $n + 1$  times higher than the costs for a (unknown) mixed FEM with  $M = N = \mathcal{O}(J)$  for the space dimension  $n = 2, 3$ . This paper introduces a stable choice of  $Y_h$  with piecewise polynomial degree at most 1 rather than  $n$  from [10,22] for  $k = 0$ .

The mentioned ultraweak formulation of the well-known Stokes equation with a volume term  $f \in L^2(\Omega; \mathbb{R}^n)$  on the right-hand side leads on the discrete level to the residual  $F(y_h) - b(x_h, y_h)$ . Throughout this paper,  $\text{dev}$  denotes the deviatoric part of a matrix,  $D_{\text{NC}}$  is the piecewise functional matrix,  $\cdot$  (resp.  $:$ ) denotes the scalar product of two vectors (resp. matrices), cf. Sect. 2 for more details. For some particular  $x_h$  and  $y_h$  in the Stokes equations below, the aforementioned residual reads as

$$\begin{aligned} & \int_{\Omega} f \cdot v_1 \, dx - \int_{\Omega} \sigma_0 : (D_{\text{NC}} v_1 + \text{dev } \tau_{\text{RT}}) \, dx \\ & - \int_{\Omega} u_0 \cdot \text{div}_{\text{NC}} \tau_{\text{RT}} \, dx + \sum_{T \in \mathcal{T}} \int_{\partial T} (t_0 \cdot v_1 + s_1 \cdot \tau_{\text{RT}} \nu) \, ds \end{aligned}$$

up to modifications for the Dirichlet boundary conditions. Therein,  $\sigma_0$  and  $u_0$  are piecewise constant functions, while  $v_1$  and  $\tau_{\text{RT}}$  are piecewise affine with respect to a triangulation  $\mathcal{T}$ . On the skeleton with respect to the sides  $\mathcal{E}$  in  $\mathcal{T}$ ,  $t_0$  is piecewise constant but,  $s_1$  is piecewise affine and globally continuous.

This paper bounds the inf-sup constants (2) and (4) for arbitrary dimension  $n$  explicitly in terms of the Friedrichs, the tr-div-dev constant, and the inf-sup constant of the  $H(\text{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n)$  mixed FEM for Stokes equations. This implies the quasi optimal convergence

$$\|x - x_h\|_X \leq \frac{\|\Pi\| \|b\|}{\beta} \min_{\xi_h \in X_h} \|x - \xi_h\|_X \quad (6)$$

for the novel low-order dPG FEM. The general a posteriori error analysis of [10] leads to the a posteriori error control for any approximation  $\xi_h \in X_h$  (so it allows an inexact solve of the discrete minimization problem)

$$\begin{aligned} \beta \|x - \xi_h\|_X & \leq \|\Pi\| \left\| F - b(\xi_h, \bullet) \right\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \\ & \leq \|b\| (\|\Pi\| + \|1 - \Pi\|) \|x - \xi_h\|_X. \end{aligned} \quad (7)$$

The residual term  $\left\| F - b(\xi_h, \bullet) \right\|_{Y_h^*}$  is computable and the remaining data approximation term  $\|F \circ (1 - \Pi)\|_{Y^*}$  involves the Fortin interpolation  $\Pi$ . In all the examples of [10] with the aforementioned larger test search spaces, this term is an oscillation and hence the data approximation term may be regarded as a higher-order term and in fact is neglected in many practical calculations.

In the novel low-order dPG scheme, this is not the case and the Fortin interpolation operator is characterized in Theorem 5.2 below. It turns out that the data approximation term is of first order and so, for quasi-uniform meshes and a singular solution possibly of higher-order. For adaptive mesh-refining, this argument is no longer valid and the a posteriori error control may fail to be efficient. This leads to the extension  $\hat{Y}_h$  of the trial search space  $Y_h \subset \hat{Y}_h \subset Y$  by three piecewise enrichments by additional cubic bubble functions, piecewise affines or first-order

Raviart-Thomas functions in the first component. The resulting overall strategy for guaranteed and effective a posteriori error control assumes an approximation  $x_h \in X_h$  computed by the proposed dPG scheme even with inexact solve from an iterative numerical linear algebra with the test search space  $Y_h$ . The a posteriori error control applies to  $\hat{Y}_h$  and computes the residual  $\|F - b(x_h, \bullet)\|_{\hat{Y}_h^*}$  in the dual norm  $\hat{Y}_h^*$  and allows for the reduced data approximation term  $\|F \circ (1 - \hat{T})\|_{Y^*}$  with respect to the Fortin interpolation  $\hat{T} : Y \rightarrow \hat{Y}_h$ . Notice that the inf-sup constant  $\hat{\beta}_h \geq \beta_h$ , where  $\hat{\beta}_h = \inf_{x_h \in X_h} \sup_{\hat{y}_h \in \hat{Y}_h} b(x_h, \hat{y}_h) / (\|x_h\|_X \|\hat{y}_h\|_Y)$ . This leads to  $\|F \circ (1 - \hat{T})\|_{Y^*}$  as oscillations, which may be negligible at least for piecewise smooth data. The analysis of this strategy and affirmative numerical examples conclude the paper.

The remaining parts of the paper are organized as follows. Section 2 recalls the necessary notation on triangulation and function spaces. Section 3 and 4 investigate the continuous and discontinuous formulation (1)–(3) related to (7) and prove the inf-sup conditions (2) and (4). Section 5 discusses the data approximation error in the two-dimensional case which contains the Fortin interpolator. Numerical experiments for benchmark problems are presented in Sect. 6. The supplement contains some remarks on the Fortin operator and on the implementation.

Standard notation applies to Lebesgue and Sobolev spaces throughout this paper,  $H^1(T)$  abbreviates  $H^1(\text{int}(T))$  for a set  $T$  with nonempty interior  $\text{int}(T)$ . Furthermore,  $a \lesssim b$  abbreviates, that there exists a generic constant  $C$  with  $a \leq Cb$ , while  $a \approx b$  abbreviates  $a \lesssim b \lesssim a$ . Given a normed linear space  $(X, \|\bullet\|_X)$ , let  $S(X) := \{x \in X : \|\bullet\|_X = 1\}$  be its unit sphere.

## 2 Notation

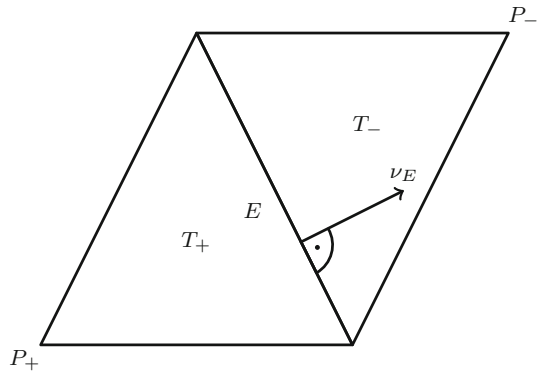
### 2.1 Vector and matrix notation

This subsection clarifies details on the overall notation of vectors and matrices. For two vectors  $a, b \in \mathbb{R}^m$ , the dot denotes the scalar product  $a \cdot b = \sum_{j=1}^m a_j b_j \in \mathbb{R}$ , while the scalar product  $A : B$  of  $m \times m$  matrices  $A, B \in \mathbb{R}^{m \times m}$  reads  $A : B = \sum_{j,k=1}^m A_{jk} B_{jk} \in \mathbb{R}$ . The dyadic product of  $a, b \in \mathbb{R}^m$  reads  $a \otimes b := ab^T \in \mathbb{R}^{m \times m}$ . Notice that  $|a \otimes b| = |a||b|$ . The identity mapping is denoted by  $\bullet$ . The notation  $|\bullet|$  is dependent on context, the norm induced by  $\cdot$  (resp.  $:$ ) on  $\mathbb{R}^n$  (resp.  $\mathbb{R}^{n \times n}$ ), the cardinality of a finite set, the  $n$ - or  $(n - 1)$ -dimensional Lebesgue measure of a subspace of  $\mathbb{R}^n$ . The linear operators deviator,  $\text{dev } A = A - 1/n (\text{tr } A) I_{n \times n}$ , and trace,  $\text{tr } A = A_{11} + \dots + A_{nn}$ , of any matrix  $A \in \mathbb{R}^{n \times n}$ , lead to  $\text{tr dev } A = 0$  and

$$\|\tau\|_{L^2(\Omega)}^2 = 1/n \|\text{tr } \tau\|_{L^2(\Omega)}^2 + \|\text{dev } \tau\|_{L^2(\Omega)}^2 \quad \text{for all } \tau \in L^2(\Omega; \mathbb{R}^{n \times n}). \quad (8)$$

(This is the theorem of Pythagoras  $|A|^2 = A : A = |\text{dev } A|^2 + 1/n(\text{tr } A)^2$  for a matrix  $A \in \mathbb{R}^{n \times n}$  based on the orthogonality of the unit matrix  $I_{n \times n}$  and the deviatoric part  $\text{dev}$ .) Let  $\mathbb{R}_{\text{dev}}^{n \times n} := \text{dev}(\mathbb{R}^{n \times n})$  denote the deviatoric (also called trace-free)  $n \times n$  matrices and note  $\text{dev } A : \text{dev } B = \text{dev } A : B = A : \text{dev } B$  for all  $A, B \in \mathbb{R}^{n \times n}$ .

**Fig. 1** Edge patch  
 $\omega_E := T_+ \cup T_-$  for an edge  
 $E \in \mathcal{E}(\Omega)$



## 2.2 Triangulation

Given a regular triangulation  $\mathcal{T}$  of  $\Omega \subseteq \mathbb{R}^n$  into closed  $n$ -simplices  $T \in \mathcal{T}$ ,  $\mathcal{E}(T)$  denotes the set of all  $n + 1$  sides ( $(n - 1)$ -simplices like edges for  $n = 2$  and faces for  $n = 3$ ) of  $T$  and  $\mathcal{N}(T)$  the set of all  $n + 1$  vertices of  $T$ . The set of all sides and nodes read

$$\mathcal{E} := \bigcup_{T \in \mathcal{T}} \mathcal{E}(T) \quad \text{and} \quad \mathcal{N} := \bigcup_{T \in \mathcal{T}} \mathcal{N}(T);$$

the set of all interior (resp. boundary) sides reads  $\mathcal{E}(\Omega)$  (resp.  $\mathcal{E}(\partial\Omega)$ ) as well as  $\mathcal{N}(\Omega)$  (resp.  $\mathcal{N}(\partial\Omega)$ ) is the set of all interior (resp. boundary) nodes. The skeleton  $\partial\mathcal{T} := \bigcup_{T \in \mathcal{T}} \partial T$  is the union of all boundaries of simplices  $T \in \mathcal{T}$ . Throughout this paper,  $h_{\mathcal{T}}$  abbreviates the piecewise constant function with  $h_{\mathcal{T}}|_T := h_T := \text{diam}(T) = \max_{x, y \in T} |x - y|$  the diameter of a simplex  $T \in \mathcal{T}$  and  $h_{\max} := \max h_{\mathcal{T}}$  its maximum.

Let  $\nu_T$  denote the outer unit normal vector field along the boundary  $\partial T$  on a fixed element  $T \in \mathcal{T}$ . Each side  $E \in \mathcal{E}$  has an assigned orientation of the unit normal  $\nu_E$ . For exterior sides  $E \in \mathcal{E}(\partial\Omega)$ ,  $\nu_E = \nu_{\Omega}$  points outwards. For an interior side  $E = \partial T_+ \cap \partial T_- \in \mathcal{E}(\Omega)$  one orientation of the unit normal  $\nu_E$  is fixed throughout this paper. The neighbouring triangles are named such that  $\nu_E$  points from  $T_+$  to  $T_-$  as in Fig. 1. In this context the following sign-function is defined  $\text{sgn}(T, E) := \nu_E \cdot \nu_T \in \{\pm 1\}$  for all  $T \in \mathcal{T}$ ,  $E \in \mathcal{E}(T)$ . Furthermore, for a function  $v \in L^2(\Omega; \mathbb{R}^{m \times n})$  the jump along an interior side  $E \in \mathcal{E}(\Omega)$  is denoted by  $[v]_E := (v|_{T_+} - v|_{T_-})|_E \in L^2(E; \mathbb{R}^{m \times n})$  and along an boundary side  $E \in \mathcal{E}(\partial\Omega)$  by  $[v]_E := v|_E \in L^2(E; \mathbb{R}^{m \times n})$ .

For each simplex  $T \in \mathcal{T}$ ,  $\text{mid}(T) := \int_T x \, dx = |T|^{-1} \int_T x \, dx = 1/(n + 1) \sum_{z \in \mathcal{N}(T)} z$  denotes the center of gravity and the function  $\bullet - \text{mid}(T) \in L^\infty(\Omega; \mathbb{R}^n)$  has the value  $x - \text{mid}(T)$  for  $x \in T \in \mathcal{T}$  and satisfies for all  $T \in \mathcal{T}$

$$\int_T x - \text{mid}(T) \, dx = 0 \quad \text{and} \quad \|\bullet - \text{mid}(T)\|_{L^\infty(\Omega)} \leq h_{\max} n/(n + 1). \quad (9)$$

### 2.3 Function spaces

Standard notation applies to  $L^2(\Omega)$ ,  $H^1(\Omega)$ ,  $H(\operatorname{div}, \Omega)$  and their vector- or matrix-valued relatives such as  $L^2(\Omega; \mathbb{R}^n)$ ,  $L^2(\Omega; \mathbb{R}^{n \times n})$ ,  $H^1(\Omega; \mathbb{R}^n)$ ,  $H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ . Let  $\mathcal{T}$  be a regular triangulation of  $\Omega$ . The test search space only exhibits certain piecewise regularity properties on  $T \in \mathcal{T}$ ,

$$H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) := \left\{ \boldsymbol{\tau} \in L^2(\Omega; \mathbb{R}^{n \times n}) : \forall T \in \mathcal{T}, 1 \leq j \leq n, \tau_j|_T \in H(\operatorname{div}, T) \right\},$$

$$H^1(\mathcal{T}; \mathbb{R}^n) := \left\{ v \in L^2(\Omega; \mathbb{R}^n) : \forall T \in \mathcal{T}, v|_T \in H^1(T; \mathbb{R}^n) \right\},$$

where  $\tau_j$  denotes the  $j$ -th row of  $\boldsymbol{\tau}$ . The piecewise application of the divergence operator  $\operatorname{div}$  and the derivative  $\mathbf{D}$  read  $\operatorname{div}_{\text{nc}}$  and  $\mathbf{D}_{\text{nc}}$  and give rise to

$$\|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T})}^2 := \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})}^2 := \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\operatorname{div}_{\text{nc}} \boldsymbol{\tau}\|_{L^2(\Omega)}^2,$$

$$\|v\|_{H^1(\mathcal{T})}^2 := \|v\|_{H^1(\mathcal{T}; \mathbb{R}^n)}^2 := \|v\|_{L^2(\Omega)}^2 + \|\mathbf{D}_{\text{nc}} v\|_{L^2(\Omega)}^2.$$

The following essential facts about trace spaces are proven in [2, 21]. For any open, bounded Lipschitz domain  $U \subseteq \mathbb{R}^n$ , there exists exactly one continuous linear mapping  $\gamma_0 : H^1(U) \rightarrow L^2(\partial U)$  with  $\gamma_0 w = w|_{\partial U}$  for all  $w \in H^1(U) \cap C^0(\bar{U})$ . Let  $H^{1/2}(\partial U) := \gamma_0(H^1(U))$  and let  $H^{-1/2}(\partial U) = (H^{1/2}(\partial U))^*$  be its dual space. Then there exists exactly one continuous linear mapping  $\gamma_\nu : H(\operatorname{div}, U) \rightarrow H^{-1/2}(\partial U)$  with  $\gamma_\nu q = (q|_{\partial U}) \cdot \nu$  for all  $q \in H(\operatorname{div}, U)$ . Moreover, for all  $q \in H(\operatorname{div}, U)$  and  $w \in H^1(U)$  it holds

$$\langle \gamma_\nu q, \gamma_0 w \rangle_{\partial U} = \int_U q \cdot \mathbf{D} w \, dx + \int_U w \operatorname{div} q \, dx. \tag{10}$$

The extension of the  $L^2$ -scalar product on the skeleton is for all  $t = (t_T)_{T \in \mathcal{T}} \in \prod_{T \in \mathcal{T}} H^{-1/2}(\partial T; \mathbb{R}^n)$  and  $s = (s_T)_{T \in \mathcal{T}} \in \prod_{T \in \mathcal{T}} H^{1/2}(\partial T; \mathbb{R}^n)$  denoted by

$$\langle t, s \rangle_{\partial \mathcal{T}} := \sum_{T \in \mathcal{T}} \langle t_T, s_T \rangle_{\partial T}.$$

Define the trace operators

$$\gamma_0^{\mathcal{T}} : H^1(\mathcal{T}; \mathbb{R}^n) \rightarrow \prod_{T \in \mathcal{T}} H^{1/2}(\partial T; \mathbb{R}^n),$$

$$\gamma_\nu^{\mathcal{T}} : H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) \rightarrow \prod_{T \in \mathcal{T}} H^{-1/2}(\partial T; \mathbb{R}^n)$$

on the skeleton  $\partial \mathcal{T}$  by  $\gamma_0^{\mathcal{T}} w := (s_T)_{T \in \mathcal{T}}$  with  $s_T := \gamma_0(w|_T)$  and  $\gamma_\nu^{\mathcal{T}} \mathbf{q} := (t_T)_{T \in \mathcal{T}}$  with  $t_T := \gamma_\nu(\mathbf{q}|_T)$  for all  $T \in \mathcal{T}$ . The associated trace spaces read

$$H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n) := \gamma_0^{\mathcal{T}} \left( H_0^1(\Omega; \mathbb{R}^n) \right), \tag{11}$$

$$H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n) := \gamma_v^{\mathcal{T}} \left( H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) \right). \quad (12)$$

These spaces are equipped with the following minimal extension norms

$$\begin{aligned} \|s\|_{H_0^{1/2}(\partial\mathcal{T})} &:= \|s\|_{H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)} := \inf_{\substack{w \in H_0^1(\Omega; \mathbb{R}^n) \\ \gamma_0^{\mathcal{T}} w = s}} \|w\|_{H^1(\Omega)}, \\ \|t\|_{H^{-1/2}(\partial\mathcal{T})} &:= \|t\|_{H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)} := \inf_{\substack{\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) \\ \gamma_v^{\mathcal{T}} \mathbf{q} = t}} \|\mathbf{q}\|_{H(\operatorname{div}, \Omega)}. \end{aligned}$$

The spaces  $H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$  and  $H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$  are subspaces of product spaces and not dual to each other in general.

**Lemma 2.1** (Duality Lemma) *It holds*

$$\begin{aligned} \|s\|_{H_0^{1/2}(\partial\mathcal{T})} &= \sup_{\boldsymbol{\tau} \in S(H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R})} \left\langle \gamma_v^{\mathcal{T}} \boldsymbol{\tau}, s \right\rangle_{\partial\mathcal{T}} \quad \text{for all } s \in H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n), \\ \|t\|_{H^{-1/2}(\partial\mathcal{T})} &= \sup_{v \in S(H^1(\mathcal{T}; \mathbb{R}^n))} \left\langle t, \gamma_0^{\mathcal{T}} v \right\rangle_{\partial\mathcal{T}} \quad \text{for all } t \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n). \end{aligned}$$

*Proof* This is contained in [11, Lemma 2.2]. □

## 2.4 Discrete function spaces

The finite-dimensional subspaces of the trial space  $X_h \subset X$  and the test search space  $Y_h \subset Y$  are piecewise polynomials. For any  $k \in \mathbb{N}_0$ , let  $P_k(\mathcal{T}; \mathbb{R}^{m \times n})$  denote polynomials of total degree at most  $k$  in each component as functions in  $L^2(\mathcal{T}; \mathbb{R}^{m \times n})$  and set

$$P_k(\mathcal{T}, \mathbb{R}^{m \times n}) := \left\{ \mathbf{q}_k \in L^\infty(\Omega; \mathbb{R}^{m \times n}) : \forall T \in \mathcal{T}, \mathbf{q}_k|_T \in P_k(\mathcal{T}; \mathbb{R}^{m \times n}) \right\}.$$

Analogous definitions apply on the skeleton, i.e.,

$$\begin{aligned} P_k(\mathcal{E}; \mathbb{R}^{m \times n}) &:= \left\{ \mathbf{q}_k \in L^\infty\left(\bigcup \mathcal{E}; \mathbb{R}^{m \times n}\right) : \forall T \in \mathcal{T}, \forall E \in \mathcal{E}(T), \right. \\ &\quad \left. \mathbf{q}_k|_E \in P_k(E; \mathbb{R}^{m \times n}) \right\}. \end{aligned}$$

Let  $\Pi_0$  be the  $L^2$  projection onto  $P_0(\mathcal{T})$  defined for  $f \in L^2(\Omega; \mathbb{R}^{m \times n})$  by  $\Pi_0 f|_T := |T|^{-1} \int_T f \, dx = \bar{f}_T$ . The continuous and piecewise finite element functions  $P_k$  on  $\mathcal{T}$  read

$$\begin{aligned} S^k(\mathcal{T}; \mathbb{R}^{m \times n}) &:= P_k(\mathcal{T}; \mathbb{R}^{m \times n}) \cap C(\bar{\Omega}), \\ S_0^k(\mathcal{T}; \mathbb{R}^{m \times n}) &:= S^k(\mathcal{T}; \mathbb{R}^{m \times n}) \cap C_0(\Omega) \end{aligned}$$



and on the skeleton

$$S^k(\mathcal{E}; \mathbb{R}^{m \times n}) := P_k(\mathcal{E}; \mathbb{R}^{m \times n}) \cap C\left(\bigcup_{T \in \mathcal{T}} \partial T\right),$$

$$S_0^k(\mathcal{E}; \mathbb{R}^{m \times n}) := \left\{ v \in S^k(\mathcal{E}; \mathbb{R}^{m \times n}) : v|_{\partial\Omega} \equiv 0 \right\}.$$

The lowest-order Raviart-Thomas functions read

$$RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) := \left\{ \mathbf{q}_{\text{RT}} \in L^\infty(\Omega; \mathbb{R}^{n \times n}) : \exists A \in P_0(\mathcal{T}; \mathbb{R}^{n \times n}), \right.$$

$$\left. \exists b \in P_0(\mathcal{T}; \mathbb{R}^n), \mathbf{q}_{\text{RT}} = A + b \otimes (\bullet - \text{mid}(T)) \right\},$$

$$RT_0(\mathcal{T}; \mathbb{R}^{n \times n}) := RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) \cap H(\text{div}, \Omega; \mathbb{R}^{n \times n}).$$

On each simplex  $T \in \mathcal{T}$  any  $\mathbf{q}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n})$  can be written as  $\mathbf{q}_{\text{RT}}|_T = A + 1/n \text{ div } \mathbf{q}_{\text{RT}} \otimes (\bullet - \text{mid}(T))$  for some  $A \in \mathbb{R}^{n \times n}$ . Then it holds by (9)

$$(1 - \Pi_0) \mathbf{q}_{\text{RT}} = 1/n \text{ div } \mathbf{q}_{\text{RT}} \otimes (\bullet - \text{mid}(T)) \perp P_0(\mathcal{T}; \mathbb{R}^{n \times n}). \tag{13}$$

It is useful to regard  $P_0(\mathcal{E}; \mathbb{R}^n)$  as a subspace of  $H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$  via the embedding

$$P_0(\mathcal{E}; \mathbb{R}^n) \hookrightarrow H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n), \quad t_0 \mapsto t = (t_T)_{T \in \mathcal{T}} \text{ with } t_T = \mathbf{q}_{\text{RT}} \nu_T|_{\partial T},$$

where  $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$  satisfies  $\mathbf{q}_{\text{RT}}|_E \nu_E = t_0|_E$  for all  $E \in \mathcal{E}$ . Notice the norm equivalence

$$\|t_0\|_{H^{-1/2}(\Omega)} \leq \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)} \leq (1 + \sqrt{1 + 4h_{\max}^2/\pi^2}) \|t_0\|_{H^{-1/2}(\Omega)}$$

from [13, Lemma 3.2].

### 3 Continuous problem

Given some  $f \in L^2(\Omega; \mathbb{R}^n)$  on some  $n$ -dimensional, bounded Lipschitz domain  $\Omega$  with polyhedral boundary  $\partial\Omega$  and Dirichlet boundary data  $g \in H^1(\partial\Omega; \mathbb{R}^n)$  with  $\int_{\partial\Omega} g \cdot \nu \, ds = 0$ , the Stokes pseudostress formulation seeks  $\sigma \in H(\text{div}, \Omega; \mathbb{R}^{n \times n})$  and  $u \in H^1(\Omega; \mathbb{R}^n)$  with

$$\text{dev } \sigma = \text{D}u, \quad f + \text{div } \sigma = 0 \text{ in } \Omega, \quad u = g \text{ along } \partial\Omega. \tag{14}$$

There exists a unique solution  $(\sigma, u)$  to (14) up to a constant multiple of the  $n \times n$  unit matrix  $I_{n \times n}$  fixed by  $\int_{\Omega} \text{tr } \sigma \, dx = 0$  written  $\sigma \in H(\text{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ . The discontinuous Petrov-Galerkin formulation (dPG) is based on a regular triangulation  $\mathcal{T}$  of  $\Omega$  from Sect. 2.2. On each simplex  $T \in \mathcal{T}$ , a multiplication of (14) with the test

functions  $\boldsymbol{\tau} \in H(\operatorname{div}, T; \mathbb{R}^{n \times n})$  and  $v \in H^1(T; \mathbb{R}^n)$  followed by an integration by parts leads to

$$\int_T \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_T u \cdot \operatorname{div} \boldsymbol{\tau} \, dx = \langle \gamma_\nu \boldsymbol{\tau}, \gamma_0 u \rangle_{\partial T},$$

$$\int_T \boldsymbol{\sigma} : \mathbf{D} v \, dx - \langle \gamma_\nu \boldsymbol{\sigma}, \gamma_0 v \rangle_{\partial T} = \int_T f \cdot v \, dx.$$

The summation over all  $T \in \mathcal{T}$  results in traces on the skeleton  $\gamma_0^T u$  and  $\gamma_\nu^T \boldsymbol{\sigma}$ . Let  $g \in H^1(\Omega; \mathbb{R}^n)$  extend the Dirichlet boundary data  $g \in H^1(\partial\Omega; \mathbb{R}^n)$ . The interface variables  $s := \gamma_0^T (u - g)$  and  $t := \gamma_\nu^T \boldsymbol{\sigma}$  circumvent the continuity conditions for  $\boldsymbol{\sigma}$  and  $u$ . The sum of the two equations leads to the dPG formulation (on the continuous level). In abstract notation, the dPG formulation seeks  $x \in X$  with

$$b(x, y) = F(y) \quad \text{for all } y \in Y. \tag{15}$$

For any  $x = (\boldsymbol{\sigma}, u, s, t) \in X$  and  $y = (\boldsymbol{\tau}, v) \in Y$  with

$$X := L^2(\Omega; \mathbb{R}^{n \times n})/\mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n) \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n), \tag{16}$$

$$Y := H(\operatorname{div}, T; \mathbb{R}^{n \times n})/\mathbb{R} \times H^1(\mathcal{T}; \mathbb{R}^n), \tag{17}$$

the bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  and the functional  $F \in Y^*$  read

$$b(x, y) := \int_\Omega \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx + \int_\Omega \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_\Omega u \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \, dx$$

$$- \langle t, \gamma_0^T v \rangle_{\partial\mathcal{T}} - \langle \gamma_\nu^T \boldsymbol{\tau}, s \rangle_{\partial\mathcal{T}}, \tag{18}$$

$$F(y) := \int_\Omega f \cdot v \, dx + \langle \gamma_\nu^T \boldsymbol{\tau}, \gamma_0^T g \rangle_{\partial\mathcal{T}}. \tag{19}$$

The remaining parts of this section establish the boundedness of  $b$ , its non-degeneracy, and the inf-sup condition (2). The weak formulation of (14) leads with  $Z := H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R} \times L^2(\Omega; \mathbb{R}^n)$  to a bilinear form  $\tilde{b} : Z \times Z \rightarrow \mathbb{R}$  defined for  $(\boldsymbol{\tau}, v), (\boldsymbol{\rho}, w) \in Z$  by

$$\tilde{b}((\boldsymbol{\tau}, v), (\boldsymbol{\rho}, w)) := \int_\Omega \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\rho} \, dx + \int_\Omega v \cdot \operatorname{div} \boldsymbol{\rho} \, dx + \int_\Omega \operatorname{div} \boldsymbol{\tau} \cdot w \, dx. \tag{20}$$

The well-posedness of (14) leads to a positive inf-sup constant [9, Thm.2.3]

$$0 < \gamma := \inf_{a \in S(Z)} \sup_{b \in S(Z)} \tilde{b}(a, b), \tag{21}$$

which allows to describe the dependence of the inf-sup constant  $\beta$  below. The bilinear form  $b$  of the ultraweak formulation is a broken form of the established bilinear form  $\tilde{b}$  with the term broken used in the sense of [11].

**Theorem 3.1** *The bilinear form  $b$  from (18) is bounded with*

$$|b(x, y)| \leq \sqrt{3} \|x\|_X \|y\|_Y \text{ for all } x \in X, y \in Y$$

and satisfies  $N = \{y \in Y : b(\bullet, y) = 0 \in X^*\} = \{0\}$  as well as

$$0 < 1/\sqrt{15/\gamma^2 + 6 + \sqrt{32 + 168/\gamma^2 + 225/\gamma^4}} \leq \beta := \inf_{x \in S(X)} \sup_{y \in S(Y)} b(x, y).$$

The constant  $\gamma$  involves the Ladyshenskaya constant [7, (11.2.3)] or the constant  $C_{\text{tdd}}$  from the following tr-dev-div lemma.

**Lemma 3.2** [6, Thm.9.1.1] *There exists a constant  $C_{\text{tdd}} < \infty$  (solely depending on  $\Omega$ ) such that any  $\tau \in H(\text{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$  satisfies*

$$\|\text{tr } \tau\|_{L^2(\Omega)} \leq C_{\text{tdd}} \left( \|\text{dev } \tau\|_{L^2(\Omega)} + \|\text{div } \tau\|_{L^2(\Omega)} \right).$$

The proof of Theorem 3.1 requires the following splitting argument from [11].

**Theorem 3.3** (splitting lemma) *Let  $X$  and  $Y$  be (real) Hilbert spaces with  $X = X_1 \times X_2$ . Let  $b_1 : X_1 \times Y \rightarrow \mathbb{R}$  and  $b_2 : X_2 \times Y \rightarrow \mathbb{R}$ , suppose the continuous bilinear form  $b : X \times Y \rightarrow \mathbb{R}$  is their sum, in the sense that for all  $x = (x_1, x_2) \in X_1 \times X_2 = X$  and all  $y \in Y$ ,*

$$b(x, y) = b_1(x_1, y) + b_2(x_2, y).$$

Set  $Y_1 := \{y \in Y : b_2(x_2, y) = 0 \text{ for all } x_2 \in X_2\}$  and suppose, that

$$0 < \beta_1 := \inf_{x_1 \in S(X_1)} \sup_{y_1 \in S(Y_1)} b_1(x_1, y_1), \tag{22}$$

$$0 < \beta_2 := \inf_{x_2 \in S(X_2)} \sup_{y \in S(Y)} b_2(x_2, y), \tag{23}$$

$$N_1 := \{y_1 \in Y_1 : b_1(x_1, y_1) = 0 \text{ for all } x_1 \in X_1\} = \{0\}. \tag{24}$$

Then it follows  $N := \{y \in Y : b(x, y) = 0 \text{ for all } x \in X\} = \{0\}$  and

$$0 < \frac{\sqrt{2}\beta_1\beta_2}{\sqrt{\beta_1^2 + \beta_2^2 + \|b_1\|^2 + \sqrt{(\beta_1^2 + \beta_2^2 + \|b_1\|^2)^2 - 4\beta_1^2\beta_2^2}}} \leq \inf_{x \in S(X)} \sup_{y \in S(Y)} b(x, y).$$

*Proof* This is essentially [11, Thm.3.1] in different notation, but the constant here is slightly better than  $\beta = \beta_1\beta_2 / \left(\beta_1^2 + \beta_2^2 + \|b_1\|^2 + 2\beta_1\|b_1\|\right)^{-1/2}$  and this requires the additional condition (24). Set  $Y_2 := Y_1^\perp$  for an orthogonal split  $Y = Y_1 \oplus Y_2$ . Then  $\beta_2$  from (23) is positive and  $b_2|_{X_2 \times Y_2}$  is non-degenerate in the sense that  $b_2(\bullet, y_2) \neq 0$  in  $X_2^*$  for all  $y_2 \in Y_2 \setminus \{0\}$ . The general theory on bilinear forms [4, Thm.2.1] guarantees, that given any  $x = (x_1, x_2) \in X$ , there exists  $y_2 \in Y_2$  with  $b_2(\bullet, y_2) = (x_2, \bullet)_{X_2}$  in  $X_2^*$ . Hence  $\beta_2\|y_2\|_Y \leq \|x_2\|_{X_2}$ . Since  $\beta_1 > 0$  and  $N_1 = \{0\}$ , there exists a unique  $y_1 \in Y_1$  such that  $b_1(\bullet, y_1) = (\bullet, x_1)_{X_1} - b_1(\bullet, y_2)$  in  $X_1^*$  and  $\beta_1\|y_1\|_Y \leq \|x_1\|_{X_1} + \|b_1\|\|y_2\|_Y$ . Then

$$b(x, y_1 + y_2) = \|x_1\|_X^2 + \|x_2\|_X^2 = \|x\|_X^2.$$

Moreover,  $y = y_1 + y_2 \in Y$  satisfies

$$\begin{aligned} \|y\|_Y^2 &= \|y_1\|_Y^2 + \|y_2\|_Y^2 \leq \beta_1^{-2} \left( \|x_1\|_{X_1} + \beta_2^{-1}\|b_1\|\|x_2\|_{X_2} \right)^2 \\ &\quad + \beta_2^{-2}\|x_2\|_{X_2}^2. \end{aligned}$$

The upper bound is recast as

$$\left( \|x_1\|_{X_1}, \|x_2\|_{X_2} \right) \begin{pmatrix} \beta_1^{-2} & \beta_1^{-2}\beta_2^{-1}\|b_1\| \\ \beta_1^{-2}\beta_2^{-1}\|b_1\| & \beta_1^{-2}\beta_2^{-2}\|b_1\|^2 + \beta_2^{-2} \end{pmatrix} \begin{pmatrix} \|x_1\|_{X_1} \\ \|x_2\|_{X_2} \end{pmatrix} \leq \Lambda \|x\|_X^2$$

for the maximal eigenvalue

$$\Lambda = \frac{\beta_1^2 + \beta_2^2 + \|b_1\|^2 + \sqrt{(\beta_1^2 + \beta_2^2 + \|b_1\|^2)^2 - 4\beta_1^2\beta_2^2}}{2\beta_1^2\beta_2^2}$$

of the displayed symmetric  $2 \times 2$  coefficient matrix. This concludes the proof.  $\square$

*Proof of Theorem 3.1* In the setting of Theorem 3.3, let (equipped with the natural norms)

$$X_1 := L^2(\Omega; \mathbb{R}^{n \times n})/\mathbb{R} \times L^2(\Omega; \mathbb{R}^n), \tag{25}$$

$$X_2 := H_0^{1/2}(\partial T; \mathbb{R}^n) \times H^{-1/2}(\partial T; \mathbb{R}^n), \tag{26}$$

$$Y_1 := H(\text{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R} \times H_0^1(\Omega; \mathbb{R}^n) \subseteq Y. \tag{27}$$

For all  $x_1 = (\sigma, u) \in X_1, x_2 = (s, t) \in X_2$  and  $y = (\tau, v) \in Y$  set

$$b_1(x_1, y) := \int_\Omega \sigma : D_{\text{NC}} v \, dx + \int_\Omega \text{dev } \sigma : \tau \, dx + \int_\Omega u \cdot \text{div}_{\text{NC}} \tau \, dx, \tag{28}$$

$$b_2(x_2, y) := -\left\langle t, \gamma_0^T v \right\rangle_{\partial T} - \left\langle \gamma_v^T \tau, s \right\rangle_{\partial T}. \tag{29}$$

For all  $x_1 = (\boldsymbol{\sigma}, u) \in X_1$  and  $y = (\boldsymbol{\tau}, v) \in Y$ , the Cauchy–Schwarz inequality proves

$$\begin{aligned} |b_1(x_1, y)| &\leq \|\boldsymbol{\sigma}\|_{L^2(\Omega)} (\|\mathbf{D}_{\text{nc}} v\|_{L^2(\Omega)} + \|\boldsymbol{\tau}\|_{L^2(\Omega)}) + \|u\|_{L^2(\Omega)} \|\operatorname{div}_{\text{nc}} \boldsymbol{\tau}\|_{L^2(\Omega)} \\ &\leq \sqrt{2} \|x_1\|_{X_1} \|y\|_Y. \end{aligned}$$

Thus,  $\|b_1\| \leq \sqrt{2}$ . Given  $x_2 = (s, t) \in X_2$  and  $y = (\boldsymbol{\tau}, v) \in Y$ . The substitution of  $s = \gamma_0^T w$  with  $w \in H_0^1(\Omega; \mathbb{R}^n)$  and  $\|w\|_{H^1(\Omega)} = \|s\|_{H^{1/2}(\partial T)}$  as well as  $t = \gamma_v^T \boldsymbol{q}$  with  $\boldsymbol{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$  and  $\|\boldsymbol{q}\|_{H(\operatorname{div})} = \|t\|_{H^{-1/2}(\partial T)}$  allow an integration by parts. Hence,

$$\begin{aligned} b_2(x_2, y) &= - \int_{\Omega} \operatorname{div}_{\text{nc}} \boldsymbol{\tau} \cdot w \, dx - \int_{\Omega} \boldsymbol{\tau} : \mathbf{D} w \, dx - \int_{\Omega} \operatorname{div} \boldsymbol{q} \cdot v \, dx \\ &\quad - \int_{\Omega} \boldsymbol{q} : \mathbf{D}_{\text{nc}} v \, dx \leq \|x_2\|_{X_2} \|y\|_Y. \end{aligned}$$

It follows,  $\|b_2\| \leq 1$  and so  $\|b\| \leq \sqrt{3}$ .

For an arbitrary  $0 \neq x_1 = (\boldsymbol{\sigma}, u) \in X_1$ , define  $\tilde{F} \in Z^*$  by

$$\tilde{F}(\boldsymbol{\rho}, w) := \int_{\Omega} (\boldsymbol{\sigma} : \boldsymbol{\rho} + u \cdot w) \, dx \quad \text{for all } (\boldsymbol{\rho}, w) \in Z. \tag{30}$$

The Cauchy–Schwarz inequality implies  $\|\tilde{F}\|_{Z^*} \leq \sqrt{2} \|x_1\|_{X_1}$ . Since the formulation (20) for the Stokes equations has unique solutions [9, Thm.2.3], there exists  $(\boldsymbol{\tau}, -v) \in Z$  such that  $\tilde{b}((\boldsymbol{\tau}, -v), \bullet) = \tilde{F}$  in  $Z^*$ . For any  $(\boldsymbol{\rho}, w) \in Z$ , this reads

$$0 = \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\rho} \, dx + \int_{\Omega} \operatorname{div} \boldsymbol{\rho} \cdot v \, dx + \int_{\Omega} (u - \operatorname{div} \boldsymbol{\tau}) \cdot w \, dx. \tag{31}$$

Since  $w \in L^2(\Omega; \mathbb{R}^n)$  and  $\boldsymbol{\rho}$  is arbitrary in  $H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ ,  $\operatorname{div} \boldsymbol{\tau} = u$  and (31) implies  $v \in H_0^1(\Omega; \mathbb{R}^n)$  with  $\mathbf{D} v = \boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}$ . This test function  $y_1 := (\boldsymbol{\tau}, v) \in Y_1$  allows for

$$b_1(x_1, y_1) = \|x_1\|_{X_1}^2.$$

Recall  $\gamma$  from (21) the inf-sup constant for  $\tilde{b}$ . Then

$$\begin{aligned} \gamma \|y_1\|_Z &= \gamma \|(\boldsymbol{\tau}, -v)\|_Z \leq \|\tilde{b}((\boldsymbol{\tau}, -v), \bullet)\|_{Z^*} \\ &= \|\tilde{F}\|_{Z^*} \leq \sqrt{2} \|x_1\|_{X_1}. \end{aligned}$$

The triangle inequality implies  $\|\mathbf{D} v\|_{L^2(\Omega)}^2 \leq 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + 2\|\boldsymbol{\tau}\|_{L^2(\Omega)}^2$ . The previous two displayed inequalities prove

$$\|y_1\|_Y^2 \leq (6\gamma^{-2} + 2) \|x_1\|_{X_1}^2.$$

Hence, for all  $x_1 = (\boldsymbol{\sigma}, u) \in X_1$  and  $y_1 := (\boldsymbol{\tau}, v) \in Y_1$  as above,

$$(6\gamma^{-2} + 2)^{-1/2} \|x_1\|_{X_1} \leq b_1(x_1, y_1) / \|y_1\|_Y \leq \sup_{y_1 \in S(Y_1)} b_1(x_1, y_1).$$

This proves (22) with  $(6\gamma^{-2} + 2)^{-1/2} \leq \beta_1$ .

The duality Lemma 2.1 shows, that any  $x_2 = (s, t) \in X_2$  satisfies

$$\begin{aligned} \|x_2\|_{X_2} &\leq \|s\|_{H^{1/2}(\partial T)} + \|t\|_{H^{-1/2}(\partial T)} \\ &= \sup_{\mathbf{q} \in S(H(\operatorname{div}, T; \mathbb{R}^{n \times n})/\mathbb{R})} \left\langle \gamma_v^T \mathbf{q}, s \right\rangle_{\partial T} + \sup_{w \in S(H^1(T; \mathbb{R}^n))} \left\langle t, \gamma_0^T w \right\rangle_{\partial T} \\ &\leq \sup_{\substack{\mathbf{q} \in S(H(\operatorname{div}, T; \mathbb{R}^{n \times n})/\mathbb{R}) \\ w \in S(H^1(T; \mathbb{R}^n))}} b_2(x_2, (\mathbf{q}, w)) \leq \sqrt{2} \sup_{y \in S(Y)} b_2(x_2, y). \end{aligned}$$

Hence, (23) holds with  $2^{-1/2} \leq \beta_2$ .

Given any  $y = (\boldsymbol{\tau}, v) \in Y$  with  $b_2(x_2, y) = -\left\langle \gamma_v^T \boldsymbol{\tau}, s \right\rangle_{\partial T} - \left\langle t, \gamma_0^T v \right\rangle_{\partial T} = 0$  for all  $x_2 = (s, t) \in X_2$ . This means that all jumps of  $v$  and (normal components) of  $\boldsymbol{\tau}$  disappear. Hence,  $y \in Y_1$  as demanded in Theorem 3.3.

Let  $y_1 = (\boldsymbol{\tau}, v) \in N_1$ . With  $x_1 = (0, u) \in X_1$  for any  $u \in C_0^\infty(\Omega; \mathbb{R}^n) \subseteq L^2(\Omega; \mathbb{R}^n)$ ,  $y_1 \in N_1$  implies  $\operatorname{div} \boldsymbol{\tau} \equiv 0$ . The boundary conditions and continuity in  $Y_1$  prove

$$0 = \int_{\partial \Omega} v \cdot \boldsymbol{\tau} \, ds = \int_{\Omega} v \cdot \operatorname{div} \boldsymbol{\tau} \, dx + \int_{\Omega} \operatorname{D} v : \boldsymbol{\tau} \, dx = \int_{\Omega} \operatorname{D} v : \boldsymbol{\tau} \, dx.$$

Furthermore, the choice  $x_1 = (\boldsymbol{\tau}, 0) \in X_1$  results in

$$0 = \int_{\Omega} \boldsymbol{\tau} : \operatorname{D} v \, dx + \int_{\Omega} \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\tau} \, dx = \int_{\Omega} \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\tau} \, dx = \|\operatorname{dev} \boldsymbol{\tau}\|_{L^2(\Omega)}^2;$$

whence  $\operatorname{dev} \boldsymbol{\tau} = 0$ . Lemma 3.2 proves  $\boldsymbol{\tau} \equiv 0$ . Further, for all  $\boldsymbol{\sigma} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$  set  $\tilde{\boldsymbol{\sigma}} := \boldsymbol{\sigma} - 1/n (\int_{\Omega} \operatorname{tr} \boldsymbol{\sigma} \, dx) \mathbf{I}_{n \times n}$  and  $x_1 = (\tilde{\boldsymbol{\sigma}}, 0) \in X_1$ . Then

$$\int_{\Omega} \boldsymbol{\sigma} : \operatorname{D} v \, dx = \int_{\Omega} \tilde{\boldsymbol{\sigma}} : \operatorname{D} v \, dx + \frac{1}{n} \left( \int_{\Omega} \operatorname{tr} \boldsymbol{\sigma} \, dx \right) \int_{\Omega} \operatorname{div} v \, dx = b_1(x_1, y_1) = 0.$$

Hence,  $\operatorname{D} v \equiv 0$  for  $v \in H_0^1(\Omega; \mathbb{R}^n)$  and so  $v \equiv 0$ . This concludes the proof. □

### 4 Discrete problem

The low-order discrete trial and test search space of the introduced method read

$$X_h := P_0(T; \mathbb{R}^{n \times n})/\mathbb{R} \times P_0(T; \mathbb{R}^n) \times S_0^1(\mathcal{E}; \mathbb{R}^n) \times P_0(\mathcal{E}; \mathbb{R}^n), \tag{32}$$

$$Y_h := RT_0^{\text{PW}}(T; \mathbb{R}^{n \times n})/\mathbb{R} \times P_1(T; \mathbb{R}^n). \tag{33}$$

Given  $b$  from (18) and  $F$  from (19), the discrete problem seeks  $x_h \in X_h$  with (3). This section establishes the discrete inf-sup condition (4) with a constant  $\beta_h$ , which depends on the Friedrichs constant  $C_F$  (with  $\|\bullet\|_{L^2(\Omega)} \leq C_F \|D\bullet\|_{L^2(\Omega)}$  in  $H_0^1(\Omega)$ ) and the tr-div-dev constant  $C_{\text{idd}}$ .

**Theorem 4.1** (inf-sup) *The discrete spaces (32)–(33) and the bilinear form  $b$  from (18) satisfy*

$$1 \lesssim \beta_h := \inf_{x_h \in X_h \setminus \{0\}} \sup_{y_h \in Y_h} \frac{b(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y}.$$

*Proof Step 1. Discrete test functions.* The discrete traces in  $S_0^1(\mathcal{E}; \mathbb{R}^n)$  (resp.  $P_0(\mathcal{E}; \mathbb{R}^n)$ ) admit a unique extension by  $S_0^1(\mathcal{T}; \mathbb{R}^n)$  (resp.  $RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ ). Thus, given  $x_h = (\sigma_0, u_0, s_1, t_0) \in X_h$  chose  $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^n)$  with  $\gamma_0^T w_c = s_1$  and  $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$  with  $\gamma_v^T \mathbf{q}_{\text{RT}} = t_0$ . The norm for the trace space in Sect. 2.3 by minimal extension fulfils

$$\begin{aligned} \|x_h\|_X^2 &= \|\sigma_0\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2 + \|\gamma_0^T w_c\|_{H^{1/2}(\partial\mathcal{T})}^2 + \|\gamma_v^T \mathbf{q}_{\text{RT}}\|_{H^{-1/2}(\partial\mathcal{T})}^2 \\ &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2 + \|w_c\|_{H^1(\Omega)}^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2. \end{aligned} \tag{34}$$

For  $x_h = (\sigma_0, u_0, \gamma_0^T w_c, \gamma_v^T \mathbf{q}_{\text{RT}}) \in X_h \setminus \{0\}$ , set  $y_h = (\tau_{\text{RT}}, v_1) \in Y_h$

$$\begin{aligned} \tau_{\text{RT}} &:= \text{dev } \sigma_0 - D w_c + 1/n (u_0 - \Pi_0 w_c) \otimes (\bullet - \text{mid}(\mathcal{T})), \\ v_1 &:= -\text{div } \mathbf{q}_{\text{RT}} + (\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}) (\bullet - \text{mid}(\mathcal{T})). \end{aligned}$$

Notice, that  $\text{div}_{\text{NC}} \tau_{\text{RT}} = u_0 - \Pi_0 w_c$ ,  $D_{\text{NC}} v_1 = \sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}$  and  $\Pi_0 v_1 = -\text{div } \mathbf{q}_{\text{RT}}$ . The side restriction  $\int_{\Omega} \text{tr } \sigma_0 \, dx = 0$  implies  $\int_{\Omega} \text{tr } \tau_{\text{RT}} \, dx = 0$ .

Furthermore, the substitution of  $s_1$  by  $\gamma_0^T w_c$  and  $t_0$  by  $\gamma_v^T \mathbf{q}_{\text{RT}}$  allows an integration by parts. Hence,  $x_h = (\sigma_0, u_0, \gamma_0^T w_c, \gamma_v^T \mathbf{q}_{\text{RT}})$  and the above test function  $y_h = (\tau_{\text{RT}}, v_1)$  satisfy

$$\begin{aligned} b(x_h, y_h) &= \int_{\Omega} \sigma_0 : D_{\text{NC}} v_1 \, dx + \int_{\Omega} \text{dev } \sigma_0 : \tau_{\text{RT}} \, dx + \int_{\Omega} u_0 \cdot \text{div}_{\text{NC}} \tau_{\text{RT}} \, dx \\ &\quad - \int_{\Omega} (v_1 \cdot \text{div } \mathbf{q}_{\text{RT}} + \mathbf{q}_{\text{RT}} : D_{\text{NC}} v_1) \, dx \\ &\quad - \int_{\Omega} (\tau_{\text{RT}} : D w_c + w_c \cdot \text{div}_{\text{NC}} \tau_{\text{RT}}) \, dx \\ &= \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)}^2 \\ &\quad + \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned} \tag{35}$$

*Step 2. Key estimates.* The test function from Step 1. and (9) prove

$$\|y_h\|_Y^2 = \|\tau_{\text{RT}}\|_{H(\text{div}, \mathcal{T})}^2 + \|v_1\|_{H^1(\mathcal{T})}^2$$

$$\begin{aligned} &\leq \|\operatorname{dev} \boldsymbol{\sigma}_0 - \mathbf{D} w_c\|_{L^2(\Omega)}^2 + \left(h_{\max}^2/(n+1)^2 + 1\right) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 \\ &\quad + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \left(h_{\max}^2 n^2/(n+1)^2 + 1\right) \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned}$$

The combination with (35) shows for  $C_{y_h} := 1 + h_{\max}^2 n^2/(n+1)^2$  holds  $\|y_h\|_Y^2 \leq C_{y_h} b(x_h, y_h)$ . The proof of  $\|x_h\|_X \lesssim b(x_h, y_h)$  requires the computation of  $C_{w_c}$  with  $\|\mathbf{D} w_c\|_{L^2(\Omega)} =: \|w_c\|^2 \leq C_{w_c} b(x_h, y_h)$ . The function

$$\tilde{\mathbf{q}}_{\text{RT}} := \mathbf{q}_{\text{RT}} - 1/n \left( \int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \mathbf{I}_{n \times n} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R} \quad (36)$$

allows an application of Lemma 3.2. Moreover,

- (i)  $\|\operatorname{dev} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} = \|\operatorname{dev} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)}$  and  $\|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} = \|\operatorname{div} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)}$ ,
  - (ii) for  $\boldsymbol{\sigma}_0 \in L^2(\Omega; \mathbb{R}^{n \times n})/\mathbb{R}$  holds  $\|\boldsymbol{\sigma}_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \leq \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}$ ,
  - (iii) for  $f \in L^2(\Omega; \mathbb{R})$  with  $\int_{\Omega} f \, dx = 0$  holds  $\int_{\Omega} f \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx = \int_{\Omega} f \operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}} \, dx$ .
- This verifies

$$C_{\text{td}}^{-1} \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \leq \|\operatorname{dev} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}. \quad (37)$$

It holds

$$\begin{aligned} \|\operatorname{dev} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} &\leq \|\operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}})\|_{L^2(\Omega)} + \|\operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{D} w_c)\|_{L^2(\Omega)} \\ &\quad + \|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)}. \end{aligned}$$

From (13) and (9) it follows

$$\|\operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}})\|_{L^2(\Omega)} \leq \|\operatorname{dev}(\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}})\|_{L^2(\Omega)} + \frac{h_{\max}}{n+1} \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}.$$

This proves

$$\begin{aligned} C_{\text{td}}^{-1} \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} &\leq \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{D} w_c)\|_{L^2(\Omega)} \\ &\quad + \|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)} + (1 + h_{\max}/(n+1)) \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}. \end{aligned}$$

On the other hand,

$$\begin{aligned} \|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)}^2 &= \int_{\Omega} (\Pi_0(\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}) - \operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{D} w_c) + \mathbf{q}_{\text{RT}}) : \operatorname{dev} \mathbf{D} w_c \, dx \\ &\leq \left( \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{D} w_c)\|_{L^2(\Omega)} \right) \|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)} \\ &\quad + \int_{\Omega} \tilde{\mathbf{q}}_{\text{RT}} : \operatorname{dev} \mathbf{D} w_c \, dx. \end{aligned}$$



The decomposition of the deviator followed by an integration by parts shows

$$\begin{aligned} \int_{\Omega} \operatorname{dev} \tilde{\mathbf{q}}_{\text{RT}} : \mathbf{D} w_c \, dx &= \int_{\Omega} \tilde{\mathbf{q}}_{\text{RT}} : \mathbf{D} w_c \, dx - 1/n \int_{\Omega} (\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}) \operatorname{div} w_c \, dx \\ &\leq - \int_{\Omega} w_c \cdot \operatorname{div} \mathbf{q}_{\text{RT}} \, dx + 1/n \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \|\operatorname{div} w_c\|_{L^2(\Omega)} \\ &\leq C_F \|w_c\| \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\quad + 1/n \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \|\operatorname{div} w_c\|_{L^2(\Omega)}. \end{aligned}$$

The combination of the aforementioned estimates leads with

$$\begin{aligned} a &:= \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{dev}(\sigma_0 - \mathbf{D} w_c)\|_{L^2(\Omega)} + C_{\text{tdd}}/n \|\operatorname{div} w_c\|_{L^2(\Omega)}, \\ b &:= C_{\text{tdd}} \left( \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{dev}(\sigma_0 - \mathbf{D} w_c)\|_{L^2(\Omega)} \right) + \|\operatorname{div} w_c\|_{L^2(\Omega)} \\ &\quad + C_{\text{tdd}}(1 + h_{\max}/(n + 1)) \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}, \text{ and } c := C_F \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ \text{to } \|w_c\|^2 &\leq a \|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)} + b/n \|\operatorname{div} w_c\|_{L^2(\Omega)} + c \|w_c\|. \end{aligned}$$

This upper bound is the scalar product in  $\mathbb{R}^3$  of the vector  $(a, b/\sqrt{n}, c)$  with  $(\|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)}, \|\operatorname{div} w_c\|_{L^2(\Omega)}/\sqrt{n}, \|w_c\|)$ . The Cauchy-Schwarz inequality leads to

$$\begin{aligned} \|w_c\|^2 &\leq \sqrt{\|\operatorname{dev} \mathbf{D} w_c\|_{L^2(\Omega)}^2 + 1/n \|\operatorname{div} w_c\|_{L^2(\Omega)}^2} + \|w_c\| \sqrt{a^2 + b^2/n + c^2} \\ &= \|w_c\| \sqrt{2} \sqrt{a^2 + b^2/n + c^2} =: C_1 \|w_c\|. \end{aligned}$$

The Cauchy-Schwarz inequality, (8), and the abbreviations  $g := \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}$ ,  $f := \|\mathbf{D} w_c - \operatorname{dev} \sigma_0\|_{L^2(\Omega)}$  and  $e := \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}$  allow to rewrite the pre-factors  $a$ ,  $b$ , and  $c$  as

$$\begin{aligned} a &\leq \sqrt{2 + C_{\text{tdd}}^2/n} \sqrt{e^2 + f^2}, \\ b &\leq \sqrt{e^2 + f^2 + g^2} \sqrt{n + C_{\text{tdd}}^2(2 + (1 + h_{\max}/(n + 1))^2)} \text{ and } c = C_F g. \end{aligned}$$

Since by (35),  $e^2 + f^2 + g^2 \leq b(x_h, y_h)$ , it follows

$$C_1^2 \leq 2 \left( \max \left\{ C_F^2, 2 + \frac{C_{\text{tdd}}^2}{n} \right\} + 1 + \frac{C_{\text{tdd}}^2}{n} \left( 2 + \left( 1 + \frac{h_{\max}}{n + 1} \right)^2 \right) \right) b(x_h, y_h).$$

Therefore, the constant  $C_{w_c}$  with  $\|w_c\|^2 \leq C_{w_c} b(x_h, y_h)$  satisfies

$$C_{w_c} \leq 2 \left( \max \left\{ C_F^2, 2 + C_{\text{tdd}}^2/n \right\} + 1 + C_{\text{tdd}}^2/n(2 + (1 + h_{\max}/(n + 1))^2) \right).$$

It remains to prove  $\|x_h\|_X^2 \lesssim b(x_h, y_h)$ . For  $\mathbf{q}_{\text{RT}} \in \text{RT}_0(\mathcal{T}; \mathbb{R}^{n \times n})$ , (13) and (9) imply

$$\|\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \leq \frac{h_{\max}^2}{(n+1)^2} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + 2\|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + 2\|\sigma_0\|_{L^2(\Omega)}^2.$$

On the other hand, the auxiliary function  $\tilde{\mathbf{q}}_{\text{RT}}$  from (36)–(37) allows for

$$\begin{aligned} \|\sigma_0\|_{L^2(\Omega)} &\leq \|\sigma_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + \frac{C_{\text{tdd}}}{\sqrt{n}} \|\text{div } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &\quad + \left(1 + \frac{C_{\text{tdd}}}{\sqrt{n}}\right) \|\text{dev } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq (2 + C_{\text{tdd}}/\sqrt{n}) \|\sigma_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_{\text{tdd}}/\sqrt{n} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\quad + (1 + C_{\text{tdd}}/\sqrt{n}) \|\text{dev } \sigma_0\|_{L^2(\Omega)}. \end{aligned}$$

Furthermore, it holds

$$\begin{aligned} \|\sigma_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} &\leq \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + h_{\max}/(n+1) \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \quad \text{and} \\ \|\text{dev } \sigma_0\|_{L^2(\Omega)} &\leq \|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)} + \|w_c\|. \end{aligned}$$

Therefore, all terms in the decomposition of  $\|x_h\|_X$  as in (34) are under control,

$$\begin{aligned} \|x_h\|_X^2 &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2 + \|w_c\|_{H^1(\Omega)}^2 \\ &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + (1 + C_F^2) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 \\ &\quad + (2 + 2C_F^2) \|w_c\|^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2. \end{aligned}$$

Careful bookkeeping reveals that

$$\begin{aligned} C_3 &:= \frac{3h_{\max}^2}{(n+1)^2} \left(\frac{C_{\text{tdd}}}{\sqrt{n}} + 2\right)^2 + \frac{6h_{\max}}{n+1} \left(\frac{C_{\text{tdd}}^2}{n} + 2\frac{C_{\text{tdd}}}{\sqrt{n}}\right) \\ &\quad + 12 \left(\frac{C_{\text{tdd}}}{\sqrt{n}} + 1\right)^2 + 6 \end{aligned}$$

satisfies

$$\begin{aligned} \|x_h\|_X^2 &\leq \left(1 + C_F^2\right) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + C_3 \|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)}^2 \\ &\quad + (2 + C_3) \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + (2 + 2C_F^2 + C_3) \|w_c\|^2 \\ &\quad + (1 + h_{\max}^2/(n+1)^2 + C_3) \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned}$$

Therefore,  $\|x_h\|_X^2 \leq C_{x_h} b(x_h, y_h)$  holds for

$$C_{x_h} := \max \left\{ 1 + C_F^2, C_3 + \max \left\{ 2, 1 + \frac{h_{\max}^2}{(n+1)^2} \right\} \right\} + (C_3 + 2C_F^2 + 2) C_{w_c}.$$

*Step 3.* Altogether, for all  $x_h = (\sigma_0, u_0, \gamma_0^T w_c, \gamma_v^T \mathbf{q}_{RT}) \in X_h \setminus \{0\}$  and  $y_h \in Y_h$  as in *Step 1*, it holds

$$\|x_h\|_X \leq \sqrt{C_{x_h} C_{y_h}} b(x_h, y_h) / \|y_h\|_Y.$$

This concludes the proof of  $0 < (C_{x_h} C_{y_h})^{-1/2} \leq \beta_h$  with  $\beta_h$  from (4). □

### 5 Data approximation error

The Fortin interpolator (5) is explicitly constructed in [10, 19, 22] with higher-order test search functions. The low order spaces in [13, 14] require a direct verification of the discrete inf-sup condition and allow explicit constants  $\|b\|$ ,  $\beta_h$  in the a posteriori error bound (7). The upper error bound involves the computable residual error  $\|F - b(\xi_h, \bullet)\|_{Y_h^*}$  and the remaining data approximation error  $\|F \circ (1 - \Pi)\|_{Y^*}$ . The latter is *not* of higher-order in general as shown in Sect. 5.1. This motivates an extension of the test search space in Sect. 5.2.

#### 5.1 Fortin interpolation

The description of the operator  $\Pi: Y \rightarrow Y_h$  with (5) in 2D with a shape-regular triangulation  $\mathcal{T}$  of the simply-connected bounded polygonal domain  $\Omega \subset \mathbb{R}^2$  into triangles requires further notation. For  $\beta \in C^1(\Omega; \mathbb{R}^2)$ , set

$$\text{Curl } \beta := \begin{pmatrix} -\partial\beta_1/\partial x_2 & \partial\beta_1/\partial x_1 \\ -\partial\beta_2/\partial x_2 & \partial\beta_2/\partial x_1 \end{pmatrix}, \quad \text{curl } \beta := \text{tr}(\text{Curl } \beta) = \partial\beta_2/\partial x_1 - \partial\beta_1/\partial x_2$$

with the piecewise version  $\text{Curl}_{\text{nc}}$  and  $\text{curl}_{\text{nc}}$  (piecewise with respect to  $\mathcal{T}$ ). Define

$$\begin{aligned} X_{\text{curl}} &:= \left\{ v_C \in S^1(\mathcal{T}; \mathbb{R}^2) : \int_{\Omega} v_C \, dx = 0, \int_{\Omega} \text{curl } v_C \, dx = 0 \right\} \\ &\equiv S^1(\mathcal{T}; \mathbb{R}^2) / \mathbb{R}^3. \end{aligned}$$

The nonconforming Crouzeix-Raviart functions space reads

$$\begin{aligned} \text{CR}^1(\mathcal{T}; \mathbb{R}^2) &:= \{v \in P_1(\mathcal{T}, \mathbb{R}^2) : v \text{ is continuous in } \text{mid}(E) \text{ for all } E \in \mathcal{E}(\Omega)\}, \\ \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2) &:= \{v \in \text{CR}^1(\mathcal{T}; \mathbb{R}^2) : v(\text{mid}(E)) = 0 \text{ for all } E \in \mathcal{E}(\partial\Omega)\}, \\ \text{CR}^1(\mathcal{T}, \mathbb{R}^2) / \mathbb{R}^2 &:= \{v \in \text{CR}^1(\mathcal{T}, \mathbb{R}^2) : \int_{\Omega} v \, dx = 0\}. \end{aligned}$$

The discrete divergence-free Crouzeix-Raviart functions

$$Z_{CR} := \{v \in CR_0^1(\mathcal{T}; \mathbb{R}^2) : \operatorname{div}_{NC} v = 0\} \subseteq CR_0^1(\mathcal{T}; \mathbb{R}^2)$$

are well known from the the nonconforming finite element analysis of the Stokes equations. Any  $w_{CR} \in CR_0^1(\mathcal{T}; \mathbb{R}^2)$  satisfies  $\int_E [w_{CR}]_E ds = 0$  along any edge  $E \in \mathcal{E}$ . The local nonconforming interpolant  $I_{NC}^{pw}$  guarantees a similar property: Define  $I_{NC}^{pw} v \in P_1(\mathcal{T}; \mathbb{R}^2)$  for  $v \in H^1(\mathcal{T}; \mathbb{R}^2)$  on any  $T \in \mathcal{T}$  via

$$(I_{NC}^{pw} v)|_T (\operatorname{mid}(E)) = \int_E v|_T ds := |E|^{-1} \int_E v|_T ds. \tag{38}$$

For all  $T \in \mathcal{T}$  and  $E \in \mathcal{E}(T)$  holds  $\int_E I_{NC}^{pw} v|_T = \int_E v|_T ds$ , whence  $\Pi_0 D_{NC} v = D_{NC} I_{NC}^{pw} v$ .

**Lemma 5.1** (discrete Helmholtz decompositions) *For a simply-connected domain  $\Omega$ , the following decompositions are orthogonal in  $L^2(\Omega; \mathbb{R}^{2 \times 2})$*

$$P_0(\mathcal{T}; \mathbb{R}^{2 \times 2}_{dev}) = D_{NC} Z_{CR} \oplus \operatorname{dev} \operatorname{Curl} X_{curl}, \tag{39}$$

$$P_0(\mathcal{T}; \mathbb{R}^{2 \times 2}) = D S_0^1(\mathcal{T}; \mathbb{R}^2) \oplus \operatorname{Curl}_{NC} CR^1(\mathcal{T}, \mathbb{R}^2)/\mathbb{R}^2. \tag{40}$$

*Proof* The paper [17] includes a proof of (39) and (40) is known from [3]. □

Based on those preliminaries, the Fortin interpolation is characterized in the sequel. Given a simply-connected domain  $\Omega$  and  $y = (\tau, v) \in Y$ . Let  $\alpha_{CR} \in CR^1(\mathcal{T}; \mathbb{R}^2)/\mathbb{R}^2$  satisfy  $\int_{\Omega} \alpha_{CR} dx = 0$  and

$$\int_{\Omega} (D_{NC} \alpha_{CR} - \Pi_0 \tau) : D_{NC} w_{CR} dx = \int_{\Omega} w_{CR} \cdot (1 - \Pi_0) \operatorname{div}_{NC} \tau dx \tag{41}$$

for all  $w_{CR} \in CR^1(\mathcal{T}; \mathbb{R}^2)$ . (This follows from one solve of the Crouzeix-Raviart FEM and  $\int_{\Omega} (1 - \Pi_0) \operatorname{div}_{NC} \tau dx = 0$ .) Let  $\alpha_0 := -1/2 \int_{\Omega} \operatorname{div}_{NC} \alpha_{CR}$ . The discrete Helmholtz decomposition (39) guarantees the existence of  $z_{CR} \in Z_{CR}$  and  $\beta_c \in X_{curl}$ , such that

$$\operatorname{dev} (\Pi_0 \tau - D_{NC} \alpha_{CR}) = D_{NC} z_{CR} + \operatorname{dev} \operatorname{Curl} \beta_c. \tag{42}$$

**Theorem 5.2** *Given  $(\tau, v) \in Y$  and  $\alpha_{CR} \in CR^1(\mathcal{T}; \mathbb{R}^2)/\mathbb{R}^2, z_{CR} \in Z_{CR}, \beta_c \in X_{curl}, \alpha_0 \in \mathbb{R}^2$  as above with (41)–(42), set*

$$\begin{aligned} \tau_{RT} &:= D_{NC} \alpha_{CR} + \operatorname{Curl} \beta_c + \alpha_0 I_{2 \times 2} + (\Pi_0 \operatorname{div}_{NC} \tau / 2) \otimes (\bullet - \operatorname{mid}(T)), \\ v_1 &:= I_{NC}^{pw} v + z_{CR}. \end{aligned}$$

*The mapping  $\Pi : Y \rightarrow Y_h, (\tau, v) \mapsto (\tau_{RT}, v_1)$  is linear, bounded, idempotent and fulfils (5). The discrete kernel  $N_h := \{y_h \in Y_h : b(x_h, y_h) = 0 \forall x_h \in X_h\}$  of  $B_{2,h} : Y_h \rightarrow X_h^*, y_h \mapsto b(\bullet, y_h)|_{Y_h}$  has dimension  $\dim(N_h) = 2(|\mathcal{T}| - 1)$  and is equal to*

$$N_h = \left\{ (\text{Curl}_{NC} \beta_{CR}, v_{CR}) \in \text{Curl}_{NC} \text{CR}^1(\mathcal{T}, \mathbb{R}^2)/\mathbb{R}^2 \times Z_{CR} : \right. \\ \left. - D_{NC} v_{CR} = \text{dev Curl}_{NC} \beta_{CR} \right\}.$$

*Proof* The design of  $\alpha_0$  leads to  $\int_{\Omega} \text{tr } \tau_{RT} \, dx = 0$ . For all  $\sigma_0 \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ , the split (42) and  $\Pi_0 D_{NC} v = D_{NC} I_{NC}^{pw} v$  prove

$$b((\sigma_0, 0, 0, 0), y - y_h) = \int_{\Omega} \sigma_0 : D_{NC}(v - v_1) \, dx + \int_{\Omega} \text{dev } \sigma_0 : (\tau - \tau_{RT}) \, dx \\ = \int_{\Omega} \sigma_0 : (\Pi_0 D_{NC}(v - v_1) + \text{dev } \Pi_0(\tau - \tau_{RT})) \, dx \\ = \int_{\Omega} \sigma_0 : (-D_{NC} z_{CR} + \text{dev}(\Pi_0 \tau - D_{NC} \alpha_{CR} - \text{Curl } \beta_c)) \, dx \\ = 0.$$

Since  $\text{div}_{NC} \tau_{RT} = \Pi_0 \text{div}_{NC} \tau$ , any  $u_0 \in P_0(\mathcal{T}; \mathbb{R}^2)$  satisfies

$$b((0, u_0, 0, 0), y - y_h) = \int_{\Omega} u_0 \cdot \text{div}_{NC}(\tau - \tau_{RT}) \, dx \\ = \int_{\Omega} u_0 \cdot \Pi_0 \text{div}_{NC}(\tau - \tau_{RT}) \, dx = 0.$$

For all  $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^2)$  on the skeleton  $\partial\mathcal{T}$ , consider the linear extension  $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2) \subseteq \text{CR}^1(\mathcal{T}, \mathbb{R}^2)$  with  $\gamma_0^T w_c = s_1$  to allow an integration by parts. Thus,  $\int_{\Omega} D w_c \, dx = 0$ , (40)–(41), and  $\text{div}_{NC} \tau_{RT} = \Pi_0 \text{div}_{NC} \tau$  prove

$$-b((0, 0, s_1, 0), y - y_h) = \int_{\Omega} D w_c : \Pi_0(\tau - \tau_{RT}) \, dx \\ + \int_{\Omega} w_c \cdot \text{div}_{NC}(\tau - \tau_{RT}) \, dx \\ = \int_{\Omega} D w_c : (\Pi_0 \tau \, dx - D_{NC} \alpha_{CR} - \text{Curl } \beta_c - \alpha_0 I_{2 \times 2}) \, dx \\ + \int_{\Omega} w_c \cdot (1 - \Pi_0) \text{div}_{NC} \tau \, dx = 0.$$

The properties of the Crouzeix-Raviart-functions and  $I_{NC}^{pw}$  (38) prove for all  $t_0 \in P_0(\mathcal{E}; \mathbb{R}^2)$ ,

$$-b((0, 0, 0, t_0), y - y_h) \\ = \sum_{E \in \mathcal{E}} t_0|_E \cdot \left( \int_E [v - I_{NC}^{pw} v]_E \, ds - \int_E [z_{CR}]_E \, ds \right) = 0.$$

For the proof, that  $\Pi$  is idempotent (hence a projection), suppose that  $(\tau, v) \in RT_0^{pw}(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^2)$  and decompose  $\Pi_0 \tau = D_{NC} a_{CR} + \text{Curl } b_c$  for unique

$a_{\text{CR}} \in CR^1(\mathcal{T}; \mathbb{R}^2)/\mathbb{R}^2$  and  $b_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$ . Since  $(1 - \Pi_0) \operatorname{div}_{\text{NC}} \boldsymbol{\tau} = 0$  a.e. in  $\Omega$ , (41) shows  $a_{\text{CR}} = \alpha_{\text{CR}}$ . Since  $b_c = 0$  along  $\partial\Omega$  and  $b_c - \int_{\Omega} b_c \, dx \in X_{\text{curl}}$ , (42) reveals that  $\beta_c = b_c - \int_{\Omega} b_c \, dx$  and  $z_{\text{CR}} = 0$ . Notice that  $0 = \int_{\Omega} \operatorname{tr} \boldsymbol{\tau} \, dx = \int_{\Omega} \operatorname{tr} \Pi_0 \boldsymbol{\tau} \, dx = \int_{\Omega} \operatorname{div}_{\text{NC}} \alpha_{\text{CR}} \, dx = 0$  implies  $\alpha_0 = 0$ . Altogether, it follows that  $\boldsymbol{\tau}_{\text{RT}} = \boldsymbol{\tau}$  and  $v_1 = v$ , i.e.,  $\Pi^2 = \Pi$ .

The discrete Friedrichs and Poincaré inequality [7, Thm.10.6.12], Lemma 3.2, and [12, Thm.4] show that the proposed mapping  $\Pi$  is bounded. A lengthy but straight forward calculation with  $C := 2 \max \left\{ 1 + C_{\text{td}}^2, 2 + 2C_{\text{df}}^2 \right\}$  reveals

$$\|\Pi\|^2 \leq (1 + C)(1 + C_{\text{dp}}^2) + \max \left\{ C, 1 + h_{\text{max}}^2/9 \right\}.$$

It remains to characterize  $N_h$ . For all  $\tilde{y}_h = (\tilde{\boldsymbol{\tau}}_{\text{RT}}, \tilde{v}_1) \in N_h$ , the condition

$$0 = b((0, 0, 0, t_0), \tilde{y}_h) = - \sum_{E \in \mathcal{E}} t_0|_E \int_E [\tilde{v}_1]_E \, ds \quad \text{for all } t_0 \in P_0(\mathcal{E}; \mathbb{R}^2)$$

implies  $\tilde{v}_1 \in CR_0^1(\mathcal{T}; \mathbb{R}^2)$ . Since

$$0 = b((0, u_0, 0, 0), \tilde{y}_h) = \int_{\Omega} u_0 \cdot \operatorname{div}_{\text{NC}} \tilde{\boldsymbol{\tau}}_{\text{RT}} \, dx \quad \text{for all } u_0 \in P_0(\mathcal{T}; \mathbb{R}^2),$$

it follows  $\operatorname{div}_{\text{NC}} \tilde{\boldsymbol{\tau}}_{\text{RT}} = 0$  and  $\Pi_0 \tilde{\boldsymbol{\tau}}_{\text{RT}} = \tilde{\boldsymbol{\tau}}_{\text{RT}}$ . The linear extension of  $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^2)$  to  $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$  with  $\gamma_0^T w_c = s_1$  and an integration by parts result in

$$0 = b((0, 0, s_1, 0), \tilde{y}_h) = - \int_{\Omega} D w_c : \tilde{\boldsymbol{\tau}}_{\text{RT}} \, dx \quad \text{for all } w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2).$$

This and the Helmholtz decomposition (40) reveal  $\tilde{\boldsymbol{\tau}}_{\text{RT}} = \operatorname{Curl}_{\text{NC}} \beta_{\text{CR}}$  for  $\beta_{\text{CR}} \in CR^1(\mathcal{T}; \mathbb{R}^2)/\mathbb{R}^2$ . For all  $\boldsymbol{\sigma}_0 \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ ,

$$0 = b((\boldsymbol{\sigma}_0, 0, 0, 0), \tilde{y}_h) = \int_{\Omega} (\boldsymbol{\sigma}_0 : D_{\text{NC}} \tilde{v}_1 + \operatorname{dev} \boldsymbol{\sigma}_0 : \tilde{\boldsymbol{\tau}}_{\text{RT}}) \, dx.$$

Hence,  $\operatorname{dev} \tilde{\boldsymbol{\tau}}_{\text{RT}} = -D_{\text{NC}} \tilde{v}_1$  and so  $v_1 \in Z_{\text{CR}}$ . This proves the asserted representation of  $N_h$ . Let  $M_h := N_h^{\perp} \subseteq Y_h$  denote the orthogonal compliment of  $N_h$  in  $Y_h$  with respect to the scalar product in  $Y$ . Then the dPG FEM is equivalent to the mixed FEM with  $x_h \in X_h$  and  $b(x_h, \bullet) = F$  in  $M_h^*$  [13, 14]. Its solvability guarantees  $\dim(M_h) = \dim(X_h) = 6|\mathcal{T}| + 2|\mathcal{N}(\Omega)| + 2|\mathcal{E}| - 1$ . This and  $\dim(Y_h) = 12|\mathcal{T}| - 1$  leads to  $\dim(N_h) = 2(|\mathcal{T}| - 1)$ . □

Given an extension  $g \in H^1(\Omega; \mathbb{R}^2)$  of the Dirichlet data  $g \in H^1(\partial\Omega; \mathbb{R}^2)$  the data approximation error contribution reads

$$\|F \circ (1 - \Pi)\|_{Y^*} = \sup_{(v, \tau) \in S(Y)} \left( \int_{\Omega} f \cdot (v - \Pi v) \, dx + \left\langle \gamma_0^T g, \gamma_v^T (1 - \Pi)\tau \right\rangle_{\partial\mathcal{T}} \right).$$

In addition, assume that  $g \in H^1(\Omega; \mathbb{R}^2) \cap H^2(\mathcal{T}; \mathbb{R}^2)$  is piecewise divergence free in that  $\Pi_0 \operatorname{div} g = 0$  in  $\Omega$ . Let  $\kappa := \sqrt{1/48 + j_{1,1}^{-2}} = 0.298234942888$  for the first root  $j_{1,1}$  of the first Bessel function and the discrete Friedrichs constant  $C_{\text{dF}}$  [7, 10.6.14].

**Theorem 5.3** *The projection  $\Pi$  from Theorem 5.2 satisfies*

$$\begin{aligned} \|F \circ (1 - \Pi)\|_{Y^*} &\leq \kappa \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \\ &\quad + \frac{h_{\max}}{j_{1,1}} \left( C_{\text{dF}} \|f\|_{L^2(\Omega)} + \left( 2 + \sqrt{1 + \kappa^2 h_{\max}^2} \|\Pi\| \right) \|g\| \right). \end{aligned}$$

*Proof* First investigate the volume contributions, i.e., the data approximation error in case  $g \equiv 0$ . The Cauchy-Schwarz and the discrete Friedrichs inequality [7, 10.6.14] with constant  $C_{\text{dF}}$  prove, for all  $y = (\tau, v) \in Y$ , that

$$\begin{aligned} \int_{\Omega} f \cdot (v - \Pi v) \, dx &= \int_{\Omega} f \cdot (v - I_{\text{NC}}^{\text{pw}} v) \, dx - \int_{\Omega} f \cdot z_{\text{CR}} \, dx \\ &\leq \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \left\| h_{\mathcal{T}}^{-1} (v - I_{\text{NC}}^{\text{pw}} v) \right\|_{L^2(\Omega)} \\ &\quad + C_{\text{dF}} \|f\|_{L^2(\Omega)} \|z_{\text{CR}}\|_{\text{NC}}. \end{aligned}$$

The first term is bounded as in [12, Thm.4] by

$$\left\| h_{\mathcal{T}}^{-1} (v - I_{\text{NC}}^{\text{pw}} v) \right\|_{L^2(\Omega)} \leq \kappa \|v - I_{\text{NC}}^{\text{pw}} v\|_{\text{NC}} \leq \kappa \|v\|_{\text{NC}}. \tag{43}$$

The choice of  $z_{\text{CR}}$  in the Helmholtz decomposition (42), (41) and the Poincaré inequality, prove for the second term

$$\begin{aligned} \|z_{\text{CR}}\|_{\text{NC}}^2 &= \int_{\Omega} D_{\text{NC}} z_{\text{CR}} : (\Pi_0 \tau - D_{\text{NC}} \alpha_{\text{CR}}) \, dx \\ &= - \int_{\Omega} z_{\text{CR}} \cdot (1 - \Pi_0) \operatorname{div}_{\text{NC}} \tau \, dx \\ &\leq h_{\max}/j_{1,1} \|z_{\text{CR}}\|_{\text{NC}} \|(1 - \Pi_0) \operatorname{div}_{\text{NC}} \tau\|_{L^2(\Omega)}. \end{aligned} \tag{44}$$

Altogether,

$$\sup_{(\tau, v) \in S(Y)} \int_{\Omega} f \cdot (v - \Pi v) \, dx \leq h_{\max} C_{\text{dF}}/j_{1,1} \|f\|_{L^2(\Omega)} + \kappa \|h_{\mathcal{T}} f\|_{L^2(\Omega)}.$$

Let  $g \in H^1(\Omega; \mathbb{R}^2) \cap H^2(\mathcal{T}; \mathbb{R}^2)$  be as above and define the nonconforming interpolant  $I_{\text{NC}g} \in \text{CR}^1(\mathcal{T}; \mathbb{R}^2)$  by  $I_{\text{NC}g}(\text{mid } E) := \int_E g \, ds$  for all  $E \in \mathcal{E}$ .

$$\begin{aligned} \left\langle \gamma_0^T g, \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} &= \left\langle \gamma_0^T (g - I_{\text{NC}}g), \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} \\ &\quad + \left\langle \gamma_0^T I_{\text{NC}}g, \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}}. \end{aligned}$$

The definition of  $\boldsymbol{\tau}_{\text{RT}}$ , (9),  $\text{div}_{\text{NC}} g = 0$ , and (41) lead to

$$\begin{aligned} \left\langle \gamma_0^T I_{\text{NC}}g, \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} &= \int_{\Omega} \mathbf{D}_{\text{NC}} I_{\text{NC}}g : (\Pi_0 \boldsymbol{\tau} - \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} - \text{Curl } \beta_c) \, dx \\ &\quad + \int_{\Omega} I_{\text{NC}}g \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx \\ &= - \int_{\Omega} \mathbf{D}_{\text{NC}} I_{\text{NC}}g : \text{Curl } \beta_c \, dx. \end{aligned}$$

Equation (42), the Cauchy-Schwarz inequality, and (41) lead to

$$\begin{aligned} &- \int_{\Omega} \mathbf{D}_{\text{NC}} I_{\text{NC}}g : \text{Curl } \beta_c \, dx \\ &= \int_{\Omega} \mathbf{D}_{\text{NC}} I_{\text{NC}}g : (\mathbf{D}_{\text{NC}} z_{\text{CR}} + \text{dev}(\mathbf{D}_{\text{NC}} \alpha_{\text{CR}} - \Pi_0 \boldsymbol{\tau})) \, dx \\ &\leq \|I_{\text{NC}}g\|_{\text{NC}} \|z_{\text{CR}}\|_{\text{NC}} + \int_{\Omega} \mathbf{D}_{\text{NC}} I_{\text{NC}}g : (\mathbf{D}_{\text{NC}} \alpha_{\text{CR}} - \Pi_0 \boldsymbol{\tau}) \, dx \\ &= \|I_{\text{NC}}g\|_{\text{NC}} \|z_{\text{CR}}\|_{\text{NC}} + \int_{\Omega} h_{\mathcal{T}}^{-1} (I_{\text{NC}}g - \Pi_0 I_{\text{NC}}g) \cdot h_{\mathcal{T}}^{+1} (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx. \end{aligned}$$

The application of the Poincaré inequality and (44) prove

$$\begin{aligned} \sup_{(\boldsymbol{\tau}, v) \in S(Y)} \left\langle \gamma_0^T I_{\text{NC}}g, \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} &\leq \frac{2h_{\max}}{j_{1,1}} \|(1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)} \|I_{\text{NC}}g\|_{\text{NC}} \\ &\leq \frac{2h_{\max}}{j_{1,1}} \|g\|. \end{aligned}$$

Finally, for all  $(\boldsymbol{\tau}, v) \in S(Y_h)$ , it holds

$$\begin{aligned} \left\langle \gamma_0^T (g - I_{\text{NC}}g), \gamma_v^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} &\leq \|g - I_{\text{NC}}g\|_{H^1(\mathcal{T})} \|\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}\|_{H(\text{div}, \mathcal{T})} \\ &\leq \sqrt{1 + \kappa^2 h_{\max}^2} \|g - I_{\text{NC}}g\|_{\text{NC}} \|1 - \Pi\| \|\boldsymbol{\tau}\|_{H(\text{div}, \mathcal{T})} \\ &\leq \sqrt{1 + \kappa^2 h_{\max}^2} \|\Pi\| h_{\max} / j_{1,1} \|g\|. \end{aligned}$$

The equality  $\|1 - \Pi\| = \|\Pi\|$  follows from Kato’s lemma [23, Lemma 4]. □

In conclusion, the data approximation term is not necessarily of higher-order, but at least controlled by  $h_{\max}$  even for non homogeneous boundary data.



### 5.2 Extensions of test spaces

The discrete inf-sup condition (4) also holds for an enlarged discrete test search space. Three examples for an enlarged test search space  $\hat{Y}_h = Y_{h,j}$  for  $j = 1, 2, 3$  allow for a decoupling of the Fortin interpolation operator  $\Pi$  and higher-order data approximation error. Here,  $\mathcal{B}_3(\mathcal{T}) := \{v \in P_3(\mathcal{T}) : v = 0 \text{ along } \partial\mathcal{T}\}$  denotes the cubic bubble functions, and

- (i)  $Y_{h,1} := Y_h \oplus (\mathcal{B}_3(\mathcal{T})\mathbb{R}_{\text{dev}}^{2 \times 2} \times \{0\})$ ,
- (ii)  $Y_{h,2} := P_1(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^2)$ ,
- (iii)  $Y_{h,3} := RT_1^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^2)$ .

Given  $\tau \in H(\text{div}, \mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ , there exists  $(\hat{\tau}_{\text{RT}}, 0) \in \hat{Y}_h := Y_{h,j}$  for  $j = 1, 2, 3$  with

$$\Pi_0 \text{div } \hat{\tau}_{\text{RT}} = \Pi_0 \text{div } \tau, \tag{45}$$

$$\langle (\hat{\tau}_{\text{RT}} - \tau)v, w_c \rangle_{\partial T} = 0 \quad \text{for all } T \in \mathcal{T} \text{ and } w_c \in P_1(T; \mathbb{R}^2). \tag{46}$$

In particular (46) implies  $\Pi_0 \text{div } \hat{\tau}_{\text{RT}} = \Pi_0 \text{div } \tau$ . Altogether, the definition  $\hat{\Pi}(\tau, v) := (\hat{\tau}_{\text{RT}}, I_{\text{NC}}^{\text{pw}} v)$  with (45)–(46) guarantees  $b(x_h, (1 - \hat{\Pi})v) = 0$  for all  $x_h \in X_h$ . In case  $\hat{Y}_h = Y_{h,2}$  and  $\hat{Y}_h = Y_{h,3}$ , (45)–(46) allow multiple choices of  $\hat{\tau}_{\text{RT}}$ .

**Lemma 5.4** *In case  $\hat{Y}_h = Y_{h,1}$ ,  $\hat{\Pi}(y) = (\hat{\tau}_{\text{RT}}, I_{\text{NC}}^{\text{pw}} v) \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^2)/R \oplus \mathcal{B}_3(\mathcal{T})\mathbb{R}_{\text{dev}}^{2 \times 2} \times P_1(\mathcal{T}; \mathbb{R}^2)$  is unique and defines a projection. A bound of  $\|\hat{\Pi}\| \leq (2 + 15.5\sqrt{\cot(\alpha_{\text{min}})}h_{\text{max}} + (3.22 + 60 \cot(\alpha_{\text{min}}))h_{\text{max}}^2)^{1/2}$  depends on  $h_{\text{max}}$  and the smallest angle of the triangulation  $\alpha_{\text{min}}$ .*

*Proof* The proof of Lemma 5.4 is given in the appendix. □

The representation  $\hat{\Pi}(\tau, v) := (\hat{\tau}_{\text{RT}}, I_{\text{NC}}^{\text{pw}} v)$  of the operator  $\hat{\Pi}$  and (43) prove

$$\sup_{(v, \tau) \in \mathcal{S}(Y)} \int_{\Omega} f \cdot (v - I_{\text{NC}}^{\text{pw}} v) \, dx \leq \kappa \|h_{\mathcal{T}} f\|_{L^2(\Omega)}. \tag{47}$$

In case of inhomogeneous boundary data, let  $g \in H^1(\Omega, \mathbb{R}^2)$  be an extension of  $g \in H^1(\partial\Omega; \mathbb{R}^2)$  with  $g|_E \in P_1(E)$  for all  $E \in \mathcal{E}(\Omega)$ . Let  $Ig \in S_1(\mathcal{E}; \mathbb{R}^n)$  denote the conforming interpolation defined by linear interpolation of the nodal values,  $Ig(z) = g(z)$  for all  $z \in \mathcal{N}$ . Hence,  $g$  and  $Ig$  coincide along any interior edge  $E \in \mathcal{E}(\Omega)$ . This choice and (46) lead to

$$\left\langle \gamma_0^{\mathcal{T}} g, \gamma_v^{\mathcal{T}} (\tau - \hat{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} = \langle \gamma_0(g - Ig), (\tau - \hat{\tau}_{\text{RT}})v \rangle_{\partial\Omega}.$$

Let  $g' := \partial g / \partial s$  denote the arc-length derivative of  $g \in H^1(\partial\Omega; \mathbb{R}^2)$  along the boundary,  $\Pi_0^{\mathcal{E}} g'$  the  $L^2(\partial\Omega)$ -orthogonal projection of  $g'$  onto  $P_0(\mathcal{E}(\partial\Omega); \mathbb{R}^2)$ , and

$h_{\mathcal{E}} \in P_0(\mathcal{E})$  the piecewise constant function with  $h_{\mathcal{E}}|_E = \text{diam}(\omega_E) = \text{diam}(T_+ \cup T_-)$  for every  $E \in \mathcal{E}$  as in Fig. 1. This allows to define the Dirichlet data oscillation

$$\text{osc}(g', \mathcal{E}(\partial\Omega)) := \|h_{\mathcal{E}}^{1/2}(1 - \Pi_0^{\mathcal{E}})g'\|_{L^2(\partial\Omega)}.$$

According to [8, Proof of Lemma 2.1], [5] there exists  $w \in H^1(\Omega; \mathbb{R}^2)$  with  $w|_{\partial\Omega} = (1 - I)g|_{\partial\Omega}$  and  $\|w\|_{H^1(\Omega)} \lesssim \text{osc}(g', \mathcal{E}(\partial\Omega))$ . Hence,

$$\begin{aligned} \left\langle \gamma_0 g, \gamma_v^T(\boldsymbol{\tau} - \hat{\boldsymbol{\tau}}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} &= \left\langle \gamma_0 w, (\boldsymbol{\tau} - \hat{\boldsymbol{\tau}}_{\text{RT}}) \nu \right\rangle_{\partial\Omega} \\ &= \int_{\Omega} w \cdot \text{div}_{\text{NC}}(\boldsymbol{\tau} - \hat{\boldsymbol{\tau}}_{\text{RT}}) \, dx + \int_{\Omega} \text{D} w : (\boldsymbol{\tau} - \hat{\boldsymbol{\tau}}_{\text{RT}}) \, dx \\ &\leq \|w\|_{H^1(\Omega)} \|(1 - \hat{\Pi})\boldsymbol{\tau}\|_{H(\text{div}, \mathcal{T})}. \end{aligned} \tag{48}$$

Therefore, for each  $\hat{Y}_h := Y_{h,j}$  for  $j = 1, 2, 3$ , it follows

$$\sup_{(v, \boldsymbol{\tau}) \in S(Y)} \left\langle \gamma_0^T g, \gamma_v^T(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \right\rangle_{\partial\mathcal{T}} \lesssim \|1 - \hat{\Pi}\| \text{osc}(g', \mathcal{E}(\partial\Omega)).$$

Hence, a slight enlargement of the test search space guarantees an higher-order data approximation error independent of the given Dirichlet data  $g$ . For  $g|_E \in H^2(E)$  for all  $E \in \mathcal{E}(\partial\Omega)$  with edgewise second surface derivative  $\partial_{\mathcal{E}}^2 g / \partial s^2$  a better estimate with explicit constants is possible. There exists  $w \in H^1(\Omega; \mathbb{R}^2)$ , such that  $w|_{\partial\Omega} = (1 - I)g$  and

$$\|w\|_{H^1(\Omega)} \leq \sqrt{c_1^2 + h_{\max, \partial\Omega}^2 c_2^2} \|h_{\mathcal{E}}^{3/2} \partial_{\mathcal{E}}^2 g / \partial s^2\|_{L^2(\Omega)}. \tag{49}$$

The constants are computed in [15, Thm.5.1] and [24, Thm.4.2.2]. They depend only on the shape of the triangles of  $\mathcal{T}$  not on the mesh-size, e.g., for right isosceles triangles  $c_1 \leq 0.4980$  and  $c_2 \leq 0.0654$ .

**Lemma 5.5** *Let  $\mathcal{T}$  consist of right isosceles triangles and  $g|_E \in H^2(E)$  for all  $E \in \mathcal{E}(\partial\Omega)$ . The data approximation error  $\|F \circ (1 - \Pi)\|_{Y^*}$  is explicitly bounded from above by*

$$\begin{aligned} &0.3 \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \\ &+ \sqrt{0.5 + 3.84 h_{\max} + 15.9 h_{\max}^2 + 0.07 h_{\max}^3 + 0.27 h_{\max}^4} \|h_{\mathcal{E}}^{3/2} \partial_{\mathcal{E}}^2 g / \partial s^2\|_{L^2(\partial\Omega)}. \end{aligned}$$

*Proof* This follows directly from (47)–(49), Kato’s Lemma and Lemma 5.4. □

## 6 Numerical examples

Three benchmark examples concern uniform and adaptive mesh-refinement with various choices of the input bulk parameter  $\theta$  in the adaptive algorithm displayed in the convergence history plots.

---

**Algorithm 1: AFEM**

---

**input:** regular initial triangulation  $\mathcal{T}_0$  and bulk parameter  $0 < \theta \leq 1$   
**for**  $\ell = 0, 1, 2, \dots$  **do**  
    **solve** Compute solution  $x_\ell$  for discrete problem (3) on  $\mathcal{T}_\ell$   
    **estimate** Compute local contributions  $\eta_\ell^2(T)$  for all  $T \in \mathcal{T}_\ell$  and global error estimator  
     $\eta_\ell^2 = \sum_{T \in \mathcal{T}} \eta_\ell^2(T)$   
    **mark** Choose  $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell$  of minimal cardinality with  $\theta \sum_{T \in \mathcal{T}_\ell} \eta_\ell^2(T) \leq \sum_{T \in \mathcal{M}_\ell} \eta_\ell^2(T)$   
    **refine** Generate minimal refinement  $\mathcal{T}_{\ell+1}$  of  $\mathcal{T}_\ell$  with  $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell \setminus \mathcal{T}_{\ell+1}$   
**output:** Series of triangulations  $\mathcal{T}_\ell$ , discrete solutions  $x_\ell$ , and error estimators  $\eta_\ell$

---

**6.1 Numerical realisation**

The implementation has been performed straightforwardly into Matlab and extends the data structures of [1]. The adaptive finite element mesh-refining runs Algorithm 1 with  $\eta_\ell^2 := \|F - b(x_\ell, \bullet)\|_{Y_\ell^*}^2$  for the discrete test search space  $Y_\ell$  on level  $\ell$  and the associated discrete solution  $x_\ell$ . Let  $Y_\ell(T) \subset Y_\ell$  denote the set of all basis functions with support  $T \in \mathcal{T}$  and use  $\eta_\ell^2(T) := \|F - b(x_\ell, \bullet)\|_{Y_\ell(T)^*}^2$  as a refinement indicator. The extended test space  $\hat{Y}_\ell := Y_\ell \oplus (\mathcal{B}_3(\mathcal{T})\mathbb{R}_{\text{dev}}^{2 \times 2} \times \{0\})$  from Sect. 5.2 leads to the error estimator  $\hat{\eta}_\ell^2 := \|F - b(x_\ell, \bullet)\|_{\hat{Y}_\ell^*}^2$  and  $\tilde{\eta}_\ell^2 := \hat{\eta}_\ell^2 + \text{osc}^2(g', \mathcal{E}(\partial\Omega))$ . The oscillations are computed with the exact derivatives of  $g \in H^1(\partial\Omega; \mathbb{R}^2)$  and numerical integration with 7 Gauss points per edge. In the examples  $f \equiv 0$ , so that  $\hat{\eta}_\ell$  is a guaranteed error estimator upto a multiplicative generic constant. Instead of the exact error  $\|x - x_\ell\|_X$  an upper bound is computed and displayed via the unique extensions  $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$  (resp.  $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^2)$ ) of  $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^2)$  (resp.  $t_0 \in P_0(\mathcal{E}; \mathbb{R}^2)$ ) as in (34),

$$\begin{aligned} \|x - x_\ell\|_X^2 &\leq \|u - u_0\|_{L^2(\Omega)}^2 + \|\sigma - \sigma_0\|_{L^2(\Omega)}^2 \\ &\quad + \|u - w_c\|_{H^1(\Omega)}^2 + \|\sigma - \mathbf{q}_{\text{RT}}\|_{H(\text{div}; \Omega)}^2. \end{aligned}$$

**6.2 Colliding flow example**

In this benchmark problem  $f \equiv 0$  in  $\Omega = (-1, 1)^2$  with given boundary data from the exact solution  $(u, p)$  with, for all  $(x_1, x_2) \in \Omega$ ,

$$\begin{aligned} u(x_1, x_2) &= 4 \left( 5x_1x_2^4 - x_1^5, 5x_1^4x_2 - x_2^5 \right), \\ p(x_1, x_2) &= 120x_1^2x_2^2 - 20(x_1^4 + x_2^4) - 16/3. \end{aligned}$$

Figure 2 presents the computed error estimator and upper bound for the exact error for uniform refinement. The estimator converges with the optimal rate 0.5 for uniform red-refinement. The exact error is dominated by the error in the pseudostress component.

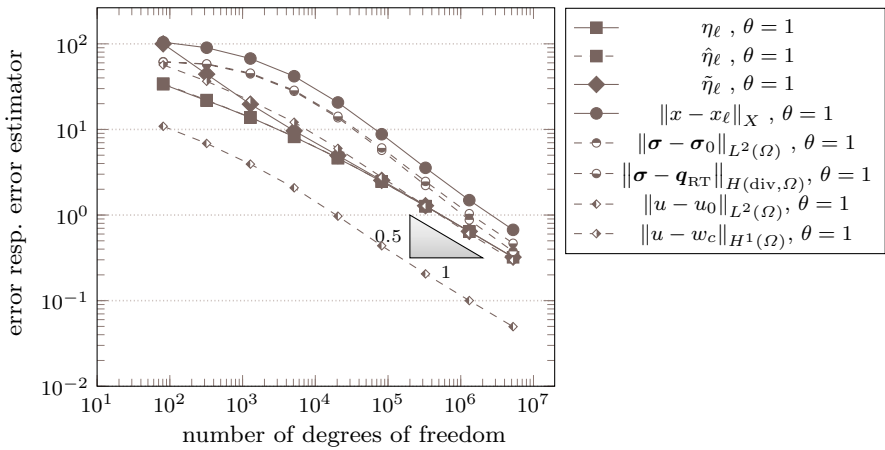


Fig. 2 Convergence history plot for uniform red-refinement for the colliding flow example

### 6.3 Example on L-shaped domain

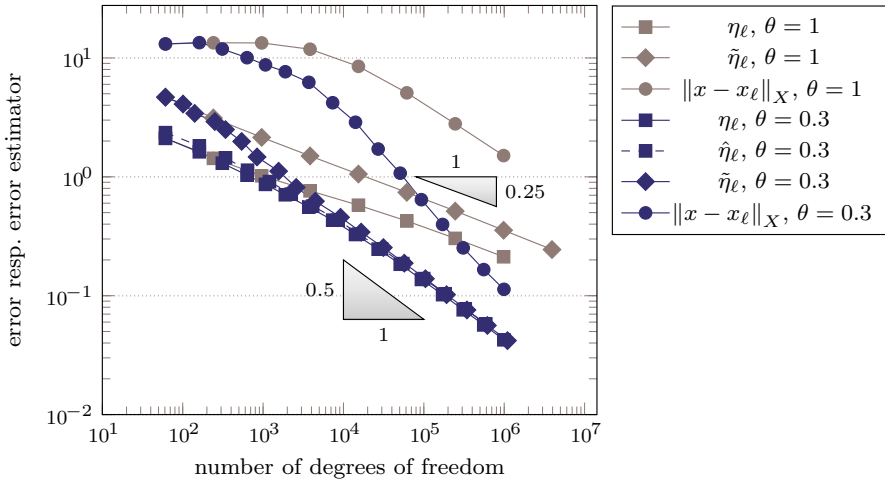
In this example  $f \equiv 0$  on the L-shaped domain,  $\Omega = (-1, 1)^2 \setminus ([0, 1] \times [-1, 0])$ , with  $\omega = 3\pi/2$ ,  $\alpha = 856399/1572864$ , and

$$w(\varphi) = \frac{\sin((1 + \alpha)\varphi) \cos(\alpha\omega)}{1 + \alpha} - \cos((1 + \alpha)\varphi) + \frac{\sin((\alpha - 1)\varphi) \cos(\alpha\omega)}{1 - \alpha} + \cos((\alpha - 1)\varphi).$$

The exact solution from [26, p.324] reads, in polar coordinates for the implicit boundary data, for all  $(r, \varphi) \in [0, \infty) \times [0, 3\pi/2]$ ,

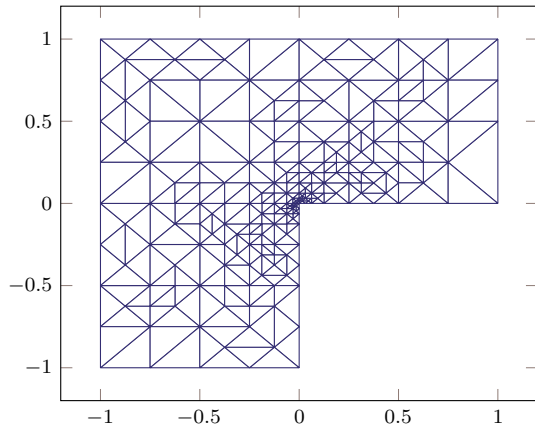
$$u(r, \varphi) = r^\alpha \left( (1 + \alpha) \sin(\varphi) w(\varphi) + \cos(\varphi) w'(\varphi), \right. \\ \left. - (1 + \alpha) \cos(\varphi) w(\varphi) + \sin(\varphi) w'(\varphi) \right), \\ p(r, \varphi) = -r^{\alpha-1} \left( (1 + \alpha)^2 w'(\varphi) + w'''(\varphi) \right) / (1 - \alpha).$$

Figure 3 shows the convergence history plot with an adaptive refinement strategy of optimal empirical convergence rate 0.5. In case of uniform refinement, as expected in a non-convex domain with singularity in the reentrant corner, the empirical convergence rate is 0.25. The computed error shows some pre-asymptotic range, which is typical also for other finite element discretizations (not displayed). The error estimator  $\tilde{\eta}_\ell$ , which includes the boundary oscillation, follows accordingly. The adaptive algorithm resolves the singularity in the reentrant corner first as depicted in Fig. 4.



**Fig. 3** Convergence history plot for uniform and adaptive refinement with  $\theta = 0.3$  for the example on the L-shaped domain

**Fig. 4** Triangulation  $\mathcal{T}_\ell$  with 3711 degrees of freedom (371 elements) for the example on the L-shaped domain from adaptive refinement with  $\theta = 0.3$

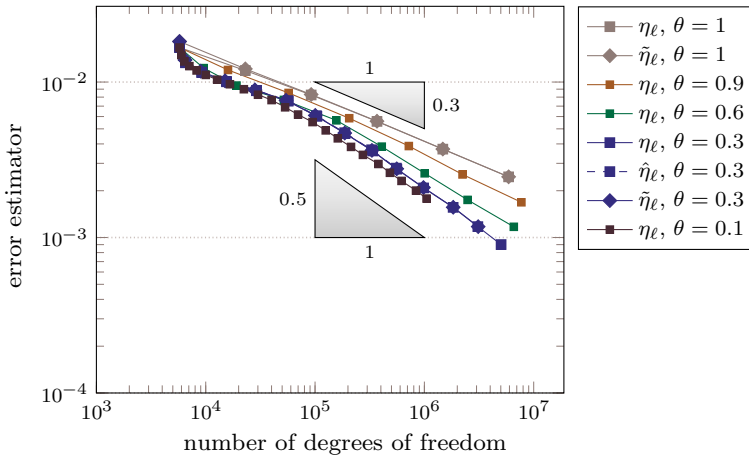


### 6.4 Backward facing step example

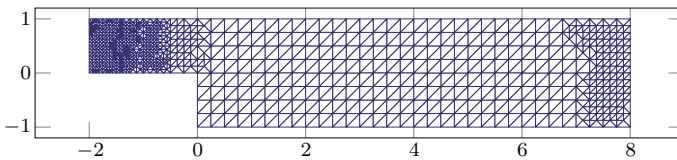
This benchmark example with  $f \equiv 0$  on a slightly deformed L-shaped domain  $\Omega = ((-2, 8) \times (-1, 1)) \setminus ((-2, 0) \times (-1, 0))$  has the Dirichlet data  $g$  for  $(x_1, x_2) \in \partial\Omega$

$$g(x_1, x_2) = \begin{cases} 1/10(-x_2(x_2 - 1), 0) & \text{for } x_1 = -2, \\ 1/80(-(x_2 - 1)(x_2 + 1), 0) & \text{for } x_1 = 8, \\ (0, 0)^\top & \text{else.} \end{cases}$$

Figure 5 presents the error estimator for varying bulk parameter  $\theta$ . Obviously, a smaller  $\theta$  leads to a better convergence rate. On the other hand, more levels are needed to



**Fig. 5** Convergence history plot with varying bulk parameter  $\theta$  for the backward facing step example



**Fig. 6** Triangulation  $\mathcal{T}_\ell$  with 15511 degrees of freedom (1551 elements) for the backward facing step example from adaptive refinement with  $\theta = 0.3$

reach a certain number of degrees of freedom and  $\theta = 0.3$  leads to the optimal empirical convergence rate. The choice of different error estimators  $\eta_\ell, \hat{\eta}_\ell, \tilde{\eta}_\ell$  does not influence the result. The inhomogeneous Dirichlet boundary conditions are resolved in the triangulation in Fig. 6 before the singularity at the reentrant corner becomes significant.

### 6.5 Conclusion

All the numerical experiments confirm the theoretical results and support the conjectured instant stability of the dPG paradigm: The systematic convergence with a clear empirical convergence rate is visible from the very beginning even for the coarsest meshes. The extensions of the test search space do not affect the approximation of the discrete solution significantly. It is not rewarding to compute with bigger test search spaces. The error estimators  $\eta_\ell$  and  $\hat{\eta}_\ell$  are almost identical though  $\tilde{\eta}_\ell$  leads to a guaranteed error bound. It is utterly an empirical observation that the associated adaptive mesh-refining algorithm improves suboptimal convergence rates in case of singular solutions.

**Acknowledgements** This work has been supported by the Deutsche Forschungsgemeinschaft (DFG) in the Priority Program 1748 ‘Reliable simulation techniques in solid mechanics. Development of non-standard

discretization methods, mechanical and mathematical analysis' under the project CA 151/22-1. The second author is supported by the Berlin Mathematical School.

### Appendix: Fortin interpolation in an extended test search space

*Proof of Lemma 5.4* For  $\hat{Y}_h := (RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^2)/\mathbb{R} \oplus \mathcal{B}_3(\mathcal{T})\mathbb{R}_{\text{dev}}^{2 \times 2}) \times P_1(\mathcal{T}; \mathbb{R}^2)$  the Fortin interpolator  $\hat{\Pi} : Y \rightarrow \hat{Y}_h$  maps  $(\boldsymbol{\tau}, v) \mapsto (\hat{\boldsymbol{\tau}}_{\text{RT}}, I_{\text{NC}}^{\text{pw}} v)$  such that (45)–(46) hold. Recall the edge-based Raviart–Thomas basis  $\psi_E$  for  $E \in \mathcal{E}(T)$ , the opposite vertex  $P_E \in \mathcal{N}(T)$ , and the component  $\kappa = 1, 2$ , let  $\Psi_{E,\kappa} := e_\kappa \otimes \psi_E := e_\kappa \otimes \chi(T)(x - P_E)|E|/(2|T|)$ . The Crouzeix–Raviart basis functions in two components  $\Phi_{E,\kappa}$  read  $\Phi_{E,\kappa} := \phi_E e_\kappa := \chi(T)(1 - 2\varphi_{P_E}) e_\kappa$  with nodal basis functions  $\varphi_z$  of  $z \in \mathcal{N}$ . Given  $T \in \mathcal{T}$ ,  $\boldsymbol{\tau} \in H(\text{div}, T; \mathbb{R}^{2 \times 2})$ , define

$$I_F \boldsymbol{\tau} := \sum_{\kappa=1,2} \sum_{E \in \mathcal{E}(T)} \left( \frac{1}{|E|} \int_{\partial T} \Phi_{E,\kappa} \cdot \boldsymbol{\tau} \nu_T \, ds \, \Psi_{E,\kappa} \right). \tag{50}$$

Since  $(\psi_E \cdot \nu_T)|_F = \delta_{EF}$  and  $\int_F \phi_E \, ds = \delta_{EF}|E|$  for  $E, F \in \mathcal{E}(T)$ ,  $I_F = I_F^2$  is a projection. Moreover,  $\langle I_F \boldsymbol{\tau} \nu_T, \Phi_{E,\kappa} \rangle_{\partial T} = \langle \boldsymbol{\tau} \nu_T, \Phi_{E,\kappa} \rangle_{\partial T}$  for any  $E \in \mathcal{E}(T)$  and  $\kappa = 1, 2$  implies (46). An integration by parts allows to rewrite  $I_F \boldsymbol{\tau}$ . The projection property implies  $I_F q = q$  for all  $q \in P_0(T)$ . This and  $\sum_{E \in \mathcal{E}(T)} \varphi_E = 1$  reveal

$$\begin{aligned} I_F \boldsymbol{\tau} &= \sum_{\kappa=1,2} \sum_{E \in \mathcal{E}(T)} \frac{1}{|E|} \left( \int_T \text{D} \Phi_{E,\kappa} : \boldsymbol{\tau} \, dx + \int_T \Phi_{E,\kappa} \cdot \text{div} \boldsymbol{\tau} \, dx \right) \Psi_{E,\kappa} \\ &= \Pi_0 \boldsymbol{\tau} + \sum_{\kappa=1,2} \sum_{E \in \mathcal{E}(T)} \frac{1}{2} \int_T \Phi_{E,\kappa} \cdot \text{div} \boldsymbol{\tau} \, dx \, e_\kappa \otimes (x - P_E) \\ &= \Pi_0 \boldsymbol{\tau} + \frac{\Pi_0 \text{div} \boldsymbol{\tau}}{2} \otimes (x - \text{mid}(T)) + Q(T), \end{aligned}$$

where  $Q(T) := \sum_{\kappa=1,2} \sum_{E \in \mathcal{E}(T)} \frac{1}{2} \left( \int_T \Phi_{E,\kappa} \cdot \text{div} \boldsymbol{\tau} \, dx \right) e_\kappa \otimes (\text{mid}(T) - P_E)$ . Given  $b_T := 60\varphi_1\varphi_2\varphi_3 \in \mathcal{B}_3(T)$ , for  $T \in \mathcal{T}$ , set

$$\Pi_\tau \boldsymbol{\tau} := I_F \boldsymbol{\tau} + b_T \text{dev} \Pi_0(\boldsymbol{\tau} - I_F \boldsymbol{\tau}) = I_F \boldsymbol{\tau} - b_T \text{dev} Q(T). \tag{51}$$

Since  $b_T|_{\partial T} = 0$  and  $\int_T b_T = 1$ , this operator  $\Pi_\tau$  satisfies (45)–(46). To compute  $\|\Pi_\tau\|$  with

$$\begin{aligned} \|\Pi_\tau(\boldsymbol{\tau})\|_{H(\text{div}, T)}^2 &= \|I_F \boldsymbol{\tau} - b_T \text{dev} Q(T)\|_{L^2(T)}^2 \\ &\quad + \|\text{div} I_F \boldsymbol{\tau} - \text{dev} Q(T)\|_{L^2(T)}^2, \end{aligned}$$

estimate the Frobenius norm  $\|Q(T)\|_F$  of  $Q(T)$  as follows. Since  $|\text{mid}(T) - P_E| \leq 2h_T/3$  and  $\|\phi_E\|_{L^2(T)}^2 = |T|/3$ , the Cauchy-Schwarz inequality allows

$$\begin{aligned} \|Q(T)\|_F^2 &= \sum_{\kappa=1,2} \sum_{k=1,2} \left( \sum_{E \in \mathcal{E}(T)} 1/2 \int_T \Phi_{E,\kappa} \cdot \operatorname{div} \boldsymbol{\tau} \, dx \, (\operatorname{mid}(T) - P_E) \cdot e_k \right)^2 \\ &\leq 2 \sum_{\kappa=1,2} \left( \sum_{E \in \mathcal{E}(T)} 1/(2|T|) \int_T |\Phi_{E,\kappa} \cdot \operatorname{div} \boldsymbol{\tau}| \, dx \, |\operatorname{mid}(T) - P_E| \right)^2 \\ &\leq 2h_T^2/9 \sum_{\kappa=1,2} \left( \|\operatorname{div} \boldsymbol{\tau} \cdot e_\kappa\|_{L^2(T)} \sum_{E \in \mathcal{E}(T)} |T|^{-1} \|\phi_E\|_{L^2(T)} \right)^2 \\ &\leq 2h_T^2/(3|T|) \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)}^2. \end{aligned}$$

Moreover,  $\|b_T\|_{L^2(T)}^2 = 10|T|/7$  and  $\|\nabla b_T\|_{L^2(T)}^2 = 20|T| \|G(T)\|_F^2$ , with

$$G(T) := \begin{pmatrix} \nabla \varphi_1^\top \\ \nabla \varphi_2^\top \\ \nabla \varphi_3^\top \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}$$

and  $\|G(T)\|_F^2 = (2|T|)^{-2} \sum_{E \in \mathcal{E}(T)} |E|^2 \leq 3h_T^2(2|T|)^{-2}$ . Let  $\alpha_T$  denote the smallest angle in  $T$  and recall  $h_T^2 \leq 4|T| \cot(\alpha_T)$ . Hence,

$$\begin{aligned} \|b_T \operatorname{dev} Q(T)\|_{L^2(T)}^2 &\leq \|Q(T)\|_F^2 \|b_T\|_{L^2(T)}^2 = 20h_T^2/21 \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)}^2, \\ \|\operatorname{dev} Q(T) \nabla b_T\|_{L^2(T)}^2 &\leq \|Q(T)\|_F^2 \|\nabla b_T\|_{L^2(T)}^2 = 15h_T^4/|T| \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)}^2 \\ &\leq 60h_T^2 \cot(\alpha_T) \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)}^2. \end{aligned}$$

Alltogether,

$$\begin{aligned} \|\Pi_\tau(\boldsymbol{\tau})\|_{H(\operatorname{div},T)}^2 &\leq \left( \|\Pi_0 \boldsymbol{\tau}\|_{L^2(T)} + \|Q(T)\|_{L^2(T)} + \|b_T \operatorname{dev} Q(T)\|_{L^2(T)} \right)^2 \\ &\quad + \left( \|\operatorname{div} \Pi_0 \boldsymbol{\tau}\|_{L^2(T)} + \|\operatorname{dev} Q(T) \nabla b_T\|_{L^2(T)} \right)^2 \\ &\leq \left( \|\boldsymbol{\tau}\|_{L^2(T)} + \left( \sqrt{2/3} + \sqrt{20/21} \right) h_T \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)} \right)^2 \\ &\quad + \left( 1 + 2\sqrt{15 \cot(\alpha_T)} h_T \right)^2 \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(T)}^2 \\ &\leq \left( 1 + 3.22h_T^2 + (1 + 7.75\sqrt{\cot(\alpha)} h_T)^2 \right) \|\boldsymbol{\tau}\|_{H(\operatorname{div},T)}^2 \\ &\leq (2 + 15.5\sqrt{\cot(\alpha_T)} h_T + (3.22 + 60 \cot(\alpha_T)) h_T^2) \|\boldsymbol{\tau}\|_{H(\operatorname{div},T)}^2. \end{aligned}$$

Set  $\hat{\Pi}(\boldsymbol{\tau}, v)|_T := (\Pi_\tau(\boldsymbol{\tau}|_T), I_{\text{nc}}^{\text{pw}} v|_T)$  as in (51). Given  $y = (\boldsymbol{\tau}, v) \in Y$ , the above computation, the abbreviation  $\alpha_{\min}$  for the smallest angle of the triangulation (which is bounded in a regular triangulation) and [12, Thm. 4] prove



$$\begin{aligned}
\|\hat{\Pi}(y)\|_Y^2 &= \|\Pi_\tau(\boldsymbol{\tau})\|_{H(\operatorname{div}, \mathcal{T})}^2 + \|I_{\text{NC}}^{\text{pw}} v\|_{H^1(\mathcal{T})}^2 \\
&\leq (2 + 15.5\sqrt{\cot(\alpha_{\min})}h_{\max} + (3.22 + 60\cot(\alpha_{\min}))h_{\max}^2) \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T})}^2 \\
&\quad + (2 + 2\kappa^2 h_{\max}^2) \|v\|_{H^1(\mathcal{T})}^2 \\
&\leq (2 + 15.5\sqrt{\cot(\alpha_{\min})}h_{\max} + (3.22 + 60\cot(\alpha_{\min}))h_{\max}^2) \|y\|_Y^2.
\end{aligned}$$

□

## References

1. Albery, J., Carstensen, C., Funken, S.: Remarks around 50 lines of Matlab: short finite element implementation. *Numer. Algor.* **20**(2–3), 117–137 (1999)
2. Alt, H.W.: *Lineare Funktionalanalysis: Eine anwendungsorientierte Einführung*, 5. Auflg (2006)
3. Arnold, D.N., Falk, R.S.: A uniformly accurate finite element method for the Reissner-Mindlin plate. *SIAM J. Numer. Anal.* **26**(6), 1276–1290 (1989)
4. Babuška, I.: Error-bounds for finite element method. *Numer. Math.* **16**, 322–333 (1970/1971)
5. Bartels, S., Carstensen, C., Dolzmann, G.: Inhomogeneous Dirichlet conditions in a priori and a posteriori finite element error analysis. *Numer. Math.* **99**(1), 1–24 (2004)
6. Boffi, D., Brezzi, F., Fortin, M.: *Mixed Finite Element Methods and Applications*. Springer Series in Computational Mathematics. Springer, Berlin, Heidelberg (2013)
7. Brenner, S., Scott, R.: *The Mathematical Theory of Finite Element Methods*. Texts in Applied Mathematics. Springer, Berlin (2008)
8. Bringmann, P., Carstensen, C.: An adaptive least-squares FEM for the Stokes equations with optimal convergence rates. *Numer. Math.* **135**(2), 459–492 (2017)
9. Cai, Z., Tong, C., Vassilevski, P.S., Wang, C.: Mixed finite element methods for incompressible flow: stationary Stokes equations. *Numer. Methods Partial Differ. Equ.* **26**(4), 957–978 (2010)
10. Carstensen, C., Demkowicz, L., Gopalakrishnan, J.: A posteriori error control for DPG methods. *SIAM J. Numer. Anal.* **52**(3), 1335–1353 (2014)
11. Carstensen, C., Demkowicz, L., Gopalakrishnan, J.: Breaking spaces and forms for the DPG method and applications including Maxwell equations. *Comput. Math. Appl.* **72**(3), 494–522 (2016)
12. Carstensen, C., Gallistl, D.: Guaranteed lower eigenvalue bounds for the biharmonic equation. *Numer. Math.* **126**(1), 33–51 (2014)
13. Carstensen, C., Gallistl, D., Hellwig, F., Weggler, L.: Low-order dPG-FEM for an elliptic PDE. *Comput. Math. Appl.* **68**(11), 1503–1512 (2014)
14. Carstensen, C., Hellwig, F.: Low-order discontinuous Petrov–Galerkin finite element methods for linear elasticity. *SIAM J. Numer. Math.* **54**(6), 3388–3410 (2016)
15. Carstensen, C., Merdon, C.: Computational survey on a posteriori error estimators for nonconforming finite element methods for the Poisson problem. *J. Comput. Appl. Math.* **249**, 74–94 (2013)
16. Carstensen, C., Merdon, C.: Computational survey on a posteriori error estimators for the Crouzeix-Raviart nonconforming finite element method for the Stokes problem. *Comput. Methods Appl. Math.* **14**(1), 35–54 (2014)
17. Carstensen, C., Peterseim, D., Rabus, H.: Optimal adaptive nonconforming FEM for the Stokes problem. *Numer. Math.* **123**(2), 291–308 (2013)
18. Demkowicz, L., Gopalakrishnan, J.: A class of discontinuous Petrov–Galerkin methods. Part I: The transport equation. *Comput. Methods Appl. Mech. Eng.* **199**(23), 1558–1572 (2010)
19. Demkowicz, L., Gopalakrishnan, J.: Analysis of the DPG method for the Poisson equation. *SIAM J. Numer. Anal.* **49**(5), 1788–1809 (2011)
20. Demkowicz, L., Gopalakrishnan, J.: A class of discontinuous Petrov–Galerkin methods. Part II. Optimal test functions. *Numer. Methods Partial Differ. Equ.* **27**(1), 70–105 (2011)
21. Girault, V., Raviart, P.A.: *Finite element methods for Navier–Stokes equations: theory and algorithms*. Springer, Berlin (1986)
22. Gopalakrishnan, J., Qiu, W.: An analysis of the practical DPG method. *Math. Comput.* **83**(286), 537–552 (2014)

23. Kato, T.: Estimation of iterated matrices, with application to the von Neumann condition. *Numer. Math.* **2**(1), 22–29 (1960)
24. Merdon, C.: Aspects of guaranteed error control in computations for partial differential equations. Ph.D. thesis, Humboldt-Universität zu Berlin (2013)
25. Roberts, N.V., Bui-Thanh, T., Demkowicz, L.: The DPG method for the Stokes problem. *Comput. Math. Appl.* **67**(4), 966–995 (2014). High-order Finite Element Approximation for Partial Differential Equations
26. Verfürth, R.: A posteriori error estimators for the Stokes equations. *Numer. Math.* **55**(3), 309–325 (1989)