

## A CUT FINITE ELEMENT METHOD WITH BOUNDARY VALUE CORRECTION

ERIK BURMAN, PETER HANSBO, AND MATS G. LARSON

**ABSTRACT.** In this contribution we develop a cut finite element method with boundary value correction of the type originally proposed by Bramble, Dupont, and Thomée in [Math. Comp. 26 (1972), 869–879]. The cut finite element method is a fictitious domain method with Nitsche-type enforcement of Dirichlet conditions together with stabilization of the elements at the boundary which is stable and enjoy optimal order approximation properties. A computational difficulty is, however, the geometric computations related to quadrature on the cut elements which must be accurate enough to achieve higher order approximation. With boundary value correction we may use only a piecewise linear approximation of the boundary, which is very convenient in a cut finite element method, and still obtain optimal order convergence. The boundary value correction is a modified Nitsche formulation involving a Taylor expansion in the normal direction compensating for the approximation of the boundary. Key to the analysis is a consistent stabilization term which enables us to prove stability of the method and a priori error estimates with explicit dependence on the meshsize and distance between the exact and approximate boundary.

### 1. INTRODUCTION

We consider a cut finite element method (CutFEM) for a second order elliptic boundary value problem with Dirichlet conditions. In standard fictitious domain CutFEM the boundary is represented on a background grid and allowed to cut through the elements in an arbitrary fashion. The Dirichlet conditions are enforced weakly using Nitsche’s method [22]. We refer to [4], [6], [8], [21], [19], for recent developments of this approach. See also the recent overview paper [7], and [20] for implementation issues.

Cut finite element methods is one way of alleviating the problem of mesh generation and allowing for more structured meshes and associated solvers. For this reason, the interest for such methods has increased significantly during the last few years; among recent contributions we mention the finite cell method of Parvizián, Düster, et al. [14, 23]; the least squares stabilized Lagrange multiplier methods of Haslinger and Renard [18], Tur et al. [25], and Baiges et al. [2]; the stabilization of Nitsche’s method by Codina and Baiges [12]; the local projection stabilization of multipliers of Barrenechea and Chouly [3] and of Amdouni, Moakher, and Renard [1].

---

Received by the editor July 11, 2015 and, in revised form, July 8, 2016 and October 31, 2016.  
2010 *Mathematics Subject Classification.* Primary 65M60; Secondary 65M85.

This research was supported in part by EPSRC, UK, Grant No. EP/J002313/1, the Swedish Foundation for Strategic Research Grant No. AM13-0029, the Swedish Research Council Grants Nos. 2011-4992, 2013-4708, and Swedish strategic research programme eSENCE.

In this contribution we develop a version of CutFEM based on the idea of boundary value correction originally proposed for standard finite element methods on an approximate domain in [5] and further developed in [13]. Using the closest point mapping to the exact boundary, or an approximation thereof, the boundary condition on the exact boundary may be weakly enforced using Nitsche's method on the boundary of the approximate domain. A Taylor expansion is used to approximate the value of the solution on the exact boundary in terms of the value and normal derivatives at the discrete approximate boundary. Key to the stability of the method is a consistent stabilization term that, also in the case of arbitrary cut elements at the boundary, provides control of the variation of the function in the vicinity of the boundary. More precisely, the stabilization ensures that the inverse inequality necessary to prove coercivity holds and that the resulting linear system of equations has the optimal condition number  $O(h^{-2})$ , where  $h$  is the mesh parameter, independent of the position of the boundary on the background grid. A different approach to the approximation of curved boundaries using extensions from subdomains was proposed by Cockburn et al. in [9–11].

We prove optimal order a priori error estimates, in the energy and  $L^2$ -norms, in terms of the error in the boundary approximation and the meshsize. Of particular practical importance is the fact that we may use a piecewise linear approximation of the boundary, which is very convenient from a computational point of view since the geometric computations are simple in this case and a piecewise linear distance function may be used to construct the discrete domain. We obtain optimal order convergence for higher order polynomial approximation of the solution if the Taylor expansion has sufficiently high order. In particular, for second and third order polynomials we obtain optimal order error estimates in the energy and  $L^2$ -norms with only one term in the Taylor expansion. Note that without boundary correction one typically requires  $O(h^{p+1})$  accuracy in the  $L^\infty$ -norm for the approximation of the domain which leads to significantly more involved computations on the cut elements for higher order elements; see [19]. However, also in the case of no boundary value correction our analysis in fact provides optimal order error estimates if the approximation of the boundary is accurate enough and thus we obtain an analysis for the standard cut finite element method with approximate boundary. Finally, we also prove estimates for the error both on the discrete domain and on the exact domain. The discrete solution on the exact domain is directly defined by the method since we may include all elements that intersect the union of the discrete and exact domains in the active mesh. Even though some active elements may not intersect the discrete domain the resulting method is stable due to the stabilization term and no auxiliary extension of the discrete solution outside of the discrete domain is necessary. We present numerical results illustrating our theoretical findings.

The outline of the paper is as follows: In Section 2 we formulate the model problem and our method, in Section 3 we present our theoretical analysis, and in Section 4 we present the numerical results.

## 2. MODEL PROBLEM AND METHOD

**2.1. The domain.** Let  $\Omega$  be a domain in  $\mathbb{R}^d$  with smooth boundary  $\partial\Omega$  and exterior unit normal  $n$ . We let  $\rho$  be the signed distance function, negative on the inside and positive on the outside, to  $\partial\Omega$  and we let  $U_\delta(\partial\Omega)$  be the tubular neighborhood

$\{x \in \mathbb{R}^d : |\rho(x)| < \delta\}$  of  $\partial\Omega$ . Then there is a constant  $\delta_0 > 0$  such that the closest point mapping  $p(x) : U_{\delta_0}(\partial\Omega) \rightarrow \partial\Omega$  is well defined and we have the identity  $p(x) = x - \rho(x)n(p(x))$ . We assume that  $\delta_0$  is chosen small enough that  $p(x)$  is well defined. See [16], Section 14.6 for further details on distance functions.

**2.2. The model problem.** We consider the problem: find  $u : \Omega \rightarrow \mathbb{R}$  such that

$$(2.1) \quad -\Delta u = f \quad \text{in } \Omega,$$

$$(2.2) \quad u = g \quad \text{on } \partial\Omega,$$

where  $f \in H^{-1}(\Omega)$  and  $g \in H^{1/2}(\partial\Omega)$  are given data. It follows from the Lax-Milgram Lemma that there exists a unique solution to this problem and we also have the elliptic regularity estimate

$$(2.3) \quad \|u\|_{H^{s+2}(\Omega)} \lesssim \|f\|_{H^s(\Omega)}, \quad s \geq -1.$$

Here and below we use the notation  $\lesssim$  to denote less or equal up to a constant.

**2.3. The mesh, discrete domains, and finite element spaces.**

- Let  $\Omega_0 \subset \mathbb{R}^d$  be a convex polygonal domain such that  $U_{\delta_0}(\Omega) \subset \Omega_0$ , where  $U_\delta(\Omega) = U_\delta(\partial\Omega) \cup \Omega$ . Let  $\mathcal{K}_{0,h}, h \in (0, h_0]$ , be a family of quasiuniform partitions, with mesh parameter  $h$ , of  $\Omega_0$  into shape regular triangles or tetrahedra  $K$ . We refer to  $\mathcal{K}_{0,h}$  as the background mesh.
- Given a subset  $\omega$  of  $\Omega_0$ , let  $\mathcal{K}_h(\omega)$  be the submesh defined by

$$(2.4) \quad \mathcal{K}_h(\omega) = \{K \in \mathcal{K}_{0,h} : \overline{K} \cap \overline{\omega} \neq \emptyset\},$$

i.e., the submesh consisting of elements that intersect  $\overline{\omega}$ , and let

$$(2.5) \quad \mathcal{N}_h(\omega) = \bigcup_{K \in \mathcal{K}_h(\omega)} K$$

be the union of all elements in  $\mathcal{K}_h(\omega)$ . Below the  $L^2$ -norm of discrete functions frequently should be interpreted as the broken norm. For example, for norms over  $\mathcal{N}_h$  we have

$$(2.6) \quad \|v\|_{\mathcal{N}_h(\omega)}^2 := \sum_{K \in \mathcal{K}_h(\omega)} \|v\|_K^2.$$

- Let the active mesh  $\mathcal{K}_h$  be defined by

$$(2.7) \quad \mathcal{K}_h := \mathcal{K}_h(\Omega \cup \Omega_h),$$

i.e., the submesh consisting of elements that intersect  $\Omega_h \cup \Omega$ , and let

$$(2.8) \quad \mathcal{N}_h := \mathcal{N}_h(\Omega \cup \Omega_h)$$

be the union of all elements in  $\mathcal{K}_h$ .

- Let  $V_{0,h}$  be the space of piecewise continuous polynomials of order  $p$  defined on  $\mathcal{K}_{0,h}$  and let the finite element space  $V_h$  be defined by

$$(2.9) \quad V_h := \{v_h : v_h := \tilde{v}_h|_{\mathcal{N}_h} \text{ for } \tilde{v}_h \in V_{0,h}\}.$$

- Let  $\Omega_h, h \in (0, h_0]$ , be a family of polygonal domains approximating  $\Omega$ , possibly independent of the computational mesh. We assume neither  $\Omega_h \subset \Omega$  nor  $\Omega \subset \Omega_h$ , instead the accuracy with which  $\Omega_h$  approximates  $\Omega$  will be crucial. To each  $\Omega_h$  we associate the functions  $\nu_h : \partial\Omega_h \rightarrow \mathbb{R}^d, |\nu_h| = 1$ , and  $\varrho_h : \partial\Omega_h \rightarrow \mathbb{R}$ , such that if  $p_h(x, \varsigma) := x + \varsigma\nu_h(x)$ , then  $p_h(x, \varrho_h(x)) \in \partial\Omega$  for all  $x \in \partial\Omega_h$ . We will also assume that  $p_h(x, \varsigma) \in U_{\delta_0}(\Omega)$  for all  $x \in \partial\Omega_h$

and all  $\varsigma$  between 0 and  $\varrho_h(x)$ . For conciseness we will drop the second argument of  $p_h$  below whenever it takes the value  $\varrho_h(x)$ . We assume that the following assumptions are satisfied:

$$(2.10) \quad \delta_h := \|\varrho_h\|_{L^\infty(\partial\Omega_h)} = o(h), \quad h \in (0, h_0],$$

and

$$(2.11) \quad \|\nu_h - n \circ p\|_{L^\infty(\partial\Omega_h)} = o(1), \quad h \in (0, h_0],$$

where  $o(\cdot)$  denotes the little ordo. We also assume that  $h_0$  is small enough to guarantee that

$$(2.12) \quad \partial\Omega_h \subset U_{\delta_0}(\partial\Omega), \quad h \in (0, h_0],$$

and that there exists  $M > 0$  such for any  $y \in U_{\delta_0}(\partial\Omega)$  the equation, find  $x \in \partial\Omega_h$  and  $|\varsigma| \leq \delta_h$  such that

$$(2.13) \quad p_h(x, \varsigma) = y,$$

has a solution set  $\mathcal{P}_h$  with

$$(2.14) \quad \text{card}(\mathcal{P}_h) \leq M$$

uniformly in  $h$ . The rationale of this assumption is to ensure that the image of  $p_h$  cannot degenerate for vanishing  $h$ .

- We note that it follows from (2.10) that

$$(2.15) \quad \|\rho\|_{L^\infty(\partial\Omega_h)} \lesssim \|\rho_h\|_{L^\infty(\partial\Omega_h)} = o(h)$$

since  $|\rho_h(x)| \geq |\rho(x)|$ ,  $x \in U_{\delta_0}(\partial\Omega)$ .

*The validity of assumption (2.14).* Assumption (2.14) will hold in any reasonable situation in practice. Here we give a proof in the special case where  $\nu_h$  is chosen constant on each element.

**Lemma 2.1.** *Assume that for all  $K \in \mathcal{N}_h(\partial\Omega_h)$ ,  $\nu_h|_{K \cap \partial\Omega_h} \in \mathbb{R}^d$ . Then there exists  $M > 0$  such that (2.14) holds uniformly in  $h$ .*

*Proof.* For a triangle  $K \in \mathcal{N}_h(\partial\Omega_h)$  define the domain  $E_K := \{x : x = x_K + \varsigma\nu_h(x_K); x_K \in K \cap \partial\Omega_h; -\delta_h \leq \varsigma \leq \delta_h\}$ . Then clearly for every  $y \in E_K$  the equation  $p_h(x, \varsigma) = y$  has a unique solution. It then suffices to show that  $\text{card}(\{K' \in \mathcal{N}_h(\partial\Omega_h) : E_K \cap E_{K'} \neq \emptyset\}) < M$ . That is  $E_K$  will have nonzero intersection with a finite number of other domains  $E_{K'}$ . To see this let  $B_{2\delta_h}(x_E)$  be a ball with radius  $2\delta_h$  centered at  $x_E \in E_K$  such that  $E_K \subset B_{2\delta_h}(x_E)$ . Then  $E_K \cap E_{K'} = \emptyset$  for any  $E_{K'}$  that has zero intersection with  $B_{2\delta_h}(x_E)$ ; this will be the case for  $E_{K'}$  for which  $K' \cap B_{3\delta_h}(x_E) = \emptyset$ . Since the mesh is quasiregular and shape optimal there exists  $M > 0$  such that  $\text{card}(\{K : K \cap B_{3\delta_h}(x_E) \neq \emptyset\}) \leq M$  uniformly in  $h$ . The claim then holds with this value on  $M$ .  $\square$

*The choice of  $\nu_h$ .* During computation, typically the quantities that are easily accessible on  $\partial\Omega_h$  are  $n_h$  and  $\rho$ . The two choices that are natural for  $\nu_h, \varrho_h$  are therefore  $\nu_h := n_h, \varrho_h := \varsigma$ , with  $\varsigma$  solution to  $\rho(p_h(x, \varsigma)) = 0$  or  $\nu_h := n \circ p$  and  $\varrho_h := \rho$ . Both cases requires the solution of nonlinear equations. The computation of  $\varrho_h$  using Newton’s method in the first case is substantially less costly than that of  $n \circ p$ , since the first quantity is a scalar and the initial guess  $\rho$  is more accurate.

Observe that if  $\nu_h := n \circ p$ , then the mapping  $p_h$  coincides with  $p(x)$ . It is therefore well defined and all the above assumptions hold by the properties of the

closest point mapping. This is not true in the general case. However, we assume that the equation  $\rho(p_h(x, \varsigma)) = 0$  has at least one solution for every  $x \in \partial\Omega_h$  and  $\varrho_h$  may then be identified with the solution of smallest magnitude. As an example consider the practically important case where  $\partial\Omega_h$  is defined by the zero level set of a piecewise linear nodal interpolant of the distance function and we choose  $\nu_h := n_h$ , with  $n_h$  denoting the normal of  $\partial\Omega_h$ . That the associated  $\varrho_h$  exists for all  $x \in \partial\Omega_h$  follows immediately from the implicit function theorem: the equation in  $\varsigma$ ,  $\rho(x + \varsigma n \circ p) = 0$ , has a solution since  $p$  is well defined and then so does  $\rho(x + \varsigma n_h) = 0$  since  $\nabla\rho \cdot n \circ p(x) > 0$  for  $h$  small enough. We show this using a fixed point argument. For  $x \in \partial\Omega_h$  let  $\varsigma_0$  solve the equation

$$\rho(\varsigma) := \rho(x + \varsigma n \circ p(x)) = 0$$

and define

$$\begin{aligned} \delta\rho &:= \frac{\partial\rho}{\partial\varsigma}(x + \varsigma n \circ p)|_{\varsigma=\varsigma_0} \\ &= \nabla\rho(x + \varsigma_0 n \circ p(x)) \cdot n \circ p(x) > 1 - C\delta_0 > 0, \text{ for } \delta_0 \text{ small enough.} \end{aligned}$$

Then consider the iterates  $\varsigma_k$  generated by

$$\varsigma_k = \varsigma_{k-1} - (\delta\rho)^{-1}\rho(x + \varsigma_{k-1}n_h), \text{ with } k \geq 1.$$

We will now show that this iteration converges. Let  $e_k = \varsigma_k - \varsigma_{k-1}$  and  $\bar{\varsigma}_k = s\varsigma_k + (1-s)\varsigma_{k-1}$  for some  $s \in [0, 1]$ , we may then write

$$e_k = (I - (\delta\rho)^{-1}\nabla\rho(x + \bar{\varsigma}_{k-1}n_h) \cdot n_h)e_{k-1}.$$

Using the mean value theorem we see that

$$\begin{aligned} \nabla\rho(x + \bar{\varsigma}_{k-1}n_h) \cdot n_h &= \delta\rho + \nabla\rho(x + \bar{\varsigma}_{k-1}n_h) \cdot (n_h - n \circ p(x)) \\ &\quad + (\bar{\varsigma}_{k-1}n_h - \varsigma_0 n)^T (\nabla \otimes \nabla\rho(\bar{x}) \cdot n \circ p(x)), \end{aligned}$$

where  $\bar{x} = x + t(\bar{\varsigma}_{k-1}n_h) + (1-t)\varsigma_0 n$  for some  $t \in [0, 1]$ . Therefore there exists  $C_\delta > 0$  such that

$$(2.16) \quad (I - (\delta\rho)^{-1}\nabla\rho(x + \bar{\varsigma}_{k-1}n_h) \cdot n_h) \leq C_\delta(h + (\bar{\varsigma}_{k-1} - \varsigma_0)).$$

Assuming that  $\|\mu_h - n \circ p\|_{L^\infty(\partial\Omega_h)} \lesssim h$  and  $|\varsigma_0| \lesssim h$  we may conclude using induction in the following way. Assume that there exists  $\tilde{C} > 0$  such that

$$(2.17) \quad e_k \leq h(\tilde{C}h)^k$$

and that  $h$  is small so that  $\tilde{C}h < 1$ . Observe that since  $|e_1| = |\delta\rho^{-1}\rho(x + \varsigma_0 n_h)| \leq \delta\rho^{-1}|\varsigma_0||n - n_h| \leq C_0 h^2$  this is true for  $e_1$  if  $h$  is chosen small enough. It follows that for  $k \geq 2$

$$\bar{\varsigma}_{k-1} - \varsigma_0 \leq \sum_{i=1}^{k-1} |e_i| \leq \tilde{C}h^2 \sum_{i=1}^{k-2} (\tilde{C}h)^{i-1} \leq \tilde{C}h^2.$$

Then considering (2.16) we obtain, with  $\tilde{C} = \max(C_0, C_\delta)$ ,

$$|e_{k+1}| \leq C_\delta(h + \tilde{C}h^2)|e_k| \leq h(\tilde{C}h)^{k+1},$$

and we conclude that (2.17) holds. As a consequence the sequence  $\{\varsigma_k\}_{k=0}^\infty$  is Cauchy, assume given two indices  $M > N$  there holds

$$|\varsigma_M - \varsigma_N| \leq \sum_{i=N+1}^M |e_i| \leq h(\tilde{C}h)^N.$$

Therefore there exists  $\varsigma$  such that  $\rho(x + \varsigma n_h) = 0$ . By Lemma 2.1 the assumption (2.14) clearly holds in this case. Moreover we have the estimates

$$(2.18) \quad \delta_h \lesssim h^2, \quad \|\nu_h - n \circ p\|_{L^\infty(\partial\Omega_h)} \lesssim h.$$

**2.4. Extensions.** There is an extension operator  $E : H^s(\Omega) \rightarrow H^s(U_{\delta_0}(\Omega))$  such that

$$(2.19) \quad \|Ev\|_{H^s(U_\delta(\Omega))} \lesssim \|v\|_{H^s(\Omega)}, \quad s \geq 0;$$

see [15]. For brevity we shall use the notation  $v$  for the extended function as well, i.e.,  $v = Ev$  on  $U_{\delta_0}(\Omega)$ .

### 2.5. The method.

*Derivation.* Let  $f = Ef$  and  $u = Eu$  be the extensions of  $f$  and  $u$  from  $\Omega$  to  $U_{\delta_0}(\Omega)$ . For  $v \in V_h$  we have using Green's formula

$$(2.20) \quad (f, v)_{\Omega_h} = (f + \Delta u, v)_{\Omega_h} - (\Delta u, v)_{\Omega_h}$$

$$(2.21) \quad = (f + \Delta u, v)_{\Omega_h \setminus \Omega} + (\nabla u, \nabla v)_{\Omega_h} - (n_h \cdot \nabla u, v)_{\partial\Omega_h},$$

where we used the fact  $f + \Delta u = 0$  on  $\Omega$ , while on  $\Omega_h \setminus \Omega$  we have  $f + \Delta u = Ef - \Delta Eu$ , which is not in general equal to zero. Now the boundary condition  $u = g$  on  $\partial\Omega$  may be enforced weakly as follows:

$$(2.22) \quad (f, v)_{\Omega_h} = (f + \Delta u, v)_{\Omega_h} + (\nabla u, \nabla v)_{\Omega_h} - (n_h \cdot \nabla u, v)_{\partial\Omega_h} \\ - (u \circ p_h - g \circ p_h, n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (u \circ p_h - g \circ p_h, v)_{\partial\Omega_h}.$$

The positive constant  $\beta$  must be chosen large enough to ensure stability; cf. below.

Since we do not have access to  $u \circ p_h$  we use a Taylor approximation in the direction  $\nu_h$ ,

$$(2.23) \quad u \circ p_h(x) \approx T_k(u)(x) := \sum_{j=0}^k \frac{D_{\nu_h}^j u(x)}{j!} \varrho_h^j(x),$$

where  $D_{\nu_h}^j$  is the  $j$ th partial derivative in the direction  $\nu_h$ . Thus it follows that the solution to (2.1)-(2.2) satisfies

$$(2.24) \quad (f, v)_{\Omega_h} = (f + \Delta u, v)_{\Omega_h} + (\nabla u, \nabla v)_{\Omega_h} - (n_h \cdot \nabla u, v)_{\partial\Omega_h} \\ - (T_k(u) - g \circ p_h, n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (T_k(u) - g \circ p_h, v)_{\partial\Omega_h} \\ - (u \circ p_h - T_k(u), n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (u \circ p_h - T_k(u), v)_{\partial\Omega_h}$$

for all  $v \in V_h$ . Rearranging the terms we arrive at

$$(2.25) \quad (\nabla u, \nabla v)_{\Omega_h} - (n_h \cdot \nabla u, v)_{\partial\Omega_h} \\ - (T_k(u), n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (T_k(u), v)_{\partial\Omega_h} \\ + (f + \Delta u, v)_{\Omega_h \setminus \Omega} \\ - (u \circ p_h - T_k(u), n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (u \circ p_h - T_k(u), v)_{\partial\Omega_h} \\ = (f, v)_{\Omega_h} - (g \circ p_h, n_h \cdot \nabla v)_{\partial\Omega_h} + \beta h^{-1} (g \circ p_h, v)_{\partial\Omega_h}$$

for all  $v \in V_h$ . The discrete method is obtained from this formulation by dropping the consistency terms of highest order, i.e., those on lines three and four of (2.25).

*Bilinear forms.* We define the forms

$$(2.26) \quad a_0(v, w) := (\nabla v, \nabla w)_{\Omega_h} - (n_h \cdot \nabla v, w)_{\partial\Omega_h} - (T_k(v), n_h \cdot \nabla w)_{\partial\Omega_h} + \beta h^{-1} (T_k(v), w)_{\partial\Omega_h},$$

$$(2.27) \quad a_h(v, w) := a_0(v, w) + j_h(v, w),$$

$$(2.28) \quad j_h(v, w) := \gamma_j \sum_{F \in \mathcal{F}_h} \sum_{l=1}^p h^{2l-1} ([D_{n_F}^l v], [D_{n_F}^l w])_F,$$

$$(2.29) \quad l_h(w) := (f, w)_{\Omega_h} - (g \circ p_h, n_h \cdot \nabla w)_{\partial\Omega_h} + \beta h^{-1} (g \circ p_h, w)_{\partial\Omega_h},$$

where  $\gamma_j$  and  $\beta$  are positive constants. Here we used the notation:

- $\mathcal{F}_h$  is the set of all internal faces to elements  $K \in \mathcal{K}_h$ , i.e., faces that are not included in the boundary of the active mesh  $\mathcal{K}_h$ , that intersect the set  $\Omega \setminus \Omega_h \cup \partial\Omega_h$ , and  $n_F$  is a fixed unit normal to  $F \in \mathcal{F}_h$ .
- $D_{n_F}^l$  is the partial derivative of order  $l$  in the direction of the normal  $n_F$  to the face  $F \in \mathcal{F}_h$ .
- $[v]_F = v_F^+ - v_F^-$ , with  $v_F^\pm = \lim_{s \rightarrow 0^\pm} v(x \mp sn_F)$ , is the jump of a discontinuous function  $v$  across a face  $F \in \mathcal{F}_h$ .
- The stabilizing term  $j_h(v, w)$  is introduced to extend the coercivity of  $a_0(\cdot, \cdot)$  to all of  $\mathcal{N}_h$  as we shall see below. Thanks to this property one may prove that the condition number is uniformly bounded independent of how  $\Omega_h$  is oriented compared to the mesh following the ideas of [6, 21].
- Observe the presence of the penalty coefficient  $\beta$  in (2.26) and (2.29). In order to guarantee coercivity  $\beta$  has to be chosen large enough and due to the Taylor expansions we also have to require that  $h \in (0, h_0]$  with  $h_0$  sufficiently small. See Section 3.3 and, in particular, Remark 3.2 for further details.

*The method.* Find:  $u_h \in V_h$  such that

$$(2.30) \quad a_h(u_h, v) = l_h(v), \quad \forall v \in V_h,$$

where  $a_h$  is defined in (2.27) and  $l_h$  in (2.29).

*Symmetric formulation in the case  $k = 1$ .* Using one term in the Taylor expansion gives the following forms:

$$(2.31) \quad a_h(v, w) = (\nabla v, \nabla w)_{\Omega_h} + j_h(v, w) - (n_h \cdot \nabla v, w)_{\partial\Omega_h} - (v, n_h \cdot \nabla w)_{\partial\Omega_h} - (\varrho_h \nu_h \cdot \nabla v, n_h \cdot \nabla w)_{\partial\Omega_h} + \beta h^{-1} (T_1(v), w)_{\partial\Omega_h},$$

$$(2.32) \quad l_h(w) = (f, w)_{\Omega_h} - (g \circ p_h, n_h \cdot \nabla w)_{\partial\Omega_h} + \beta h^{-1} (g \circ p_h, w)_{\partial\Omega_h}.$$

We see that only the term of the fourth line of (2.31) violate the symmetry of the formulation. To make it symmetric we choose  $\nu_h := n_h$ , assuming that the discrete approximation  $\Omega_h$  is such that this is a valid choice and also symmetrize the penalty term in the fourth line by replacing  $w$  in the right-hand slot by  $T_1(w)$ . A similar

perturbation is added to the right-hand side to keep consistency. The forms of the resulting symmetric formulation read:

$$\begin{aligned}
 (2.33) \quad a_h(v, w) &= (\nabla v, \nabla w)_{\Omega_h} + j_h(v, w) \\
 &\quad - (n_h \cdot \nabla v, w)_{\partial\Omega_h} - (v, n_h \cdot \nabla w)_{\partial\Omega_h} \\
 &\quad - (\varrho_h n_h \cdot \nabla v, n_h \cdot \nabla w)_{\partial\Omega_h} \\
 &\quad + \beta h^{-1} (T_1(v), T_1(w))_{\partial\Omega_h},
 \end{aligned}$$

$$(2.34) \quad l_h(w) = (f, w)_{\Omega_h} - (g \circ p_h, n_h \cdot \nabla w)_{\partial\Omega_h} + \beta h^{-1} (g \circ p_h, T_1(w))_{\partial\Omega_h}.$$

The analysis presented below covers this important special case. Also observe that if more terms are included in the Taylor series the resulting nonsymmetric part of the matrix is expected to be small, relative to the symmetric part, and the reduced symmetric form is likely to be a good preconditioner.

*Remark 2.1.* In principle it is possible to formulate a method using  $v \circ p_h$  instead of  $T_k(v)$  in the second and third lines of equation (2.26). Such a choice, however, may lead to nonstandard couplings in the system matrix corresponding to the form  $a_0(\cdot, \cdot)$  whenever  $p_h$  extends over an element boundary. Moreover the resulting method cannot be symmetrized.

### 3. A PRIORI ERROR ESTIMATES

**3.1. The energy norm.** Let the energy norm be defined by

$$(3.1) \quad |||v|||_h^2 := \|\nabla v\|_{\Omega_h}^2 + |||v|||_{j_h}^2 + h \|n_h \cdot \nabla v\|_{\partial\Omega_h}^2 + h^{-1} \|v\|_{\partial\Omega_h}^2,$$

where

$$(3.2) \quad |||v|||_{j_h}^2 = j_h(v, v).$$

**3.2. Consistency.** In view of (2.25) we obtain the identity

$$\begin{aligned}
 (3.3) \quad a_h(u - u_h, v) &= (u \circ p_h - T_k(u), n_h \cdot \nabla v)_{\partial\Omega_h} - \beta h^{-1} (u \circ p_h - T_k(u), v)_{\partial\Omega_h} \\
 &\quad + (f + \Delta u, v)_{\Omega_h \setminus \Omega}, \quad \forall v \in V_h,
 \end{aligned}$$

and thus we conclude that

$$\begin{aligned}
 (3.4) \quad |a_h(u - u_h, v)| &\leq \|u \circ p_h - T_k(u)\|_{\partial\Omega_h} \left( \|n_h \cdot \nabla v\|_{\partial\Omega_h} + h^{-1} \beta \|v\|_{\partial\Omega_h} \right) \\
 &\quad + \|f + \Delta u\|_{\Omega_h \setminus \Omega} \|v\|_{\Omega_h \setminus \Omega} \\
 (3.5) \quad &\leq h^{-1/2} \|u \circ p_h - T_k(u)\|_{\partial\Omega_h} |||v|||_h \\
 &\quad + \|f + \Delta u\|_{\Omega_h \setminus \Omega} \|v\|_{\Omega_h \setminus \Omega}, \quad \forall v \in V_h.
 \end{aligned}$$

*Estimate of the error in the Taylor approximation.* The Taylor polynomial  $T_k(u)(x)$  provides an approximation of  $u \circ p_h(x)$  and we have the error estimate

$$(3.6) \quad |v \circ p_h(x) - T_k(v)(x)| \lesssim \left| \int_0^{\varrho_h(x)} D_{\nu_h}^{k+1} v(x(s)) (\varrho_h(x) - s)^k ds \right|$$

$$(3.7) \quad \lesssim \|D_{\nu_h}^{k+1} v\|_{I_x} \|(\varrho_h(x) - s)^k\|_{I_x}$$

$$(3.8) \quad \lesssim \|D_{\nu_h}^{k+1} v\|_{I_x} |\varrho_h(x)|^{k+1/2},$$



where  $I_x$  is the line segment between  $x$  and  $p_h(x)$ . Combining (3.4) and (3.8) and recalling the assumption (2.14) we arrive at the estimate

$$(3.9) \quad \|v \circ p_h - T_k(v)\|_{\partial\Omega_h}^2 \lesssim \int_{\partial\Omega_h} \|D_{\nu_h}^{k+1}v\|_{I_x}^2 |\varrho_h(x)|^{2k+1} dx$$

$$(3.10) \quad \lesssim \int_{\partial\Omega_h} \|D_{\nu_h}^{k+1}v\|_{I_{\delta_h}}^2 |\varrho_h(x)|^{2k+1} dx$$

$$(3.11) \quad \lesssim \delta_h^{2k+1} \|D^{k+1}v\|_{U_{\delta_h}(\partial\Omega_h)}^2.$$

Here we handled the possible overlap of the contributions from different polygonal sides of  $\partial\Omega_h$  by using the fact that by assumption (2.14) such an overlap must have a finite number of contributions uniformly in  $h$  and by dropping the directional derivative, effectively including the derivatives of order  $k + 1$  in all directions.

With slightly stronger control of the regularity,  $v \in H^{k+\frac{3}{2}}(\Omega_0)$ , we obtain the estimate

$$(3.12) \quad \|v \circ p - T_k(v)\|_{\partial\Omega_h} \lesssim \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1}v\|_{\partial\Omega_t},$$

where  $\partial\Omega_t = \{x \in \Omega : \rho(x) = t\}$  is the levelset with distance  $t$  to the boundary  $\partial\Omega$ .

*Estimate of the residual on  $\Omega_h \setminus \Omega$ .* Suppose that

$$(3.13) \quad f + \Delta u \in H^{l+\frac{1}{2}+\epsilon}(U_{\delta_0}(\Omega))$$

with  $\epsilon > 0$  for  $l = 0$  and  $\epsilon = 0$  for  $l \geq 1$ , which, in view of (2.3) and (2.19), holds if  $f \in H^{l+\frac{1}{2}+\epsilon}(\Omega)$ . Using (3.13) and the fact that  $f + \Delta u = 0$  in  $\Omega$ , we obtain the estimate

$$(3.14) \quad \|f + \Delta u\|_{\Omega_h \setminus \Omega} \lesssim \delta_h^l \|D_n^l(f + \Delta u)\|_{\Omega_h \setminus \Omega} \lesssim \delta_h^{l+1/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t},$$

where we used the fact that  $\Omega_h \setminus \Omega \subset U_\delta(\partial\Omega)$ , where  $\delta \sim \delta_h$ .

*Estimates of the consistency error.* Combining (3.12), (3.14), and (3.16), we obtain the estimate

$$(3.15) \quad \begin{aligned} |a_h(u - u_h, v)| &\leq \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1}u\|_{\partial\Omega_t} \left( \|n_h \cdot \nabla v\|_{\partial\Omega_h} + h^{-1}\beta \|v\|_{\partial\Omega_h} \right) \\ &\quad + \delta_h^{l+1/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} \|v\|_{\Omega_h \setminus \Omega}, \quad \forall v \in V_h. \end{aligned}$$

This estimate will be used when we derive an  $L^2$  estimate of the error while for the energy error estimate we continue the estimation using the bound (for a proof see the Appendix)

$$(3.16) \quad \|v\|_{\Omega_h \setminus \Omega} \lesssim h^{1/2} \delta_h^{1/2} \|v\|_h, \quad \forall v \in V_h.$$

This leads to

$$(3.17) \quad \begin{aligned} |a_h(u - u_h, v)| &\leq \left( h^{-1/2} \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1}u\|_{\partial\Omega_t} \right. \\ &\quad \left. + h^{1/2} \delta_h^{l+1} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} \right) \|v\|_h, \quad \forall v \in V_h. \end{aligned}$$

*Remark 3.1.* We may upper bound the right-hand sides further using global trace inequalities leading to

$$(3.18) \quad \sup_{0 \leq t \leq \delta_0} \|D^{k+1}u\|_{\partial\Omega_t} \lesssim \|u\|_{H^{k+2}(\Omega)} \lesssim \|f\|_{H^k(\Omega)}$$

and

$$(3.19) \quad \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} \lesssim \|f\|_{H^{l+1}(\Omega)} + \|\Delta u\|_{H^{l+1}(\Omega)} \lesssim \|f\|_{H^{l+1}(\Omega)}.$$

The constants in the above inequalities depend on the regularity of the domain.

**3.3. Coercivity and continuity.** A key element of the analysis is that the addition of the stabilization operator  $j_h(\cdot, \cdot)$  allows us to prove coercivity of the bilinear form, independent of how the approximate domain  $\Omega_h$  intersects the computational mesh. This draws on previous results from [6, 21]. In particular the following results hold:

$$(3.20) \quad \|\nabla v\|_{\mathcal{N}_h}^2 \lesssim \|\nabla v\|_{\Omega_h}^2 + \|v\|_{j_h}^2, \quad \forall v \in V_h,$$

and

$$(3.21) \quad \|v\|_{\mathcal{N}_h}^2 \lesssim \|v\|_{\Omega_h}^2 + h^2 \|v\|_{j_h}^2, \quad \forall v \in V_h.$$

Below we will use the notation

$$(3.22) \quad T_{1,k}(v) = T_k(v) - v$$

and some inverse estimates that we collect in the following technical lemma.

**Lemma 3.1.** *For all  $v \in V_h$  there holds*

$$(3.23) \quad h^{1/2} \|n_h \cdot \nabla v\|_{\partial\Omega_h} \lesssim \|\nabla v\|_{\mathcal{N}_h(\partial\Omega_h)},$$

$$(3.24) \quad h^{-1/2} \|T_{1,k}(v)\|_{\partial\Omega_h} \lesssim \gamma(h) \|\nabla v\|_{\mathcal{N}_h(\partial\Omega_h)},$$

where  $\gamma(h) \rightarrow 0$  as  $h \rightarrow 0$ .

*Proof.* Inequality (3.23) then follows using a standard trace inequality elementwise, followed by an inverse inequality. To prove the inequality (3.24) observe that by the inverse inequality  $\|D^j v\|_K \lesssim h^{1-j} \|\nabla v\|_K$  there holds

$$\|T_{1,k}(v)\|_{\mathcal{N}_h(\partial\Omega_h)} \leq \sum_{j=1}^k \frac{\varrho_h^j}{j!} \|D^j v\|_{\mathcal{N}_h(\partial\Omega_h)} \lesssim h \left( \sum_{j=1}^k \frac{\varrho_h^j}{j! h^j} \right) \|\nabla v\|_{\mathcal{N}_h(\partial\Omega_h)}$$

and, consequently, since  $\varrho_h \leq \delta_h$ ,

$$(3.25) \quad h^{-1/2} \|T_{1,k}(v)\|_{\partial\Omega_h} \lesssim h^{-1} \|T_{1,k}(v)\|_{\mathcal{N}_h(\partial\Omega_h)}$$

$$(3.26) \quad \lesssim \underbrace{\left( \sum_{j=1}^k \frac{\delta_h^j}{h^j} \right)}_{\lesssim \gamma(h) \sim h^{-1} o(h)} \|\nabla v\|_{\mathcal{N}_h(\partial\Omega_h)}$$

$$(3.27) \quad \lesssim \gamma(h) \|\nabla v\|_{\mathcal{N}_h(\partial\Omega_h)}. \quad \square$$

The property  $\delta_h = o(h)$  is necessary to guarantee that  $\gamma(h) \rightarrow 0$  as  $h \rightarrow 0$ , which is important for the proof of the following result.

**Proposition 3.1.** *Assume that  $h_0$  is small enough and  $\beta$  is large enough then there holds*

$$(3.28) \quad |||v|||_h^2 \lesssim a_h(v, v), \quad \forall h \in (0, h_0] \text{ and } \forall v \in V_h.$$

*Proof.* Taking  $w = v$  in (2.27) we obtain

$$(3.29) \quad \begin{aligned} a_h(v, v) &= (\nabla v, \nabla v)_{\Omega_h} + j_h(v, v) - 2(n_h \cdot \nabla v, v)_{\partial\Omega_h} + \beta h^{-1}(v, v)_{\partial\Omega_h} \\ &\quad + \beta h^{-1}(T_{1,k}(v), v)_{\partial\Omega_h} - (T_{1,k}(v), n_h \cdot \nabla v)_{\partial\Omega_h} \\ (3.30) \quad &\geq \|\nabla v\|_{\Omega_h}^2 + |||v|||_{j_h}^2 - 2h^{1/2}\|n_h \cdot \nabla v\|_{\partial\Omega_h} h^{-1/2}\|v\|_{\partial\Omega_h} + \beta h^{-1}\|v\|_{\partial\Omega_h}^2 \\ &\quad - \beta h^{-1/2}\|T_{1,k}(v)\|_{\partial\Omega_h} h^{-1/2}\|v\|_{\partial\Omega_h} \\ &\quad - h^{-1/2}\|T_{1,k}(v)\|_{\partial\Omega_h} h^{1/2}\|n_h \cdot \nabla v\|_{\partial\Omega_h}. \end{aligned}$$

Using (3.20) we have

$$(3.31) \quad \|\nabla v\|_{\mathcal{N}_h}^2 \lesssim \|\nabla v\|_{\Omega_h}^2 + |||v|||_{j_h}^2.$$

Next we apply the inverse bounds (3.23) and (3.24) and the arithmetic-geometric inequality to deduce

$$(3.32) \quad h^{1/2}\|n_h \cdot \nabla v\|_{\partial\Omega_h} h^{-1/2}\|v\|_{\partial\Omega_h} \lesssim \varepsilon^{-1}\beta^{-1}\|\nabla v\|_{\mathcal{N}_h}^2 + \varepsilon\beta h^{-1}\|v\|_{\partial\Omega_h}^2,$$

$$(3.33) \quad \beta h^{-1/2}\|T_{1,k}(v)\|_{\partial\Omega_h} h^{-1/2}\|v\|_{\partial\Omega_h} \lesssim \varepsilon^{-1}\beta\gamma^2(h)\|\nabla v\|_{\mathcal{N}_h}^2 + \varepsilon\beta h^{-1}\|v\|_{\partial\Omega_h}^2,$$

$$(3.34) \quad h^{-1/2}\|T_{1,k}(v)\|_{\partial\Omega_h} h^{1/2}\|n_h \cdot \nabla v\|_{\partial\Omega_h} \lesssim \gamma(h)\|\nabla v\|_{\mathcal{N}_h}^2.$$

Using these relations we have, for positive constants  $c_1$  and  $c_2$ ,

$$(3.35) \quad a_h(v, v) \geq (c_1 - c_2(\gamma(h) + \gamma^2(h)\beta + \varepsilon^{-1}\beta^{-1}))\|\nabla v\|_{\mathcal{N}_h}^2 + \beta(1 - c_2\varepsilon)h^{-1}\|v\|_{\partial\Omega_h}^2.$$

To conclude, fix  $\varepsilon$  small enough so that  $(1 - c_2\varepsilon) > 0$ , and then observe that  $(c_1 - c_2(\gamma(h) + \gamma^2(h)\beta + \varepsilon^{-1}\beta^{-1})) > 0$  if  $\beta$  is large enough and, since  $\gamma(h) \rightarrow 0$  as  $h \rightarrow 0$ , for  $h \in (0, h_0]$ , with  $h_0$  small enough.  $\square$

*Remark 3.2.* Considering the practically relevant case when  $\Omega_h$  is a piecewise linear approximation of  $\Omega$  such that  $\delta_h \lesssim h^2$  we have  $\gamma(h) \leq c_3h \leq c_3h_0$  for  $h \in (0, h_0]$ . First taking  $\varepsilon = 1/(2c_2)$ , we get  $\beta(1 - c_2\varepsilon) = \beta/2$ . Next we have

$$(3.36) \quad c_1 - c_2(\gamma(h) + \gamma^2(h)\beta + \varepsilon^{-1}\beta^{-1}) \geq c_1 - c_2c_3h_0 - c_2c_3h_0^2\beta - 2c_2\beta^{-1} \geq c_1/2,$$

where we choose  $\beta$  and  $h_0$  such that each of the three negative factors have absolute value less or equal to  $c_1/6$ . These choices are

$$(3.37) \quad \beta = 12c_2/c_1, \quad h_0 = \min(c_1/(6c_2c_3), c_1/(\sqrt{72}c_2c_3)) = c_1/(\sqrt{72}c_2c_3).$$

Define the space  $V$  on which the functional  $V \ni v \mapsto a_h(v, w) \in \mathbb{R}$ , for a fixed  $w \in V_h$  and fixed  $h \in (0, h_0]$  is bounded,

$$(3.38) \quad V = H^{k+1/2}(\mathcal{N}_h) \cap H^{3/2}(\mathcal{N}_h) \cap H^{p+1/2}(\mathcal{N}_h).$$

Then we may write the continuity of  $a_h(\cdot, \cdot)$ .

**Proposition 3.2.** *Let  $v \in V + V_h$  and  $w \in V_h$ , then there holds*

$$(3.39) \quad a_h(v, w) \lesssim \left( |||v|||_h + h^{-1/2}\|T_{1,k}(v)\|_{\partial\Omega_h} \right) |||w|||_h, \quad \forall v \in V + V_h, w \in V_h.$$

*Proof.* The continuity estimate (3.39) follows directly from the Cauchy-Schwarz inequality applied term by term to the definition of  $a_h(\cdot, \cdot)$ , (2.27).  $\square$

**3.4. Interpolation estimates.** Let

$$(3.40) \quad \pi_h : H^1(\Omega) \ni u \mapsto \pi_{SZ,h} E u \in V_h,$$

where  $E$  is the extension operator introduced in Section 2.4, and  $\pi_{SZ,h}$  is the Scott-Zhang interpolation operator. The following error estimate for the Scott-Zhang interpolant is well known [24]:

$$(3.41) \quad \|u - \pi_{SZ,h} u\|_{H^m(K)} \lesssim h^{s-m} \|u\|_{H^s(\mathcal{N}_h(K))}, \quad 0 \leq m \leq s \leq p + 1, \quad K \in \mathcal{K}_h.$$

Using the properties of the extension operator we then immediately deduce this interpolation error estimate for (3.40):

$$(3.42) \quad \| \|u - \pi_h u\|_h + h^{-1/2} \|T_{1,k}(u - \pi_h u)\|_{\partial\Omega_h} \| \lesssim h^p \|u\|_{H^{p+1}(\Omega)}.$$

*Verification of (3.42).* The first term in (3.42) has four contributions (see (3.1)). The energy-norm contribution is bounded directly by (3.41). For the two last contributions of (3.1) using the trace inequality

$$(3.43) \quad \|v\|_{\partial\Omega_h \cap K}^2 \lesssim h^{-1} \|v\|_K^2 + h \|\nabla v\|_K^2, \quad K \in \mathcal{K}_h$$

(see [17]), followed by the interpolation estimate (3.41) and stability of the extension operator (2.19) we get the desired result. Finally to estimate  $\| \|u - \pi_h u\|_{j_h}$  observe that on each simplex we have

$$\|v\|_{\partial K}^2 \lesssim h^{-1} \|v\|_K^2 + h \|\nabla v\|_K^2$$

and we proceed elementwise as before using (3.41) and stability of the extension operator (2.19).

Again using the trace inequality (3.43) the second term in (3.42) can be estimated as

$$(3.44) \quad h^{-1/2} \|T_{1,k}(u - \pi_h u)\|_{\partial\Omega_h} \lesssim h^{-1} \|T_{1,k}(u - \pi_h u)\|_{\mathcal{N}_h(\partial\Omega_h)} + \|\nabla T_{1,k}(u - \pi_h u)\|_{\mathcal{N}_h(\partial\Omega_h)}$$

$$(3.45) \quad \lesssim h^p \|u\|_{H^{p+1}(\Omega)},$$

where finally we used the fact that  $\delta_h \lesssim h$  and the estimate

$$(3.46) \quad h^{m-1} \|\nabla^m T_{1,k}(u - \pi_h u)\|_K \lesssim \sum_{j=1}^k \delta_h^j h^{m-1} \|(u - \pi_h u)\|_{H^{j+m}(K)}$$

$$(3.47) \quad \lesssim \sum_{j=1}^k h^j h^{m-1} h^{p+1-(j+m)} \|u\|_{H^{p+1}(\mathcal{N}(K))}$$

$$(3.48) \quad \lesssim h^p \|u\|_{H^{p+1}(\mathcal{N}(K))}$$

for  $m = 0, 1$  and  $K \in \mathcal{K}_h(\partial\Omega_h)$ . □

**3.5. Error estimates.**

**Theorem 3.1.** *If  $\delta_h = o(h)$ , then the following estimate holds:*

$$(3.49) \quad \| \|u - u_h\|_h \| \lesssim h^p \|u\|_{H^{p+1}(\Omega)} + h^{-1/2} \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1} u\|_{\partial\Omega_t} + h^{1/2} \delta_h^{l+1} \sup_{-\delta_0 \leq t < 0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t}.$$

*Proof.* We first note that adding and subtracting an interpolant and using the triangle inequality and the interpolation estimate (3.42), we obtain

$$(3.50) \quad |||u - u_h|||_h \lesssim |||u - \pi_h u|||_h + |||\pi_h u - u_h|||_h$$

$$(3.51) \quad \lesssim h^p |||u|||_{H^{p+1}(\Omega)} + |||\pi_h u - u_h|||_h.$$

For the second term on the right-hand side we have the estimates

$$(3.52)$$

$$(3.53) \quad \begin{aligned} |||\pi_h u - u_h|||_h^2 &\lesssim a_h(\pi_h u - u_h, \pi_h u - u_h) \\ &= a_h(\pi_h u - u, \pi_h u - u_h) + a_h(u - u_h, \pi_h u - u_h) \end{aligned}$$

$$(3.54) \quad \begin{aligned} &\lesssim \left( |||\pi_h u - u|||_h + h^{-1/2} \|T_{1,k}(\pi_h u - u)\|_{\partial\Omega_h} \right) |||\pi_h u - u_h|||_h \\ &\quad + h^{-1/2} \|u \circ p_h - T_k(u)\|_{\partial\Omega_h} |||\pi_h u - u_h|||_h \\ &\quad + \|f + \Delta u\|_{\Omega_h \setminus \Omega_h} \|\pi_h u - u_h\|_{\Omega_h \setminus \Omega} \end{aligned}$$

$$(3.55) \quad \begin{aligned} &\lesssim h^p |||u|||_{H^{p+1}(\Omega)} |||\pi_h u - u_h|||_h \\ &\quad + h^{-1/2} \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1}u\|_{\partial\Omega_t} |||\pi_h u - u_h|||_h \\ &\quad + h^{1/2} \delta_h^{l+1} \sup_{-\delta_0 \leq t < 0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} |||\pi_h u - u_h|||_h, \end{aligned}$$

where we used coercivity (3.28), added and subtracted the exact solution  $u$ , estimated the first term using continuity (3.39) followed by the interpolation estimate (3.42) and the second using the consistency estimate (3.6). Combining estimates (3.51) and (3.55) concludes the proof.  $\square$

**Theorem 3.2.** *If  $\delta_h \lesssim h^2$ , then the following estimate holds:*

$$(3.56) \quad \begin{aligned} \|u - u_h\|_{\Omega_h} &\lesssim h^{p+1} |||u|||_{H^{p+1}(\Omega)} \\ &\quad + \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1}u\|_{\partial\Omega_t} \\ &\quad + \delta_h^{l+3/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t}. \end{aligned}$$

*Proof.* Let  $\phi \in H_0^1(\Omega)$  be the solution to the dual problem

$$(3.57) \quad a(v, \phi) = (v, \psi)_\Omega, \quad v \in H_0^1(\Omega),$$

where  $\psi : \Omega_h \cup \Omega \mapsto \mathbb{R}$  takes the values  $\psi = u - u_h$  on  $\Omega_h$  and  $\psi = 0$  on  $\Omega \setminus \Omega_h$ . We may then extend  $\phi$  using the extension operator to  $U_{\delta_0}(\Omega)$ , using the same notation for the extended function. By standard regularity theory we have the stability estimate

$$(3.58) \quad \|\phi\|_{H^2(\Omega)} \lesssim \|\psi\|_{\Omega \cap \Omega_h}.$$

We obtain the following representation formula for the error  $e = u - u_h$ :

$$(3.59) \quad \|e\|_{\Omega_h}^2 = (e, \psi + \Delta\phi)_{\Omega_h} - (e, \Delta\phi)_{\Omega_h}$$

$$(3.60) \quad = (e, \psi + \Delta\phi)_{\Omega_h \setminus \Omega} + (\nabla e, \nabla\phi)_{\Omega_h} - (e, n_h \cdot \nabla\phi)_{\partial\Omega_h}$$

$$(3.61) \quad = (e, \psi + \Delta\phi)_{\Omega_h \setminus \Omega} + a_0(e, \phi) + b_h(e, \phi)$$

$$(3.62) \quad = I + II + III,$$

where

$$(3.63) \quad III = (T_k(e) - e, n_h \cdot \nabla \phi)_{\partial\Omega_h} - \beta h^{-1} (T_k(e), \phi)_{\partial\Omega_h} + (n_h \cdot \nabla e, \phi)_{\partial\Omega_h}$$

$$(3.64) \quad = (T_{1,k}(e), n \cdot \nabla \phi)_{\partial\Omega_h} - \beta h^{-1} (e, \phi)_{\partial\Omega_h} \\ - \beta h^{-1} (T_{1,k}(e), \phi)_{\partial\Omega_h} + (n_h \cdot \nabla e, \phi)_{\partial\Omega_h}.$$

*Term I.* We have

$$(3.65) \quad |I| = |(e, \psi + \Delta\phi)_{\Omega_h \setminus \Omega}|$$

$$(3.66) \quad \lesssim \|e\|_{\Omega_h \setminus \Omega} \|\psi + \Delta\phi\|_{\Omega_h \setminus \Omega}$$

$$(3.67) \quad \lesssim \left( \delta_h^2 \|n \cdot \nabla e\|_{\Omega_h \setminus \Omega}^2 + \delta_h \|e\|_{\partial\Omega_h}^2 \right)^{1/2} \left( \|\psi\|_{\Omega_h \setminus \Omega} + \|\Delta\phi\|_{\Omega_h \setminus \Omega} \right)$$

$$(3.68) \quad \lesssim \left( (\delta_h^2 + h\delta_h) \|e\|_h^2 \right)^{1/2} \left( \|e\|_{\Omega_h \setminus \Omega} + \|\phi\|_{H^2(\Omega)} \right)$$

$$(3.69) \quad \lesssim \underbrace{(h^{-2}\delta_h + h^{-1}\delta_h)^{1/2}}_{\lesssim 1} h \|e\|_h \|e\|_{\Omega_h}.$$

Here we used the estimate

$$(3.70) \quad \|v\|_{\Omega_h \setminus \Omega}^2 \lesssim \delta_h^2 \|n \cdot \nabla v\|_{\Omega_h \setminus \Omega}^2 + \delta_h \|v\|_{\partial\Omega_h}^2, \quad v \in H^1(\Omega_h),$$

with  $v = e$ , the definition of the energy norm to conclude that  $h^{-1} \|e\|_{\partial\Omega_h}^2 \lesssim \|e\|_h^2$ , the stability (2.19) of the extension operator, the stability (3.58) of the dual problem and the assumption that  $\delta_h \lesssim h^2$ .

*Term II.* Adding and subtracting an interpolant we obtain

$$(3.71) \quad |II| = |a_h(e, \phi - \pi_h \phi) + a_h(e, \pi_h \phi)|$$

$$(3.72) \quad \lesssim \|e\|_h \|\phi - \pi_h \phi\|_h + |a_h(e, \pi_h \phi)|$$

$$(3.73) \quad \lesssim h \|e\|_h \|\phi\|_{H^2(\Omega)} + |a_h(e, \pi_h \phi)|$$

$$(3.74) \quad \lesssim h \|e\|_h \|e\|_{\Omega_h} + |a_h(e, \pi_h \phi)|.$$

To estimate the second term on the right-hand side we employ (3.15), with  $v = \pi_h \phi$ ,

$$(3.75) \quad |a_h(e, \pi_h \phi)| \leq \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1} u\|_{\partial\Omega_t} \left( \|n_h \cdot \nabla \pi_h \phi\|_{\partial\Omega_h} + h^{-1} \beta \|\pi_h \phi\|_{\partial\Omega_h} \right) \\ + \delta_h^{l+1/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l (f + \Delta u)\|_{\partial\Omega_t} \|\pi_h \phi\|_{\Omega_h \setminus \Omega}.$$

Here we have the estimates

$$(3.76)$$

$$\|n_h \cdot \nabla \pi_h \phi\|_{\partial\Omega_h} + h^{-1} \|\pi_h \phi\|_{\partial\Omega_h} \lesssim \|n_h \cdot \nabla (\pi_h \phi - \phi)\|_{\partial\Omega_h} + h^{-1} \|\pi_h \phi - \phi\|_{\partial\Omega_h} \\ + \|n_h \cdot \nabla \phi\|_{\partial\Omega_h} + h^{-1} \|\phi\|_{\partial\Omega_h}$$

$$(3.77) \quad \lesssim h^{-1/2} \|\pi_h \phi - \phi\|_h \\ + \|\phi\|_{H^2(\Omega_h)} h^{-1} \delta_h^{1/2} \|\phi\|_{H^1(U_{\delta_h}(\partial\Omega))}$$

$$(3.78) \quad \lesssim h^{1/2} \|\phi\|_{H^2(\Omega)} + \|\phi\|_{H^2(\Omega_h)} + h^{-1} \delta_h^{1/2} \|\phi\|_{H^1(U_{\delta_h}(\partial\Omega))}$$

$$(3.79) \quad \lesssim (h^{1/2} + 1 + h^{-1} \delta_h^{1/2}) \|e\|_{\Omega_h}$$

$$(3.80) \quad \lesssim \|e\|_{\Omega_h}$$

and

$$(3.81) \quad \|\pi_h \phi\|_{\Omega_h \setminus \Omega} \leq \|\pi_h \phi - \phi\|_{\Omega_h \setminus \Omega} + \|\phi\|_{\Omega_h \setminus \Omega}$$

$$(3.82) \quad \lesssim h^2 \|\phi\|_{H^2(\Omega)} + \delta_h \|\nabla \phi\|_{U_{\delta_h}(\partial\Omega)}$$

$$(3.83) \quad \lesssim (h^2 + \delta_h) \|e\|_{\Omega_h}$$

$$(3.84) \quad \lesssim \delta_h \|e\|_{\Omega_h},$$

where, in both estimates, we used the assumption  $\delta_h \lesssim h^2$ , as well as the following bounds:

$$(3.85) \quad \|\phi\|_{\partial\Omega_h} \lesssim \delta_h^{1/2} \|n \cdot \nabla \phi\|_{U_{\delta_h}(\partial\Omega)},$$

$$(3.86) \quad \|\phi\|_{\Omega_h \setminus \Omega} \lesssim \|\phi\|_{U_{\delta_h}(\partial\Omega)} \lesssim \delta_h \|n \cdot \nabla \phi\|_{U_{\delta_h}(\partial\Omega)};$$

see the Appendix for the proof of these estimates. Combining estimates (3.75), (3.76), and (3.81), we arrive at

$$(3.87) \quad |a_h(e, \pi_h \phi)| \lesssim \left( \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1} u\|_{\partial\Omega_t} + \delta_h^{l+3/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} \right) \|e\|_{\Omega_h}$$

which together with (3.74) gives

$$(3.88) \quad |III| \lesssim \left( h \|e\|_h + \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1} u\|_{\partial\Omega_t} + \delta_h^{l+3/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l(f + \Delta u)\|_{\partial\Omega_t} \right) \|e\|_{\Omega_h}.$$

*Term III.* Using the Cauchy-Schwarz inequality we get

$$(3.89) \quad |III| = |b_h(e, \phi)|$$

$$(3.90) \quad \lesssim \|T_{1,k}(e)\|_{\partial\Omega_h} \|n_h \cdot \nabla \phi\|_{\partial\Omega_h} + \beta h^{-1} \|e\|_{\partial\Omega_h} \|\phi\|_{\partial\Omega_h} + \beta h^{-1} \|T_{1,k}(e)\|_{\partial\Omega_h} \|\phi\|_{\partial\Omega_h} + \|n_h \cdot \nabla e\|_{\partial\Omega_h} \|\phi\|_{\partial\Omega_h}$$

$$(3.91) \quad \lesssim \|T_{1,k}(e)\|_{\partial\Omega_h} \left( h^{-1} \|\phi\|_{\partial\Omega_h} + \|n_h \cdot \nabla \phi\|_{\partial\Omega_h} \right) + \|e\|_h h^{-1/2} \|\phi\|_{\partial\Omega_h}$$

$$(3.92) \quad \lesssim \left( \|T_{1,k}(e)\|_{\partial\Omega_h} + h^{-1/2} \delta_h \|e\|_h \right) \|e\|_{\Omega_h}$$

$$(3.93) \quad \lesssim \left( h^{p+1} \|u\|_{H^{p+1}(\Omega)} + (h^{-3/2} \delta_h) h \|e\|_h \right) \|e\|_{\Omega_h},$$

where we used (3.85) and (3.86) followed by the stability estimate for the dual problem (3.58), and at last the estimate

$$(3.94) \quad \|T_{1,k}(e)\|_{\partial\Omega_h} \lesssim h^{p+1} \|u\|_{H^{p+1}(\Omega)} + (h^{-3/2} \delta_h) h \|e\|_h.$$

*Verification of (3.94).* We have

$$(3.95) \quad \|T_{1,k}(e)\|_{\partial\Omega_h} \lesssim \sum_{j=1}^k \delta_h^j \|D_{\nu_h}^j e\|_{\partial\Omega_h}$$

and for each of the terms  $\|D_{\nu_h}^j e\|_{\partial\Omega_h}$ ,  $j = 1, \dots, k$ , we obtain by adding and subtracting an interpolant, using the interpolation estimate (3.41) for the first term and an inverse estimate for the second, the estimates

(3.96)

$$\begin{aligned} \|D_{\nu_h}^j e\|_{\partial\Omega_h}^2 &\lesssim h^{-1} \|D_{\nu_h}^j e\|_{\mathcal{N}_h(\partial\Omega_h)}^2 + h \|\nabla D_{\nu_h}^j e\|_{\mathcal{N}_h(\partial\Omega_h)}^2 \\ (3.97) \quad &\lesssim h^{-1} \|D_{\nu_h}^j (u - \pi_h u)\|_{\mathcal{N}_h(\partial\Omega_h)}^2 + h \|\nabla D_{\nu_h}^j (u - \pi_h u)\|_{\mathcal{N}_h(\partial\Omega_h)}^2 \\ &\quad + h^{-1} \|D_{\nu_h}^j (\pi_h u - u_h)\|_{\mathcal{N}_h(\partial\Omega_h)}^2 + h \|\nabla D_{\nu_h}^j (\pi_h u - u_h)\|_{\mathcal{N}_h(\partial\Omega_h)}^2 \end{aligned}$$

$$(3.98) \quad \lesssim h^{2p+1-2j} \|u\|_{H^{p+1}(\mathcal{N}_h(\mathcal{N}_h(\partial\Omega_h)))}^2 + h^{1-2j} \|\nabla(\pi_h u - u_h)\|_{\mathcal{N}_h(\partial\Omega_h)}^2$$

$$(3.99) \quad \lesssim h^{2p+1-2j} \|u\|_{H^{p+1}(\Omega)}^2 + h^{1-2j} \|\nabla e\|_{\mathcal{N}_h(\partial\Omega_h)}^2$$

which leads to

$$(3.100) \quad \delta_h^{2j} \|D_{\nu_h}^j e\|_{\partial\Omega_h}^2 \lesssim h^{-1} (\delta_h/h)^{2j} h^{2(p+1)} \|u\|_{H^{p+1}(\Omega)}^2 + h (\delta_h/h)^{2j} \|\nabla e\|_{\mathcal{N}_h(\partial\Omega_h)}^2$$

$$(3.101) \quad \lesssim (h^{-3} \delta_h^2) h^{2(p+1)} \|u\|_{H^{p+1}(\Omega)}^2 + (h^{-3} \delta_h^2) h^2 \|\nabla e\|_{\mathcal{N}_h(\partial\Omega_h)}^2,$$

where we used (2.10) and the fact  $\delta_h/h^2 \lesssim 1$ . Thus we have

(3.102)

$$\|T_{1,k}(e)\|_{\partial\Omega_h} \lesssim \sum_{j=1}^k \delta_h^j \|D_{\nu_h}^j e\|_{\partial\Omega_h} \lesssim (h^{-3/2} \delta_h) \left( h^{p+1} \|u\|_{H^{p+1}(\Omega)} + h \|e\|_h \right).$$

*Conclusion of the proof.* Collecting the bounds (3.69), (3.88), and (3.93), of Terms I, II, and III, we obtain

$$\begin{aligned} (3.103) \quad \|e\|_{\Omega_h} &\lesssim h \|e\|_h \\ &\quad + h^{p+1} \|u\|_{H^{p+1}(\Omega)} \\ &\quad + \delta_h^{k+1} \sup_{0 \leq t \leq \delta_0} \|D^{k+1} u\|_{\partial\Omega_t} \\ &\quad + \delta_h^{l+3/2} \sup_{0 \leq t \leq \delta_0} \|D_n^l (f + \Delta u)\|_{\partial\Omega_t} \end{aligned}$$

which together with the energy norm error estimate (3.49) concludes the proof.  $\square$

**Theorem 3.3.** *The following estimates hold:*

$$(3.104) \quad \|\nabla e\|_{\Omega} \lesssim h^p \|u\|_{H^{p+1}(\Omega)} + \|e\|_h$$

and

$$(3.105) \quad \|e\|_{\Omega} \lesssim h^{p+1} \|u\|_{H^{p+1}(\Omega)} + \|e\|_{\Omega_h} + h \|e\|_h.$$

*Proof.* Adding and subtracting an interpolant, using the interpolation estimate (3.41), and the inverse inequality (3.20) or (3.21), we obtain, for  $m = 0, 1$ ,

(3.106)

$$\begin{aligned} \|\nabla^m e\|_{\Omega \setminus \Omega_h} &\lesssim \|\nabla^m (u - \pi_h u)\|_{\Omega \setminus \Omega_h} + \|\nabla^m (\pi_h u - u_h)\|_{\Omega \setminus \Omega_h} \\ (3.107) \quad &\lesssim h^{p+1-m} \|u\|_{H^{p+1}(\Omega)} + \|\nabla^m (\pi_h u - u_h)\|_{\Omega_h} + h^{1-m} \|\pi_h u - u_h\|_{j_h} \end{aligned}$$

$$(3.108) \quad \lesssim h^{p+1-m} \|u\|_{H^{p+1}(\Omega)} + \|\nabla^m e\|_{\Omega_h} + h^{1-m} \|e\|_{j_h}$$

$$(3.109) \quad \lesssim h^{p+1-m} \|u\|_{H^{p+1}(\Omega)} + \|\nabla^m e\|_{\Omega_h} + h^{1-m} \|e\|_h$$

which concludes the proof.  $\square$



*Remark 3.3.* We conclude from Theorems 3.1 and 3.2 that the precise convergence of the scheme depends on a balance between how well  $\Omega_h$  approximates  $\Omega$  and how many terms are considered in the Taylor series. A poor accuracy in  $\Omega_h$  can be compensated for by increasing the number of Taylor terms. For instance if the domain approximation is no better than  $\delta_h = o(h)$ ,  $k = 1$  is needed for optimality, even if piecewise affine approximation is used for  $u_h$ . In Tables 1 and 2 we detail the asymptotics of the different error contribution for the important case where  $\delta_h = O(h^2)$ , corresponding to a piecewise affine approximation of the boundary.

TABLE 1. The order of the terms in the energy error estimate under the assumption  $\delta_h \lesssim h^2$ . We conclude that we obtain optimal order of convergence for  $p = 2, 3$ , with one term,  $k = 1$ , in the Taylor expansion and for  $p = 4, 5$ , with two terms,  $k = 2$ .

$p$	$h^p$	$k$	$h^{-1/2}\delta_h^{k+1}$	1	$h^{1/2}\delta_h^{l+1}$
1	$h^1$	0	$h^{1.5}$	0	$h^{2.5}$
2	$h^2$	1	$h^{3.5}$	1	$h^{4.5}$
3	$h^3$	2	$h^{5.5}$	2	$h^{6.5}$
4	$h^4$	3	$h^{7.5}$	3	$h^{8.5}$

TABLE 2. The order of the terms in the  $L^2$ -error estimate under the assumption that  $\delta_h \lesssim h^2$ . We conclude that we obtain optimal order of convergence for  $p = 2, 3$ , with one term,  $k = 1$ , in the Taylor expansion and for  $p = 4, 5$ , with two terms,  $k = 2$ .

$p$	$h^{p+1}$	$k$	$\delta_h^{k+1}$	1	$\delta_h^{l+3/2}$
1	$h^2$	0	$h^2$	0	$h^3$
2	$h^3$	1	$h^4$	1	$h^5$
3	$h^4$	2	$h^6$	2	$h^7$
4	$h^5$	3	$h^8$	3	$h^9$

*Remark 3.4.* If for a given  $p$  the lowest values of  $k$  and  $l$  are chosen so that optimal convergence is obtained, it is straightforward to use a trace inequality (see (3.18) and (3.19)) to show that

$$\|u - u_h\|_{(\Omega_h)} + h\|u - u_h\| \lesssim h^{p+1}(\|f\|_{H^{p-1}(\Omega)} + \|u\|_{H^{p+1}(\Omega)}).$$

Therefore the regularities required for optimality of the consistency error of the boundary approximation are always optimal compared to the polynomial approximation.

*Remark 3.5.* We note that we obtain, as a special case, optimal order error estimates for the standard cut Nitsche method with approximate domains by assuming  $k = 0$  and

$$(3.110) \quad \delta_h \lesssim h^{p+1/2}$$

for the energy norm estimate and

$$(3.111) \quad \delta_h \lesssim h^{p+1}$$

for the  $L^2$ -norm estimate. The latter assumption is comparable with the geometric approximation accuracy achieved by standard isoparametric finite elements of order  $p$ .

#### 4. NUMERICAL EXAMPLES

In the numerical examples, we use implicitly defined boundaries by use of zero isolines to predefined functions. Two examples have been considered, one with both convex and concave boundaries, so that cut elements can have parts outside the actual domain, and one example with nonzero boundary conditions where we also compare setting the boundary condition on the exact boundary to setting them on computational boundary. In all examples the stabilization parameters were set to  $\gamma_j = 1/10$ ,  $\beta = 100$ .

**4.1. Convex and concave boundaries.** In our first example we consider a ring-shaped domain. In Figure 1 we show the zero isoline of the function  $\phi = (R - 1/4)(R - 3/4)$ ,  $R = \sqrt{x^2 + y^2}$ , used to implicitly define the domain, and the resulting mesh after removing the cut part. On this ring, we used a load corresponding to the exact solution being a square function in  $R$ ,

$$(4.1) \quad u = 20(3/4 - R)(R - 1/4),$$

with zero boundary conditions on the outside as well as inside boundaries. The elements on the inside of the ring are partially outside the computational domain; outside the domain the load was extended by zero and the exact solution (in the convergence study) by (4.1).

We show an elevation of the approximate solution on one of the meshes in a sequence in Figure 2. In Figures 3 and 4 we show the convergence rates obtained using the symmetric method (2.33)–(2.34) for  $P^2$  and  $P^3$  elements (polynomial orders  $p = 2$  and  $p = 3$ ), respectively. We also show the suboptimal convergence rates of the original Nitsche method. Note in particular that the optimal rate is attained also for  $p = 3$  even though only the first two terms in the Taylor series are accounted for.

**4.2. Nonzero boundary conditions.** The domain for the second example lies inside the ellipse defined by the zero isoline to  $\phi = x^2/(3/4)^2 + y^2/(1/2)^2 - 1$ . In Figure 5 we show the zero isoline of this function and the resulting mesh after removing the cut part. On this domain we use the right-hand side

$$f = \pi^2 \cos(\pi x/2) \cos(\pi y/2)$$

corresponding to the exact solution  $u = \cos(\pi x/2) \cos(\pi y/2)$ . This function also defines the boundary conditions on the cut boundary. An elevation of an approximate solution on one of the meshes in a sequence is given in Figure 6.

In Figure 7 we show the observed  $L^2$  convergence with a  $P^3$  approximation using four different approaches:

- The symmetric method (2.33)–(2.34).
- The unsymmetric Taylor expansion with two terms.
- The unsymmetric Taylor expansion with three terms.

- Prescribing the boundary condition on the cut boundary (using the fact that the exact solution is known).

In all cases the rate of convergence is 4, which is optimal. The error constant is slightly better if we prescribe the boundary condition on the cut boundary, which is to be expected since this does not introduce any approximations of the boundary condition. The difference between the other three methods is negligible.

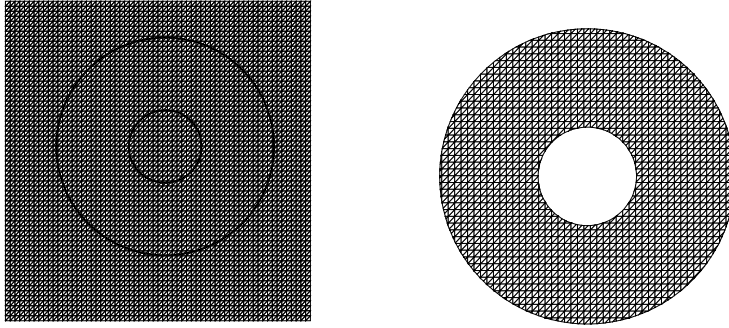


FIGURE 1. Background mesh with the boundary of  $\Omega$  indicated, and the corresponding computational mesh.

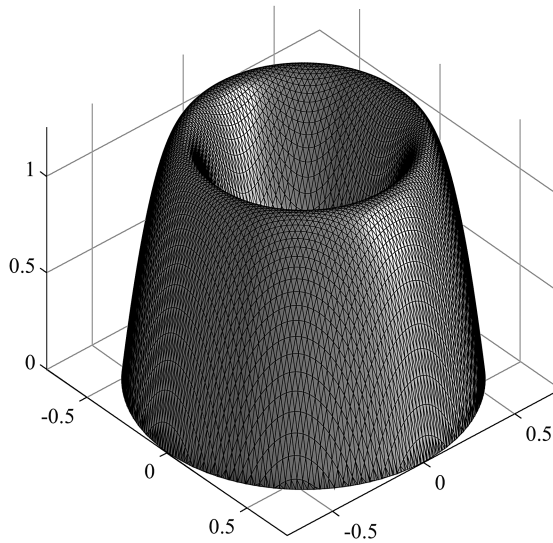


FIGURE 2. Elevation of the approximate solution on one of the meshes in a sequence.

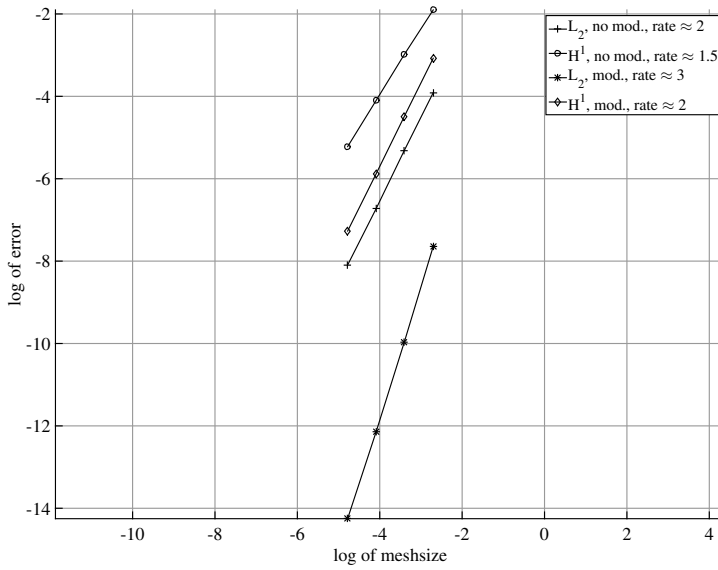


FIGURE 3. Convergence using  $P^2$  elements, symmetric form (log denotes the natural logarithm)

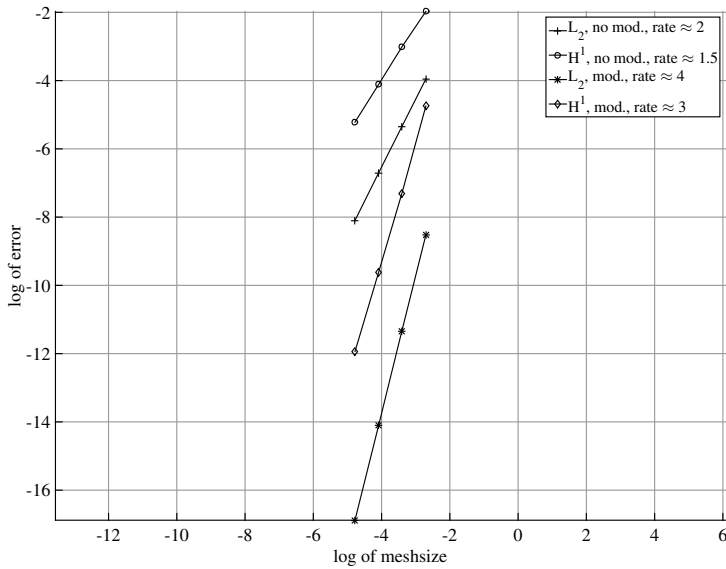


FIGURE 4. Convergence using  $P^3$  elements, symmetric form (log denotes the natural logarithm).

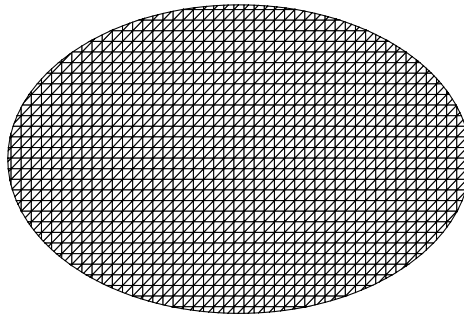
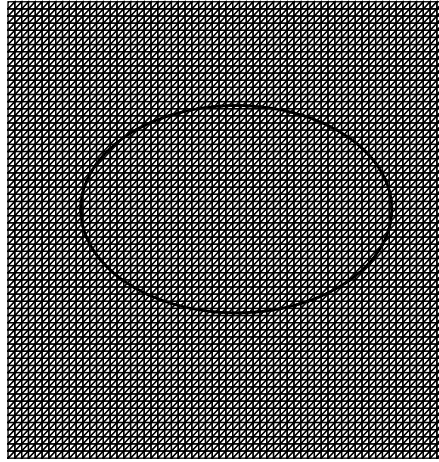


FIGURE 5. Background mesh with the boundary of  $\Omega$  indicated, and the corresponding computational mesh.

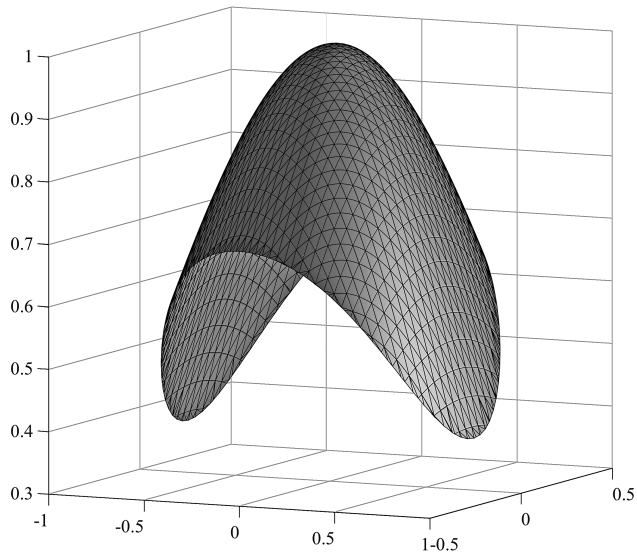


FIGURE 6. Elevation of the approximate solution on one of the meshes in a sequence.

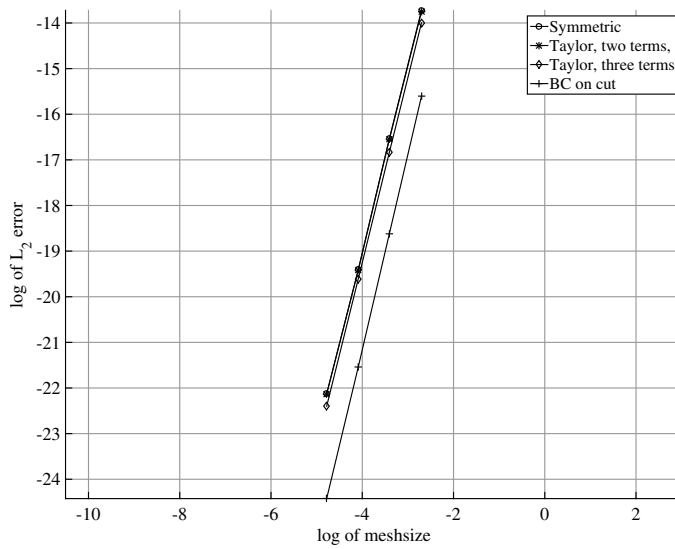


FIGURE 7. Convergence in  $L^2$  for four different approaches (log denotes the natural logarithm).

APPENDIX: VERIFICATION OF SOME ESTIMATES

**Estimates (3.85) and (3.86).** We first note that for each  $x \in U_\delta(\Gamma)$ ,  $0 < \delta \leq \delta_0$ , we have the representation

$$(A.1) \quad \phi(x) = \phi(p(x)) + \int_0^1 \nabla\phi(sx + (1-s)p(x)) \cdot (x - p(x)) ds.$$

Using the Cauchy-Schwarz inequality we obtain

$$(A.2) \quad |\phi(x)|^2 \lesssim |\phi(p(x))|^2 + \delta \|n \cdot \nabla\phi\|_{I_{x,p(x)}}^2$$

$$(A.3) \quad \lesssim |\phi(p(x))|^2 + \delta \|n \cdot \nabla\phi\|_{I_\delta(p(x))}^2,$$

where  $I_{x,p(x)}$  is the line segment between  $x$  and  $p(x)$ , and  $I_\delta(p(x))$  is the line segment between the points  $p(x) \pm \delta n(p(x))$ .

(3.85). We recall that  $\partial\Omega_h \subset U_{\delta_h}(\partial\Omega)$ . Setting  $\delta = \delta_h$  in (A.3) and integrating over  $\partial\Omega_h$  we obtain

$$(A.4) \quad \|\phi\|_{\partial\Omega_h}^2 \lesssim \int_{\partial\Omega_h} |\phi \circ p(x)|^2 dx + \int_{\partial\Omega_h} \delta_h \|n \cdot \nabla\phi\|_{I_{\delta_h}(p(x))}^2 dx$$

$$(A.5) \quad \lesssim \int_{\partial\Omega} |\phi(y)|^2 dy + \int_{\partial\Omega} \delta_h \|n \cdot \nabla\phi\|_{I_{\delta_h}(y)}^2 dy$$

$$(A.6) \quad \lesssim \|\phi\|_{\partial\Omega}^2 + \delta_h \|n \cdot \nabla\phi\|_{U_{\delta_h}(\partial\Omega)}^2,$$

where we first changed the domain of integration from  $\partial\Omega_h$  to  $\partial\Omega$  and then from the tubular coordinates to the Euclidian coordinates. This concludes the proof of (3.85), observing that where it is applied  $\phi \in H_0^1(\Omega)$ , so that  $\|\phi\|_{\partial\Omega}^2 = 0$ .

(3.86). Again setting  $\delta = \delta_h$  in (A.3) and integrating over  $I_{\delta_h}(y)$ , with  $y = p(x) \in \partial\Omega$ , we obtain

$$(A.7) \quad \|\phi\|_{I_{\delta_h}(y)}^2 \lesssim \delta_h |\phi(y)|^2 + \delta_h^2 \|n \cdot \nabla\phi\|_{I_{\delta_h}(y)}^2.$$

Using appropriate changes of coordinates we obtain

$$(A.8) \quad \|\phi\|_{U_{\delta_h}(\partial\Omega)}^2 \lesssim \int_{\partial\Omega} \|\phi\|_{I_{\delta_h}(y)}^2 dy$$

$$(A.9) \quad \lesssim \int_{\partial\Omega} \delta_h |\phi(y)|^2 dy + \int_{\partial\Omega} \delta_h^2 \|n \cdot \nabla\phi\|_{I_{\delta_h}(y)}^2 dy$$

$$(A.10) \quad \lesssim \delta_h \|\phi\|_{\partial\Omega}^2 + \delta_h^2 \|n \cdot \nabla\phi\|_{U_{\delta_h}(\partial\Omega)}^2$$

which proves (3.86).

**Estimate (3.16).** We shall prove that

$$(A.11) \quad \|v\|_{\Omega_h \setminus \Omega} \lesssim h^{1/2} \delta_h^{1/2} \|v\|_h, \quad \forall v \in V_h.$$

Let  $x \in \partial\Omega_h \setminus \Omega$ , i.e.,  $x$  belongs to the part of  $\partial\Omega_h$  that reside outside of  $\Omega$ . For  $y \in I_{x,p(x)}$  we have the representation formula

$$(A.12) \quad v(y) = v(x) + \int_0^1 \nabla v(sy + (1-s)x) \cdot (y - x) ds.$$

Estimating the right-hand side using the Cauchy-Schwarz inequality we obtain

$$(A.13) \quad v^2(y) \lesssim v^2(x) + \left( \int_0^1 \nabla v(sy + (1-s)x) \cdot (y-x) ds \right)^2$$

$$(A.14) \quad \lesssim v^2(x) + |y-x| \|\nabla v\|_{I_{x,y}}^2$$

$$(A.15) \quad \lesssim v^2(x) + \delta_h \|\nabla v\|_{I_{x,p(x)}}^2$$

which leads to

$$(A.16) \quad \|v\|_{I_{x,p(x)}}^2 = \int_{I_{x,p(x)}} v^2(y) dy$$

$$(A.17) \quad \lesssim \int_{I_{x,p(x)}} \left( v^2(x) + \delta_h \|\nabla v\|_{I_{x,p(x)}}^2 \right) dy$$

$$(A.18) \quad \lesssim \delta_h v^2(x) + \delta_h^2 \|\nabla v\|_{I_{x,p(x)}}^2.$$

Integrating over the parts of  $\partial\Omega_h$  that reside outside of  $\Omega$  we obtain

$$(A.19) \quad \|v\|_{\Omega_h \setminus \Omega}^2 \lesssim \int_{\partial\Omega_h \setminus \Omega} \int_{I_{x,p(x)}} v^2(y) dy dx$$

$$(A.20) \quad \lesssim \int_{\partial\Omega_h \setminus \Omega} \|v\|_{I_{x,p(x)}}^2 dx$$

$$(A.21) \quad \lesssim \delta_h \|v\|_{\partial\Omega_h}^2 + \delta_h^2 \|\nabla v\|_{\Omega_h \setminus \Omega}^2$$

$$(A.22) \quad \lesssim \delta_h h (h^{-1} \|v\|_{\partial\Omega_h}^2) + \delta_h^2 \|\nabla v\|_{\Omega_h}^2$$

$$(A.23) \quad \lesssim (\delta_h h + \delta_h^2) \|v\|_h^2$$

and thus (A.11) follows since we assume that  $\delta_h = O(h)$  and therefore  $h\delta_h + \delta_h^2 \lesssim h\delta_h$ .

## REFERENCES

- [1] S. Amdouni, M. Moakher, and Y. Renard, *A local projection stabilization of fictitious domain method for elliptic boundary value problems*, Appl. Numer. Math. **76** (2014), 60–75, DOI 10.1016/j.apnum.2013.10.002. MR3131863
- [2] J. Baiges, R. Codina, F. Henke, S. Shahmiri, and W. A. Wall, *A symmetric method for weakly imposing Dirichlet boundary conditions in embedded finite element meshes*, Internat. J. Numer. Methods Engrg. **90** (2012), no. 5, 636–658, DOI 10.1002/nme.3339. MR2913313
- [3] G. R. Barrenea and F. Chouly, *A local projection stabilized method for fictitious domains*, Appl. Math. Lett. **25** (2012), no. 12, 2071–2076, DOI 10.1016/j.aml.2012.04.020. MR2967791
- [4] R. Becker, E. Burman, and P. Hansbo, *A Nitsche extended finite element method for incompressible elasticity with discontinuous modulus of elasticity*, Comput. Methods Appl. Mech. Engrg. **198** (2009), no. 41–44, 3352–3360, DOI 10.1016/j.cma.2009.06.017. MR2571349
- [5] J. H. Bramble, T. Dupont, and V. Thomée, *Projection methods for Dirichlet’s problem in approximating polygonal domains with boundary-value corrections*, Math. Comp. **26** (1972), 869–879. MR0343657
- [6] E. Burman, *Ghost penalty*, C. R. Math. Acad. Sci. Paris **348** (2010), no. 21–22, 1217–1220, DOI 10.1016/j.crma.2010.10.006. MR2738930
- [7] E. Burman, S. Claus, P. Hansbo, M. G. Larson, and A. Massing, *CutFEM: discretizing geometry and partial differential equations*, Internat. J. Numer. Methods Engrg. **104** (2015), no. 7, 472–501, DOI 10.1002/nme.4823. MR3416285
- [8] E. Burman and P. Hansbo, *Fictitious domain finite element methods using cut elements: II. A stabilized Nitsche method*, Appl. Numer. Math. **62** (2012), no. 4, 328–341, DOI 10.1016/j.apnum.2011.01.008. MR2899249



- [9] B. Cockburn, W. Qiu, and M. Solano, *A priori error analysis for HDG methods using extensions from subdomains to achieve boundary conformity*, *Math. Comp.* **83** (2014), no. 286, 665–699, DOI 10.1090/S0025-5718-2013-02747-0. MR3143688
- [10] B. Cockburn and M. Solano, *Solving Dirichlet boundary-value problems on curved domains by extensions from subdomains*, *SIAM J. Sci. Comput.* **34** (2012), no. 1, A497–A519, DOI 10.1137/100805200. MR2890275
- [11] B. Cockburn and M. Solano, *Solving convection-diffusion problems on curved domains by extensions from subdomains*, *J. Sci. Comput.* **59** (2014), no. 2, 512–543, DOI 10.1007/s10915-013-9776-y. MR3188451
- [12] R. Codina and J. Baiges, *Approximate imposition of boundary conditions in immersed boundary methods*, *Internat. J. Numer. Methods Engrg.* **80** (2009), no. 11, 1379–1405, DOI 10.1002/nme.2662. MR2582494
- [13] T. Dupont,  *$L_2$  error estimates for projection methods for parabolic equations in approximating domains*, in C. de Boor, editor, *Mathematical Aspects of Finite Elements in Partial Differential Equations*, Academic Press, New York, 1974, pp. 313–352.
- [14] A. Düster, J. Parvizian, Z. Yang, and E. Rank, *The finite cell method for three-dimensional problems of solid mechanics*, *Comput. Methods Appl. Mech. Engrg.* **197** (2008), no. 45–48, 3768–3782, DOI 10.1016/j.cma.2008.02.036. MR2458114
- [15] G. B. Folland, *Introduction to Partial Differential Equations*, 2nd ed., Princeton University Press, Princeton, NJ, 1995. MR1357411
- [16] D. Gilbarg and N. S. Trudinger, *Elliptic Partial Differential Equations of Second Order*, *Classics in Mathematics*, Springer-Verlag, Berlin, 2001. Reprint of the 1998 edition. MR1814364
- [17] A. Hansbo, P. Hansbo, and M. G. Larson, *A finite element method on composite grids based on Nitsche’s method*, *M2AN Math. Model. Numer. Anal.* **37** (2003), no. 3, 495–514, DOI 10.1051/m2an:2003039. MR1994314
- [18] J. Haslinger and Y. Renard, *A new fictitious domain approach inspired by the extended finite element method*, *SIAM J. Numer. Anal.* **47** (2009), no. 2, 1474–1499, DOI 10.1137/070704435. MR2497337
- [19] A. Johansson and M. G. Larson, *A high order discontinuous Galerkin Nitsche method for elliptic problems with fictitious boundary*, *Numer. Math.* **123** (2013), no. 4, 607–628, DOI 10.1007/s00211-012-0497-1. MR3032318
- [20] A. Massing, M. G. Larson, and A. Logg, *Efficient implementation of finite element methods on nonmatching and overlapping meshes in three dimensions*, *SIAM J. Sci. Comput.* **35** (2013), no. 1, C23–C47, DOI 10.1137/11085949X. MR3033077
- [21] A. Massing, M. G. Larson, A. Logg, and M. E. Rognes, *A stabilized Nitsche fictitious domain method for the Stokes problem*, *J. Sci. Comput.* **61** (2014), no. 3, 604–628, DOI 10.1007/s10915-014-9838-9. MR3268662
- [22] J. Nitsche, *Über ein Variationsprinzip zur Lösung von Dirichlet-Problemen bei Verwendung von Teilräumen, die keinen Randbedingungen unterworfen sind* (German), *Abh. Math. Sem. Univ. Hamburg* **36** (1971), 9–15. MR0341903
- [23] J. Parvizian, A. Düster, and E. Rank, *Finite cell method: h- and p-extension for embedded domain problems in solid mechanics*, *Comput. Mech.* **41** (2007), no. 1, 121–133, DOI 10.1007/s00466-007-0173-y. MR2377802
- [24] L. R. Scott and S. Zhang, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, *Math. Comp.* **54** (1990), no. 190, 483–493, DOI 10.2307/2008497. MR1011446
- [25] M. Tur, J. Albelda, E. Nadal, and J. J. Ródenas, *Imposing Dirichlet boundary conditions in hierarchical Cartesian meshes by means of stabilized Lagrange multipliers*, *Internat. J. Numer. Methods Engrg.* **98** (2014), no. 6, 399–417, DOI 10.1002/nme.4629. MR3195261

DEPARTMENT OF MATHEMATICS, UNIVERSITY COLLEGE LONDON, LONDON, UK-WC1E 6BT, UNITED KINGDOM

DEPARTMENT OF MECHANICAL ENGINEERING, JÖNKÖPING UNIVERSITY, S-551 11 JÖNKÖPING, SWEDEN

DEPARTMENT OF MATHEMATICS AND MATHEMATICAL STATISTICS, UMEÅ UNIVERSITY, SE-90187 UMEÅ, SWEDEN

*E-mail address:* mats.larson@math.umu.se