

# Optimal error estimates for a Discontinuous Galerkin method on curved boundaries

Adérito Araújo<sup>1</sup> and Milene Santos<sup>1\*</sup>

<sup>1</sup>Centre for Mathematics, University of Coimbra, 3000-143, Coimbra, Portugal.

\*Corresponding author(s). E-mail(s): [milene@mat.uc.pt](mailto:milene@mat.uc.pt);  
Contributing authors: [alma@mat.uc.pt](mailto:alma@mat.uc.pt);

## Abstract

We consider a discontinuous Galerkin method to solve boundary value problems in curved boundary domains in two-dimensional. The question that arises concerns the reduction of the order of convergence of numerical methods when considering the approximation of the domain by a polygonal mesh. Unless the boundary conditions can be accurately transferred from the physical boundary to the computational boundary, the isoparametric element method is usually employed to recover the optimal convergence orders. However, this technique involves more complex algebra and additional computational costs when compared to approaches using polygonal meshes, which are widely used due to their simplicity in many applications. In this paper, we present and analyse a higher-order strategy that achieves the optimal convergence order on polygonal approximations of domains with smooth boundaries. The boundary approximation error is corrected by means of polynomial reconstructions of the boundary conditions. We present a study on the existence and uniqueness of the solution and derive error estimates for a two-dimensional linear reaction-diffusion boundary-value problem with homogeneous Dirichlet boundary conditions in convex and non-convex domains. We prove that the numerical solution exhibits an optimal convergence rate under certain regularity conditions on the solution. A numerical benchmark is provided to illustrate the theoretical results proven in this work.

**Keywords:** Arbitrary curved boundaries, Discontinuous Galerkin method, Reconstruction for off-site data method, Error estimate

**MSC Classification:** 65N12 , 65N15 , 65N30

# 1 Introduction

In this work, we present an approach for solving boundary-value problems posed in a curved boundary domain of arbitrary shape in the context of discontinuous Galerkin (DG) methods. The study of boundary value problems in curved boundary domains is a subject of growing interest in the numerical analysis community. One of the major problems is the reduction in the order of convergence of numerical methods when considering the approximation of the domain by a polygonal mesh. In particular, the DG solutions are highly sensitive to the accuracy of approximations of the curved boundaries [5]. It has been shown that given homogeneous Dirichlet boundary conditions on a curved boundary domain  $\Omega$ , if these conditions are imposed on the polygonal domain  $\Omega_h$ , any finite element method will be at most second-order accurate [31]. This highlights the importance of the boundary condition treatment since the errors in the boundary may pollute the solution inside the domain.

Over the past few decades, several techniques have been developed to remedy this loss of accuracy. There are two main strategies to address this issue. The isoparametric finite element method [5] and the isogeometric analysis [16] aim to reduce the geometric error without modifying the variational form. Therefore this technique requires the construction of a mesh with curved elements on the boundary, which is a challenging geometric problem where ineligible cells can be produced. Moreover, this approach also raises some numerical challenges since it considers non-constant Jacobian transformations from the reference element.

Another strategy considers a polygonal approximation domain  $\Omega_h$  and focuses on a modified variational formulation. There has been a growing body of research focused on correcting the error that results from the approximation of the physical boundary  $\partial\Omega$  by a polygonal boundary  $\partial\Omega_h$ , by modifying the boundary condition. In [22], the authors consider a computational polygonal domain in place of the physical domain and modify the normal vector involved in the wall boundary condition. However, this method can only be formulated for slip-wall boundary conditions and the work is limited to 2D geometries. In [33], the author proposes a modified DG scheme defined on polygonal meshes that avoids integrals inside curved elements. However, integrations along boundary curve segments are still necessary. This approach was extended to solving three-dimensional Euler equations and it was simplified by considering the relation between the normal vector of the computational domain and the surface Jacobian [32]. In the Shifted Boundary Method (SBM), the location where the boundary conditions are applied is shifted from the true boundary to an approximate (surrogate) boundary. The value of boundary conditions is modified by means of Taylor expansion, in order to reflect this displacement (see [4] and the references therein).

In [30] we developed a strategy called DG-ROD (Reconstruction for Off-site Data) method, which is based on a polynomial reconstruction of the boundary condition imposed on the computational domain. The main advantage of this approach relies on the use of polygonal meshes without losing the accuracy of the method by considering polynomial reconstructions to correct the error resulting from the approximation of the curved boundary with a polygonal boundary. The ROD method has been proposed in the context of the finite volume (FV) method [10–14] and it has been later extended for the finite difference (FD) method on Cartesian grids [9]. Despite the numerical

evidence in the context of the FV, FD, and DG methods, there is no theoretical evidence on the proof of the convergence of the method. Thus, the main contribution of this work is to establish error estimates for a two-dimensional linear reaction-diffusion problem with homogeneous Dirichlet boundary conditions concerning the DG-norm and the  $L^2$ -norm, and hence, fill a theoretical gap in the analysis of the DG-ROD method for boundary value problems. The overall DG-ROD method can be obtained by considering two different approaches: we can consider an iterative procedure of the DG method and the polynomial reconstruction or we can consider a global system where we only have to solve the problem once. In this work, we address the last approach.

This document is organized as follows. Section 2 is devoted to introductory concepts related to mesh notations and the space of discontinuous functions, and the formulation of the problem to be considered. In Section 3, we start by analysing some basic properties of the method, namely, we show the boundedness of the bilinear form and we prove a weak coercivity. Moreover, we present a study on the existence and uniqueness of the solution for the reaction-diffusion problem with homogeneous Dirichlet boundary conditions, following the work developed within the framework of the classical finite element method [26, 27]. The core of this work is represented by Section 4, where we derive error estimates for the method introduced in this Section 2 for convex and non-convex domains. For the first case, we prove that the DG-ROD solution exhibits an optimal  $\mathcal{O}(h^{N+1})$  convergence rate in the  $L^2$ -norm when  $N$ -degree piecewise polynomials are used, under certain regularity conditions on the solution. Finally, the numerical experiments and results are reported in Section 5. In Section 6, we summarize the results and present some final comments and perspectives for future work. The last part of the paper is an appendix that contains technical results and upper bounds estimates used in the analysis of the method.

## 2 The DG-ROD Method

This section addresses the DG-ROD formulation for a two-dimensional linear boundary-value problem on a curved boundary domain, which is discretised with piecewise linear elements. This method has the advantage of overcoming the difficulties inherent to curved mesh approaches by discretising the physical domain with polygonal meshes constructed from the conventional meshing algorithms, where piecewise linear elements approximate the arbitrary curved boundary. The main idea of the DG method is based on the use of discontinuous functions to obtain an approximate solution. Additionally, the DG-ROD method employs specific polynomial reconstructions for the prescribed boundary conditions on the physical boundary. In order to allow an easier description of our methodology, thereby avoiding non-essential technical details, we consider a two-dimensional reaction-diffusion problem. The first step is the definition of the mesh and the broken polynomial spaces. After providing an overview of some basic ideas related to computational meshes, we present the primal formulation of the method, which incorporates the modification derived from the polynomial reconstruction of the boundary conditions.

## 2.1 Model problem

The methodology for dealing with curved boundary domains studied in this work can be applied to different equations. However, to avoid non-essential technical details, we consider the reaction-diffusion equation in a two-dimensional physical domain  $\Omega$  with arbitrary smooth curved physical boundary  $\partial\Omega$ , considering the Cartesian coordinate system  $\mathbf{x} = (x, y)$ . We seek function  $u = u(\mathbf{x})$ , solution of the reaction-diffusion problem

$$-\Delta u(\mathbf{x}) + c(\mathbf{x}) u(\mathbf{x}) = f(\mathbf{x}), \quad \mathbf{x} \in \Omega, \quad (1)$$

$$u(\mathbf{x}) = 0, \quad \mathbf{x} \in \partial\Omega, \quad (2)$$

where  $c \in C(\bar{\Omega})$ , with  $c(\mathbf{x}) \geq 0$ , for  $\mathbf{x} \in \bar{\Omega}$ , and  $f \in L^2(\Omega)$ . The Lebesgue space  $L^2(\Omega)$  is defined as a space of measurable functions  $u : \Omega \rightarrow \mathbb{R}$  such that  $\|u\|_{L^2(\Omega)}^2 < +\infty$ , equipped with norm  $\|u\|_{L^2(\Omega)}^2 = (u, u)_{L^2(\Omega)}$  and inner product

$$(u, w)_{L^2(\Omega)} = \int_{\Omega} u(\mathbf{x})w(\mathbf{x}) \, d\mathbf{x}.$$

## 2.2 Definition of the mesh

The physical domain  $\Omega$  is meshed with  $K$  non-overlapping straight-sided triangles  $T^k, k = 1, \dots, K$ , leading to an approximate computational domain  $\Omega_h$  given as

$$\Omega_h = \bigcup_{k=1}^K T^k. \quad (3)$$

The triangulation  $\mathcal{T}_h = \{T^k, k = 1, \dots, K\}$  is assumed to be conformed where the intersection of two elements is either a complete edge, a vertex, or the empty set. We assume that no element  $T^k$  has more than one edge on  $\partial\Omega_h$  and all the vertexes of the polygon lie on  $\partial\Omega$ . The space parameter  $h$  represents the maximum element diameter, namely

$$h = \max_{T^k \in \mathcal{T}_h} \{h_k\}, \quad h_k = \sup_{P_1, P_2 \in T^k} \|P_1 - P_2\|.$$

The triangulation is also assumed to be regular [19] in the sense that there is a constant  $\rho > 0$  such that

$$\forall T^k \in \mathcal{T}_h, \quad \frac{h_k}{\rho_k} \leq \rho, \quad (4)$$

where  $\rho_k$  denotes the maximum radius of a ball inscribed in  $T^k$ .

Let  $\mathcal{E}_h$  denote all edges of elements in  $\mathcal{T}_h$  and  $\mathcal{E}_0$  denotes all interior edges. We assume that exists a positive constant  $\mu$  such that for every element  $T^k \in \mathcal{T}_h$  and  $e \in \mathcal{E}_h \cap \partial T^k$ , we have [24]

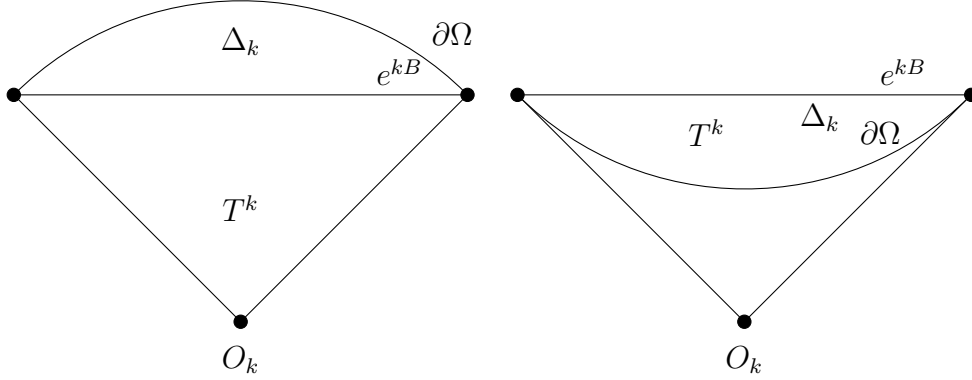
$$\mu h_k \leq h_e, \quad (5)$$

where  $h_e$  denotes the length of the edge  $e$ . Assume that exists a positive constant  $\tilde{\rho}$ , such that

$$\frac{h}{h_{min}} \leq \tilde{\rho}, \quad (6)$$

where  $h_{min} = \min_{T^k \in \mathcal{T}_h} \{h_k^{min}\}$ , with  $h_k^{min} = \inf_{P_1, P_2 \in T^k} \|P_1 - P_2\|$ .

Given an element  $T^k$ , denote as  $I^k$  the index set of the elements  $T^\ell$  that share a common edge  $e^{k\ell}$  and by  $I^B$  the index set of elements which have an edge on the boundary,  $e^{kB}$ . Normal vector  $\mathbf{n}^{k\ell}$ ,  $\ell \in I^k$ , is pointed outward of element  $T^k$  and  $\mathbf{n}^{\ell k} = -\mathbf{n}^{k\ell}$ . For each element  $T^k$ ,  $k \in I^B$ , denote as  $\Delta_k$  the closed set delimited by  $\partial\Omega$  and the edge  $e^{kB}$  (see Figure 1). Consider  $\mathcal{Q}^B$  a subset of  $I^B$  such that  $\mathcal{Q}^B$  denote the index set of elements that have an edge on the boundary and  $T^k \setminus \Omega$  is not restricted to a pair of vertexes of  $\partial\Omega$ .



**Fig. 1:** Element  $T^k$  with an edge  $e^{kB}$  on the computational boundary  $\partial\Omega_h$ , for the convex case where  $T^k \subset \Omega$  (left panel) and for the concave case, where  $T^k \not\subset \Omega$  (right panel).

### 2.3 Space of discontinuous functions

The discontinuous Galerkin method is based on the use of discontinuous approximations. Thus, we introduce the so-called broken Sobolev spaces  $H^l(\mathcal{T}_h)$ , with  $l = 1, 2$ , as

$$H^l(\mathcal{T}_h) = \{w \in L^2(\Omega_h) : w|_{T^k} \in H^l(T^k) \forall T^k \in \mathcal{T}_h\}.$$

Note that

$$(w, v)_{L^2(\Omega_h)} = \sum_{k=1}^K (w, v)_{L^2(T^k)}, \quad \forall w, v \in L^2(\Omega_h),$$

where  $(u, v)_{L^2(T^k)}$  denotes the usual inner product on  $L^2(T^k)$ . For  $w \in H^l(\mathcal{T}_h)$ , with  $l = 1, 2$ , we define the norm

$$\|w\|_{H^l(\mathcal{T}_h)}^2 = \sum_{k=1}^K \|w\|_{H^l(T^k)}^2,$$

where  $\|w\|_{H^l(T^k)}$  denotes the usual  $H^l$ -norm on the element  $T^k$ . We define the space of discontinuous piecewise polynomial functions

$$S_{hN} = \left\{ v \in L^2(\Omega_h) : v|_{T^k} \in \mathcal{P}_N(T^k) \forall T^k \in \mathcal{T}_h \right\},$$

with  $\mathcal{P}_N(T^k)$  denoting the space of polynomials of degree less than or equal to  $N$  in element  $T^k$ . We also introduce broken operators by restriction to each element  $T^k \in \mathcal{T}_h$  as follows:

- The broken gradient operator  $\nabla_h : H^1(\mathcal{T}_h) \rightarrow [L^2(\Omega_h)]^2$  is defined by  $(\nabla_h v)|_{T^k} = \nabla(v|_{T^k})$ , for  $T^k \in \mathcal{T}_h$ ,  $v \in H^1(\mathcal{T}_h)$ .
- The broken divergence operator  $\nabla_h \cdot : [H^1(\mathcal{T}_h)]^2 \rightarrow L^2(\Omega_h)$  is defined by  $(\nabla_h \cdot \mathbf{q})|_{T^k} = \nabla \cdot (\mathbf{q}|_{T^k})$ , for  $T^k \in \mathcal{T}_h$ ,  $\mathbf{q} \in [H^1(\mathcal{T}_h)]^2$ .

Let  $\Gamma = \cup_{T^k \in \mathcal{T}_h} \partial T^k$  and  $\Gamma_0 = \Gamma \setminus \partial\Omega_h$ , the traces of functions in  $H^1(\mathcal{T}_h)$  belong to  $T(\Gamma) = \prod_{T^k \in \mathcal{T}_h} L^2(\partial T^k)$ . Note that  $v$  may be double-valued on  $\Gamma_0$  and is single-valued on  $\partial\Omega_h$ .

We introduce some operators that will be useful for manipulating the numerical fluxes and obtaining the primal formulation. Let  $e^{k\ell}$  be an edge shared by the elements  $T^k$  and  $T^\ell$ . For  $\mathbf{q} \in [T(\Gamma)]^2$  and  $u \in T(\Gamma)$ , we define the averages  $\{\{\mathbf{q}_h\}\}^{k\ell}$  and  $\{\{u_h\}\}^{k\ell}$  and the jumps  $[[\mathbf{q}_h]]^{k\ell}$  and  $[[u_h]]^{k\ell}$  as follows:

$$\begin{aligned} \{\{\mathbf{q}_h\}\}^{k\ell} &= \frac{\mathbf{q}_h^k + \mathbf{q}_h^\ell}{2}, & \{\{u_h\}\}^{k\ell} &= \frac{u_h^k + u_h^\ell}{2}, \\ [[\mathbf{q}_h]]^{k\ell} &= \mathbf{n}^{k\ell} \cdot \mathbf{q}_h^k + \mathbf{n}^{\ell k} \cdot \mathbf{q}_h^\ell, & [[u_h]]^{k\ell} &= \mathbf{n}^{k\ell} u_h^k + \mathbf{n}^{\ell k} u_h^\ell. \end{aligned}$$

For a boundary edge  $e^{kB}$ , we define

$$\{\{\mathbf{q}_h\}\}^{kB} = \mathbf{q}_h^k, \quad \{\{u_h\}\}^{kB} = u_h^k, \quad [[\mathbf{q}_h]]^{kB} = \mathbf{n}^{kB} \cdot \mathbf{q}_h^k, \quad [[u_h]]^{kB} = \mathbf{n}^{kB} u_h^k.$$

When it is clear which edge we are referring to, we usually omit the superscript  $k\ell$  and simply write  $\{\{\cdot\}\}$  and  $[[\cdot]]$ . A convenient norm with which to carry out the analysis of the method is the following [3]

$$\|u\|^2 = \sum_{k=1}^K \left( \|u\|_{H^1(T^k)}^2 + h_k^2 |u|_{H^2(T^k)}^2 \right) + \sum_{e \in \mathcal{E}_h} h_e^{-1} \|[[u]]\|_{L^2(e)}^2, \quad (7)$$

for  $u \in H^2(\mathcal{T}_h)$ . For notational convenience, let

$$|v|_*^2 = \sum_{e \in \mathcal{E}_h} h_e^{-1} \| \llbracket v \rrbracket \|_{L^2(e)}^2, \quad (8)$$

for  $v \in L^2(\Omega_h)$ . Using a inequality ([7], Lemma 4.5.3), we may prove that, for  $T^k \in \mathcal{T}_h$ ,

$$h_k |v|_{H^2(T^k)} \leq C |v|_{H^1(T^k)}. \quad (9)$$

Thus,

$$\begin{aligned} \|v\|^2 &= \sum_{k=1}^K \left( \|v\|_{L^2(T^k)}^2 + |v|_{H^1(T^k)}^2 + h_k^2 |v|_{H^2(T^k)}^2 \right) + |v|_*^2 \\ &\leq (1 + C^2) \left( \|v\|_{H^1(\mathcal{T}_h)}^2 + |v|_*^2 \right). \end{aligned} \quad (10)$$

For each element  $T^k$ , with  $k \in I^B$ , let  $I^{kB}$  be the index set of the discontinuous Galerkin nodes different from the vertexes that belong to the boundary edge  $e^{kB}$  (see left panel of Figure 2).

Now, we introduce two spaces  $\mathcal{V}_h$  and  $\mathcal{W}_h$  associated to  $\mathcal{T}_h$ . The space  $\mathcal{V}_h$  is defined by

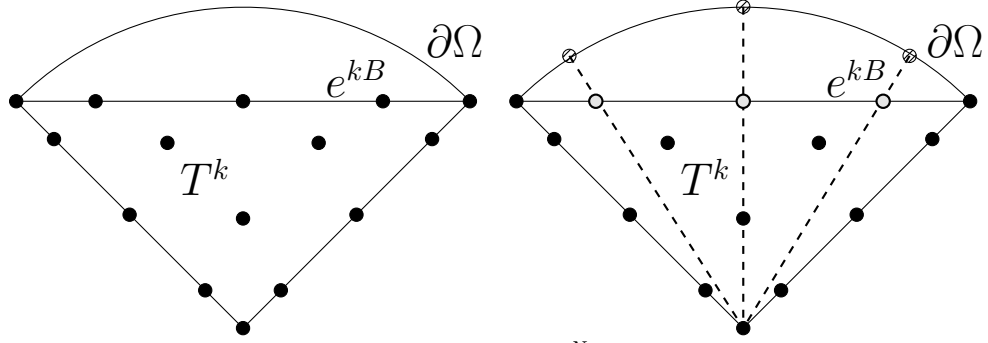
$$\mathcal{V}_h = \left\{ v \in H^2(\mathcal{T}_h) : v|_{\partial\Omega_h} = 0, v|_{T^k} \in \mathcal{P}_N(T^k), \forall T^k \in \mathcal{T}_h \right\}.$$

For convenience, we extend by 0 every function  $v \in \mathcal{V}_h$  to  $\Omega \setminus \Omega_h$ .  $\mathcal{W}_h$  is the space defined in  $\Omega_h$  that satisfies the following properties for  $w \in \mathcal{W}_h$

- (1)  $w|_{T^k} \in \mathcal{P}_N(T^k), \forall T^k \in \mathcal{T}_h$ ;
- (2)  $w \in H^2(\mathcal{T}_h)$ ;
- (3) The expression of  $w$  is extended to  $\Omega \setminus \Omega_h$  in such a way that its polynomial expression in  $T^k, k \in I^B$ , also applies in  $\Delta_k$ ;
- (4)  $w$  vanishes at the vertexes of  $\partial\Omega_h$  and  $w(P_r^k) = 0, r = 1, \dots, N-1, \forall T^k \in \mathcal{T}_h, k \in I^B$  (where each point  $P_r^k$  is chosen to be the nearest intersection with the physical boundary  $\partial\Omega$  of the line passing through the vertex  $O_k$  of  $T^k$  not belonging to  $\partial\Omega$  and one of  $N-1$  discontinuous Galerkin nodes  $\mathbf{x}_i^k, i \in I^{kB}$ , lying on the associated boundary edge,  $e^{kB}$ ). Thus,  $w$  vanishes at  $N+1$  points on  $\partial\Omega$ .

For notation proposes, assume that the vertexes of the element  $T^k, k \in I^B$ , on  $\partial\Omega$  are denoted by  $P_N^k$  and  $P_{N+1}^k$ . Thus, according to property (4), we may write  $w(P_r^k) = 0, r = 1, \dots, N+1, \forall T^k \in \mathcal{T}_h, k \in I^B$ . An example of the nodes associated with  $\mathcal{W}_h$  is reported in Figure 2. Namely, the discontinuous Galerkin nodal set and the points  $P_r^k, r = 1, \dots, N-1$ , resulting from a projection of the nodal points lying on the boundary edge  $e^{kB}$ . For the non-convex case, the points  $P_r^k$  are obtained using the same approach.

For each element  $T^k, k \in I^B$ , let  $m_N = N(N+1)/2$  be the number of nodal points that do not lie in the interior of the edge  $e^{kB}$ . In other words,  $m_N = N_p - (N-1)$ , with  $N_p = (N+1)(N+2)/2$ . The next lemma establishes that  $\mathcal{W}_h$  is a non-empty



**Fig. 2:** Discontinuous Galerkin nodal set  $\{\mathbf{x}_i^k\}_{i=1}^{N_p}$  denoted by the black dots (left panel) and points  $P_r^k$ ,  $r = 1, \dots, N-1$ , denoted by the dots with diagonal lines pattern (right panel).

finite-dimensional space and the proof of this result follows the same arguments as in [28].

**Lemma 2.1.** *Let  $\mathcal{P}_N(T^k)$  be the space of polynomials defined in  $T^k$ ,  $k \in I^B$ , of degree less than or equal to  $N$ . Provided  $h$  small enough  $\forall T^k$ ,  $k \in I^B$ , given a set of  $m_N$  real values  $\gamma_i^k$ ,  $i = 1, \dots, m_N$ , there exists a unique function  $w \in \mathcal{P}_N(T^k)$  that vanishes at both vertex of  $T^k$  located on  $\partial\Omega$  and at the points  $P_r^k$  of  $\partial\Omega$ ,  $r = 1, \dots, N-1$ , and takes value  $\gamma_i^k$  respectively at the  $m_N$  nodes of  $T^k$  not located on  $\partial\Omega_h$ .*

## 2.4 Variational formulation

In order to use a mixed formulation, consider the vector function  $\mathbf{q} = (q_x, q_y)^T$  defined as the gradient of  $u$ , i.e.  $\mathbf{q}(\mathbf{x}) = \nabla u(\mathbf{x})$ . Thus, we may write  $\Delta u(\mathbf{x}) = \nabla \cdot \mathbf{q}(\mathbf{x})$ . Replacing this expression in (1), the solution is sought for the equivalent problem

$$-\nabla \cdot \mathbf{q}(\mathbf{x}) + c(\mathbf{x})u(\mathbf{x}) = f(\mathbf{x}), \quad (11)$$

$$\mathbf{q}(\mathbf{x}) = \nabla u(\mathbf{x}). \quad (12)$$

Consider that numerical solution  $u_h$  has the following decomposition

$$u_h(\mathbf{x}) = \bigoplus_{k=1}^K u_h^k(\mathbf{x}) \in \mathcal{W}_h. \quad (13)$$

In each element  $T^k$ , the local solution  $u_h^k$  has a polynomial decomposition with the two-dimensional Lagrange polynomials

$$\mathbf{x} \in T^k \in \mathcal{T}_h : \quad u_h^k(\mathbf{x}) = \sum_{i=1}^{N_p} u_i^k \ell_i^k(\mathbf{x}), \quad (14)$$



where  $u_i^k = u_h^k(\mathbf{x}_i^k)$  are the nodal values of the Lagrange polynomials basis  $\ell_i^k(\mathbf{x})$  at points  $\mathbf{x}_i^k \in T^k$ ,  $i = 1, \dots, N_p$ . Vector  $\mathbf{u}^k = (u_1^k, \dots, u_{N_p}^k)^\top$  gathers the  $N_p$  nodal values. The DG discretisation of vector function  $\mathbf{q}$  is also introduced by taking  $q_{h,x}^k, q_{h,y}^k \in \mathcal{W}_h$  and the auxiliary variable discretisation expressed as  $\mathbf{q}_h^k = (q_{h,x}^k, q_{h,y}^k)^\top$  with

$$\mathbf{x} \in T^k \in \mathcal{T}_h : \quad q_{h,\zeta}^k(\mathbf{x}) = \sum_{i=1}^{N_p} q_{i,\zeta}^k \ell_i^k(\mathbf{x}), \quad \zeta = x, y. \quad (15)$$

Vectors  $\mathbf{q}_\zeta^k = (q_{1,\zeta}^k, \dots, q_{N_p,\zeta}^k)^\top$ ,  $\zeta = x, y$ , gather the nodal values of polynomials  $q_{h,\zeta}^k$ .

A discrete solution  $(u_h, \mathbf{q}_h)$  is sought for (11) and (12) that satisfy

$$-\left(\nabla \cdot \mathbf{q}_h^k, \phi_h^k\right)_{L^2(T^k)} + \left(cu_h^k, \phi_h^k\right)_{L^2(T^k)} = \left(f, \phi_h^k\right)_{L^2(T^k)}, \quad (16)$$

$$\left(\mathbf{q}_h^k, \mathbf{\Pi}_h^k\right)_{L^2(T^k)} - \left(\nabla u_h^k, \mathbf{\Pi}_h^k\right)_{L^2(T^k)} = 0, \quad (17)$$

where  $\phi_h^k = \phi_{h|_{T^k}} \in \mathcal{V}_h$  and  $\mathbf{\Pi}_h^k = \mathbf{\Pi}_{h|_{T^k}} \in [\mathcal{V}_h]^2$ .

Now, if we integrate (16) and (17) by parts, following the same arguments as in [30], and if we add over all the elements of the mesh, we get

$$\left(\mathbf{q}_h, \nabla_h \phi_h\right)_{L^2(\Omega_h)} = -\left(cu_h, \phi_h\right)_{L^2(\Omega_h)} + \left(f, \phi_h\right)_{L^2(\Omega_h)} + \sum_{T^k \in \mathcal{T}_h} \int_{\partial T^k} \mathbf{n}^{k\ell} \cdot \mathbf{q}_h^{*k\ell} \phi_h^{k\ell} ds, \quad (18)$$

$$\left(\mathbf{q}_h, \mathbf{\Pi}_h\right)_{L^2(\Omega_h)} = -\left(u_h, \nabla_h \cdot \mathbf{\Pi}_h\right)_{L^2(\Omega_h)} + \sum_{T^k \in \mathcal{T}_h} \int_{\partial T^k} \mathbf{n}^{k\ell} \cdot \mathbf{\Pi}_h^{k\ell} u_h^{*k\ell} ds, \quad (19)$$

where  $\mathbf{q}_h^{*k\ell} = \mathbf{q}_h^{*\ell k}$  and  $u_h^{*k\ell} = u_h^{*\ell k}$  are symmetric numerical fluxes defined on the interface  $e^{k\ell}$ , and  $\phi_h^{k\ell}$  and  $\mathbf{\Pi}_h^{k\ell}$  are the polynomials  $\phi_h^k$  and  $\mathbf{\Pi}_h^k$  defined on the edge  $e^{k\ell}$ .

Using the average and jump operators, note that ([21], Lemma 7.9)

$$\sum_{T^k \in \mathcal{T}_h} \int_{\partial T^k} u_h^{k\ell} \mathbf{n}^{k\ell} \cdot \mathbf{\Pi}_h^{k\ell} ds = \int_\Gamma \llbracket u_h \rrbracket^{k\ell} \cdot \{ \{ \mathbf{\Pi}_h \} \}^{k\ell} ds + \int_{\Gamma_0} \{ \{ u_h \} \}^{k\ell} \llbracket \mathbf{\Pi}_h \rrbracket^{k\ell} ds. \quad (20)$$

Integrating by parts we get

$$-\int_{\Omega_h} u_h \nabla_h \cdot \mathbf{\Pi}_h d\mathbf{x} = \int_{\Omega_h} \nabla_h u_h \cdot \mathbf{\Pi}_h d\mathbf{x} - \int_\Gamma \llbracket u_h \rrbracket^{k\ell} \cdot \{ \{ \mathbf{\Pi}_h \} \}^{k\ell} ds - \int_{\Gamma_0} \{ \{ u_h \} \}^{k\ell} \llbracket \mathbf{\Pi}_h \rrbracket^{k\ell} ds. \quad (21)$$

Thus, applying identity (20), we can rewrite (18) and (19) as

$$\begin{aligned} (\mathbf{q}_h, \nabla_h \phi_h)_{L^2(\Omega_h)} &= -(cu_h, \phi_h)_{L^2(\Omega_h)} + (f, \phi_h)_{L^2(\Omega_h)} + \int_{\Gamma} \llbracket \phi_h \rrbracket \cdot \{\{\mathbf{q}_h^*\}\} ds \\ &\quad + \int_{\Gamma_0} \{\{\phi_h\}\} \llbracket \mathbf{q}_h^* \rrbracket ds, \end{aligned} \quad (22)$$

$$(\mathbf{q}_h, \mathbf{\Pi}_h)_{L^2(\Omega_h)} = -(u_h, \nabla_h \cdot \mathbf{\Pi}_h)_{L^2(\Omega_h)} + \int_{\Gamma} \llbracket u_h^* - u_h \rrbracket \cdot \{\{\mathbf{\Pi}_h\}\} ds + \int_{\Gamma_0} \{\{u_h^*\}\} \llbracket \mathbf{\Pi}_h \rrbracket ds. \quad (23)$$

We omit the superscript  $k\ell$  for brevity of notation in the expressions above. Now, using the identity (21) in (23), we get

$$(\mathbf{q}_h, \mathbf{\Pi}_h)_{L^2(\Omega_h)} = \int_{\Omega_h} \nabla_h u_h \cdot \mathbf{\Pi}_h \, d\mathbf{x} + \int_{\Gamma} \llbracket u_h^* - u_h \rrbracket \cdot \{\{\mathbf{\Pi}_h\}\} ds + \int_{\Gamma_0} \{\{u_h^* - u_h\}\} \llbracket \mathbf{\Pi}_h \rrbracket ds. \quad (24)$$

Taking  $\mathbf{\Pi}_h = \nabla_h \phi_h$  and combining (22) and (24), we obtain

$$\begin{aligned} &(\nabla_h u_h, \nabla_h \phi_h)_{L^2(\Omega_h)} + (cu_h, \phi_h)_{L^2(\Omega_h)} + \int_{\Gamma} (\llbracket u_h^* - u_h \rrbracket \cdot \{\{\nabla_h \phi_h\}\} - \llbracket \phi_h \rrbracket \cdot \{\{\mathbf{q}_h^*\}\}) ds \\ &+ \int_{\Gamma_0} (\{\{u_h^* - u_h\}\} \llbracket \nabla_h \phi_h \rrbracket - \{\{\phi_h\}\} \llbracket \mathbf{q}_h^* \rrbracket) ds = (f, \phi_h)_{L^2(\Omega_h)}. \end{aligned} \quad (25)$$

The numerical flux is defined by considering the internal penalty fluxes given by

$$\mathbf{q}_h^{*k\ell} = \{\{\nabla u_h\}\}^{k\ell} - \tau \llbracket u_h \rrbracket^{k\ell}, \quad u_h^{*k\ell} = \{\{u_h\}\}^{k\ell}, \quad \text{on } \Gamma_0 \quad (26)$$

$$\mathbf{q}_h^{*kB} = \nabla u_h^k - \tau \mathbf{n}^{kB} (u_h^k - g_D), \quad u_h^{*kB} = g_D, \quad \text{on } \partial\Omega_h \quad (27)$$

where  $\ell$  corresponds to the index of an adjacent element in the case of an inner interface and  $\ell = B$  for a boundary element. Parameter  $\tau = \eta/h_e$ , where  $\eta$  is some large positive constant. Moreover,  $g_D$  defines the boundary condition imposed in  $\partial\Omega_h$ . Then, replacing the numerical flux in Eq. (25) and attending that  $\llbracket \{\{\cdot\}\} \rrbracket = 0$ ,  $\llbracket \llbracket \cdot \rrbracket \rrbracket = 0$ ,  $\{\{\{\{\cdot\}\}\}\} = \{\{\cdot\}\}$  and  $\{\{\llbracket \cdot \rrbracket\}\} = \llbracket \cdot \rrbracket$ , we may write

$$\begin{aligned} &(\nabla_h u_h, \nabla_h \phi_h)_{L^2(\Omega_h)} + (cu_h, \phi_h)_{L^2(\Omega_h)} - \int_{\Gamma} \left( \llbracket u_h \rrbracket \cdot \{\{\nabla_h \phi_h\}\} + \llbracket \phi_h \rrbracket \cdot \{\{\nabla_h u_h\}\} \right) ds \\ &+ \int_{\Gamma} \llbracket \phi_h \rrbracket \cdot \frac{\eta}{h_e} \llbracket u_h \rrbracket ds = (f, \phi_h)_{L^2(\Omega_h)} - \sum_{k \in I^B} \int_{e^{kB}} \left( g_D^k \mathbf{n}^{kB} \cdot \nabla_h \phi_h^k - \frac{\eta}{h_e} g_D^k \phi_h^k \right) ds, \end{aligned} \quad (28)$$

where  $g_D^k = g_D|_{\tau^k}$ , for  $k \in I^B$ . Thus, scheme (28) corresponds to the interior penalty Galerkin method [21].

Instead of imposing homogeneous Dirichlet boundary conditions on  $\partial\Omega_h$ , i.e.  $g_D = 0$ , we consider a new boundary condition  $g_{ROD}$  determined by the ROD method. The polynomial  $g_{ROD}$  takes into account the geometrical mismatch between  $\partial\Omega$  and  $\partial\Omega_h$

and it has the following decomposition

$$g_{ROD}(\mathbf{x}; \mathbf{a}) = \bigoplus_{k \in I^B} g^k(\mathbf{x}; \mathbf{a}^k),$$

and in each element  $T^k$ , with  $k \in I^B$ , the local polynomial  $g^k$  has a polynomial decomposition with the two-dimensional Lagrange polynomials

$$\mathbf{x} \in T^k \in \mathcal{T}_h : \quad g^k(\mathbf{x}; \mathbf{a}^k) = \sum_{i=1}^{N_p} a_i^k \ell_i^k(\mathbf{x}), \quad (29)$$

where vector  $\mathbf{a}^k = (a_1^k, \dots, a_{N_p}^k)^\top$  gathers the  $N_p$  nodal values,  $a_i^k$ , and  $\mathbf{a}$  gathers all the vectors  $\mathbf{a}^k$ ,  $k \in I^B$ . From [30], recall that the coefficients of each polynomial  $g^k$ ,  $k \in I^B$ , are determined by solving the following system

$$\begin{bmatrix} \mathbf{I}_{N_p} & \mathbf{B}^k \\ (\mathbf{B}^k)^\top & \mathbf{0}_{N+1} \end{bmatrix} \begin{bmatrix} \mathbf{a}^k \\ \boldsymbol{\lambda}^k \end{bmatrix} = \begin{bmatrix} \mathbf{u}^k \\ \mathbf{0} \end{bmatrix}, \quad (30)$$

where  $\mathbf{I}_{N_p}$  is the identity matrix in  $\mathbb{R}^{N_p \times N_p}$ ,  $\mathbf{0}_{N+1}$  is the null matrix in  $\mathbb{R}^{(N+1) \times (N+1)}$ ,  $\mathbf{0}$  is the null vector in  $\mathbb{R}^{(N+1) \times 1}$ , and  $\mathbf{B}^k = [\mathbf{B}_1^k, \dots, \mathbf{B}_{N+1}^k]$  in  $\mathbb{R}^{N_p \times (N+1)}$ , with

$$\mathbf{B}_r^k = \left[ \ell_j^k \left( P_r^k \right) \right]_{j=1}^{N_p},$$

for  $r = 1, \dots, N+1$ , including the vertexes of  $T^k$  on  $\partial\Omega_h$ . Thus,

$$\begin{bmatrix} \mathbf{a}^k \\ \boldsymbol{\lambda}^k \end{bmatrix} = \begin{bmatrix} \mathbf{I}_{N_p} & \mathbf{B}^k \\ (\mathbf{B}^k)^\top & \mathbf{0}_{N+1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{u}^k \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} C_1 & C_2 \\ C_3 & C_4 \end{bmatrix} \begin{bmatrix} \mathbf{u}^k \\ \mathbf{0} \end{bmatrix},$$

where  $C_i$  denotes the  $i$ -th block of the inverse matrix. Noticing that  $\mathbf{I}_{N_p}$  is invertible, the inverse of the  $N_p \times N_p$  matrix  $C_1$  is given by [6]

$$\begin{aligned} C_1 &= \mathbf{I}_{N_p}^{-1} + \mathbf{I}_{N_p}^{-1} \mathbf{B}^k \left( \mathbf{0}_{N+1} - (\mathbf{B}^k)^\top \mathbf{I}_{N_p}^{-1} \mathbf{B}^k \right)^{-1} (\mathbf{B}^k)^\top \mathbf{I}_{N_p}^{-1} \\ &= \mathbf{I}_{N_p} - \mathbf{B}^k \left( (\mathbf{B}^k)^\top \mathbf{B}^k \right)^{-1} (\mathbf{B}^k)^\top. \end{aligned}$$

Thus, we get

$$\mathbf{a}^k = \mathbf{u}^k - \mathbf{B}^k \left( (\mathbf{B}^k)^\top \mathbf{B}^k \right)^{-1} (\mathbf{B}^k)^\top \mathbf{u}^k.$$

Note that

$$\left(\mathbf{B}^k\right)^T \mathbf{u}^k = \begin{bmatrix} \ell_1^k(P_1^k) & \cdots & \ell_{N_p}^k(P_1^k) \\ \ell_1^k(P_2^k) & \cdots & \ell_{N_p}^k(P_2^k) \\ \vdots & \ddots & \vdots \\ \ell_1^k(P_{N+1}^k) & \cdots & \ell_{N_p}^k(P_{N+1}^k) \end{bmatrix} \begin{bmatrix} u_1^k \\ u_2^k \\ \vdots \\ u_{N_p}^k \end{bmatrix} = \begin{bmatrix} u_h^k(P_1^k) \\ u_h^k(P_2^k) \\ \vdots \\ u_h^k(P_{N+1}^k) \end{bmatrix} = \mathbf{0},$$

since  $u_h^k \in \mathcal{W}_h$ . Then  $\mathbf{a}^k = \mathbf{u}^k$  and  $g^k = u_h^k$ , with  $k \in I^B$ . Replacing  $g_D^k$  by  $u_h^k$  in (28), we get for  $v \in \mathcal{V}_h$

$$\begin{aligned} & (\nabla_h u_h, \nabla_h v)_{L^2(\Omega_h)} + (cu_h, v)_{L^2(\Omega_h)} - \int_{\Gamma} \left( \llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\} + \llbracket v \rrbracket \cdot \{\{\nabla_h u_h\}\} - \frac{\eta}{h_e} \llbracket v \rrbracket \cdot \llbracket u_h \rrbracket \right) ds \\ &= (f, v)_{L^2(\Omega_h)} - \sum_{k \in I^B} \int_{e^{k_B}} \left( u_h^k \mathbf{n}^{k_B} \cdot \nabla_h v|_{\tau^k} - \frac{\eta}{h_e} u_h^k v|_{\tau^k} \right) ds \\ &= (f, v)_{L^2(\Omega_h)} - \int_{\partial\Omega_h} \llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\} ds + \int_{\partial\Omega_h} \frac{\eta}{h_e} \llbracket u_h \rrbracket \cdot \llbracket v \rrbracket ds. \end{aligned} \quad (31)$$

Attending that  $v = 0$  on  $\partial\Omega_h$  and considering (31), we obtain

$$\begin{aligned} & (\nabla_h u_h, \nabla_h v)_{L^2(\Omega_h)} + (cu_h, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\} ds - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{\{\nabla_h u_h\}\} ds \\ &+ \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket v \rrbracket \cdot \llbracket u_h \rrbracket ds = (f, v)_{L^2(\Omega_h)}. \end{aligned} \quad (32)$$

Then, the variational problem of (1)–(2) can be reformulated as follows: find  $u_h \in \mathcal{W}_h$  such that

$$a_h(u_h, v) = (f, v)_{L^2(\Omega_h)}, \quad \forall v \in \mathcal{V}_h, \quad (33)$$

where the bilinear form  $a_h(\cdot, \cdot)$  is defined as

$$\begin{aligned} a_h(u_h, v) &= (\nabla_h u_h, \nabla_h v)_{L^2(\Omega_h)} + (cu_h, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\} ds \\ &- \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{\{\nabla_h u_h\}\} ds + \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket v \rrbracket \cdot \llbracket u_h \rrbracket ds. \end{aligned} \quad (34)$$

We call (33)–(34) the primal formulation of the method and the bilinear form  $a_h(\cdot, \cdot)$  the primal form.

### 3 Existence and Uniqueness of the Solution

In this section, we prove the existence and uniqueness of the numerical solution. We start by analysing some basic properties of the method, namely, we show the boundedness of the bilinear form  $a_h(\cdot, \cdot)$  defined by (34). Then, we prove a weak coercivity in connection with finite-dimensional subspaces, with  $\dim(\mathcal{W}_h) = \dim(\mathcal{V}_h)$ .

Now, we present the Generalized Lax-Milgram Theorem as stated by Brezzi [8].

**Theorem 3.1** ([8]). *Let  $X$  and  $Y$  be two Hilbert spaces and consider a continuous real bilinear form defined on  $X \times Y$ . The following properties are equivalent:*

1.  *$a$  is weakly coercive, i.e., the following conditions are satisfied:*
  - (i)  $\exists \alpha > 0$  such that  $\inf_{u \in X \setminus \{0_X\}} \sup_{v \in Y \setminus \{0_Y\}} \frac{a(u,v)}{\|u\|_X \|v\|_Y} \geq \alpha$ ;
  - (ii)  $\forall v \in Y \setminus \{0_Y\}, \exists u \in X$ , such that  $a(u,v) \neq 0$ .
2.  $\forall L \in Y'$  the problem  $a(u,v) = L(v)$  has a unique solution  $u \in X$  which satisfies the stability condition

$$\|u\|_X \leq \frac{\|L\|_{Y'}}{\alpha},$$

where  $\alpha$  is the co-norm of  $a$ , i.e. the maximum of all  $\alpha$  satisfying (i).

In practice, the subspaces  $X$  and  $Y$  are often finite-dimensional. The following corollary establishes a result of the weak coercivity in connection with finite-dimensional subspaces. In particular, condition (ii) can be replaced with  $\dim(X) = \dim(Y)$  for bilinear forms associated with finite-dimensional spaces.

**Corollary 3.1** ([15]). *If  $X$  and  $Y$  are finite-dimensional spaces, the bilinear form  $a$  is weakly coercive over  $X \times Y$  if and only if either:*

- condition (i) holds and  $\dim(X) = \dim(Y)$ ;
- matrix  $A$  associated with the bilinear form  $a$  is a square invertible matrix;

both conditions being equivalent.

Thus, to prove the existence and uniqueness of the solution, we may prove that the bilinear form (34) is bounded, the inf-sup condition (i) holds, and  $\dim(\mathcal{W}_h) = \dim(\mathcal{V}_h)$ . We start by discussing the boundedness of the bilinear form  $a_h(\cdot, \cdot)$ .

### 3.1 Boundedness

We show that the bilinear form  $a_h(\cdot, \cdot)$  is continuous on  $H^2(\mathcal{T}_h) \times H^2(\mathcal{T}_h)$  equipped with the norm  $\|\cdot\|$ , i.e., there exists a positive real number  $C_b$  such that

$$|a_h(u_h, v)| \leq C_b \|u_h\| \|v\|, \quad \forall u_h \in H^2(\mathcal{T}_h), \forall v \in H^2(\mathcal{T}_h). \quad (35)$$

In particular, note that  $\mathcal{W}_h \subset H^2(\mathcal{T}_h)$  and  $\mathcal{V}_h \subset H^2(\mathcal{T}_h)$ . Recall that

$$\begin{aligned} |a_h(u_h, v)| &\leq \left| (\nabla_h u_h, \nabla_h v)_{L^2(\Omega_h)} \right| + \left| (c u_h, v)_{L^2(\Omega_h)} \right| + \int_{\Gamma_0} |[[u_h]] \cdot \{\{\nabla_h v\}\}| \, ds \\ &\quad + \int_{\Gamma_0} |[[v]] \cdot \{\{\nabla_h u_h\}\}| \, ds + \int_{\Gamma_0} \left| \frac{\eta}{h_e} [[v]] \cdot [[u_h]] \right| \, ds. \end{aligned} \quad (36)$$

We show that each term in (36) can be bounded by the  $\|\cdot\|$ -norm. For the first and second terms, we use the Cauchy-Schwarz inequality and obtain

$$\left| (\nabla_h u_h, \nabla_h v)_{L^2(\Omega_h)} \right| + \left| (c u_h, v)_{L^2(\Omega_h)} \right| \leq \sum_{k=1}^K \left| (\nabla_h u_h, \nabla_h v)_{L^2(T^k)} \right| + \sum_{k=1}^K \left| (c u_h, v)_{L^2(T^k)} \right|$$

$$\begin{aligned}
&\leq \sum_{k=1}^K \|\nabla u_h\|_{L^2(T^k)} \|\nabla v\|_{L^2(T^k)} + \|c\|_{L^\infty(\Omega_h)} \sum_{k=1}^K \|u_h\|_{L^2(T^k)} \|v\|_{L^2(T^k)} \\
&\leq 2 \left(1 + \|c\|_{L^\infty(\Omega_h)}\right) \|u_h\| \|v\|. \tag{37}
\end{aligned}$$

Now, we bound the third term in (36). Using the definition of jump and average, a straightforward computation shows that

$$\int_e \llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\} ds = \int_e \llbracket \mathbf{n} u_h \rrbracket \{\{\nabla_h v \cdot \mathbf{n}\}\} ds.$$

Thus,

$$\int_e |\llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\}| ds = \int_e |\llbracket \mathbf{n} u_h \rrbracket \{\{\nabla_h v \cdot \mathbf{n}\}\}| ds \leq \|\llbracket \mathbf{n} u_h \rrbracket\|_{L^2(e)} \left\| \left\{ \left\{ \frac{\partial v}{\partial n} \right\} \right\} \right\|_{L^2(e)}. \tag{38}$$

Employing the trace inequality ([2], equation (2.5)) for  $e \in \partial T^k$ ,  $T^k \in \mathcal{T}_h$ , we get

$$\left\| \frac{\partial v}{\partial n} \right\|_{L^2(e)}^2 \leq C_T^2 \left( h_e^{-1} |v|_{H^1(T^k)}^2 + h_e |v|_{H^2(T^k)}^2 \right), \quad v \in H^2(T^k). \tag{39}$$

Using the inequality above and the inequality (9), we conclude

$$\begin{aligned}
\sum_{e \in \mathcal{E}_0} \int_e |\llbracket u_h \rrbracket \cdot \{\{\nabla_h v\}\}| ds &\leq \sum_{e \in \mathcal{E}_h} \|\llbracket u_h \rrbracket\|_{L^2(e)} \left\| \left\{ \left\{ \frac{\partial v}{\partial n} \right\} \right\} \right\|_{L^2(e)} \\
&\leq \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \|\llbracket u_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2} \left( \sum_{e \in \mathcal{E}_h} h_e \left\| \left\{ \left\{ \frac{\partial v}{\partial n} \right\} \right\} \right\|_{L^2(e)}^2 \right)^{1/2} \\
&\leq |u_h|_* \left( \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e \left\| \frac{\partial v}{\partial n} \right\|_{L^2(e)}^2 \right) \right)^{1/2} \\
&\leq |u_h|_* \left( \sum_{k=1}^K C_T^2 \left( |v|_{H^1(T^k)}^2 + h_k^2 |v|_{H^2(T^k)}^2 \right) \right)^{1/2} \\
&\leq C_T \sqrt{1 + C^2} |u_h|_* |v|_{H^1(\Omega_h)} \leq C_T \sqrt{1 + C^2} |u_h|_* \|v\|. \tag{40}
\end{aligned}$$

Similarly, we can bound the fourth term in (36) as

$$\int_{\Gamma_0} |\llbracket v \rrbracket \cdot \{\{\nabla_h u_h\}\}| ds \leq C_T \sqrt{1 + C^2} |v|_* |u_h|_{H^1(\Omega_h)} \leq C_T \sqrt{1 + C^2} |v|_* \|u_h\|. \tag{41}$$

Finally, for the last term, we have

$$\int_{\Gamma_0} \left| \llbracket u_h \rrbracket \cdot \frac{\eta}{h_e} \llbracket v \rrbracket \right| ds \leq \eta \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \|\llbracket u_h \rrbracket\|_{L^2(e)}^2 \right)^{1/2} \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \|\llbracket v \rrbracket\|_{L^2(e)}^2 \right)^{1/2} = \eta |u_h|_* |v|_*. \quad (42)$$

Combining (37), (40)–(42), for  $u_h \in H^2(\mathcal{T}_h)$  and  $v \in H^2(\mathcal{T}_h)$ , we get (35) with  $C_b = \max \left\{ 2 \left( 1 + \|c\|_{L^\infty(\Omega_h)} \right), C_T \sqrt{1 + C^2}, \eta \right\}$ .

### 3.2 Weak coercivity

Now, we address the weak coercivity of the bilinear form. The following theorem establishes an inf-sup condition for the DG-ROD formulation (33)–(34).

**Theorem 3.2.** *Consider the bilinear form  $a_h(\cdot, \cdot)$  defined in (34). Given  $h$  sufficiently small and  $\eta$  sufficiently large, there exists a constant  $\alpha > 0$  independent of  $h$  such that*

$$\forall w \in \mathcal{W}_h \setminus \{0\}, \quad \sup_{v \in \mathcal{V}_h \setminus \{0\}} \frac{a_h(w, v)}{\|w\| \|v\|} \geq \alpha. \quad (43)$$

*Proof.* Let  $w \in \mathcal{W}_h$ . Let  $v \in \mathcal{V}_h$  such that  $v$  coincide with  $w$  at all mesh nodes, except those located on  $\partial\Omega_h$  that are not mesh vertexes. Thus,  $(w - v)|_{T^k} = 0$ , for every element  $T^k$  that does not have an edge on  $\partial\Omega_h$ . Taking advantage of the relationship between  $w$  and  $v$ , we may write

$$\begin{aligned} a_h(w, v) &= \sum_{k=1}^K \left( (\nabla w, \nabla w)_{L^2(T^k)} + (cw, w)_{L^2(T^k)} \right) + \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket w \rrbracket \cdot \llbracket w \rrbracket ds - 2 \int_{\Gamma_0} \llbracket w \rrbracket \cdot \{\{\nabla_h w\}\} ds \\ &\quad - \sum_{k \in I^B} \left( (\nabla w, \nabla r^k(w))_{L^2(T^k)} + (cw, r^k(w))_{L^2(T^k)} \right) - \sum_{e \in \mathcal{E}_0^B} \int_e \frac{\eta}{h_e} \llbracket r(w) \rrbracket \cdot \llbracket w \rrbracket ds \\ &\quad + \sum_{e \in \mathcal{E}_0^B} \int_e \llbracket w \rrbracket \cdot \{\{\nabla_h r(w)\}\} ds + \sum_{e \in \mathcal{E}_0^B} \int_e \llbracket r(w) \rrbracket \cdot \{\{\nabla_h w\}\} ds, \end{aligned} \quad (44)$$

where  $r(w) = \bigoplus_{k \in I^B} r^k(w)$ ,  $r^k(w) = (w - v)|_{T^k} = \sum_{i \in I^{kB}} w(\mathbf{x}_i^k) \ell_i^k(x)$ ,  $\mathcal{E}_h^B = \bigcup_{k \in I^B} \partial T^k$  and  $\mathcal{E}_0^B = \mathcal{E}_h^B \cap \Gamma_0$ .

We aim to estimate bounds for each of the nine terms in (44). Using the definition of norm, we may write

$$\sum_{k=1}^K (\nabla w, \nabla w)_{L^2(T^k)} = \|\nabla_h w\|_{L^2(\Omega_h)}^2, \quad (45)$$

$$\sum_{e \in \mathcal{E}_0} \int_e \frac{\eta}{h_e} \llbracket w \rrbracket \cdot \llbracket w \rrbracket ds = \eta \left( \sum_{e \in \mathcal{E}_0} h_e^{-1} \|\llbracket w \rrbracket\|_{L^2(e)}^2 \right) = \eta C_w |w|_*^2, \quad (46)$$

with  $0 < C_w \leq 1$ . Now, using the inequalities of the boundedness of subsection 3.1 (see (40)) and for  $\epsilon_1 > 0$  we can estimate the following upper bound for the fourth term in (44)

$$\begin{aligned} 2 \sum_{e \in \mathcal{E}_0} \int_e \llbracket w \rrbracket \cdot \{\{\nabla_h w\}\} ds &\leq 2 \sum_{e \in \mathcal{E}_0} \int_e \left| \llbracket w \rrbracket \cdot \{\{\nabla_h w\}\} \right| ds \\ &\leq 2C_T \sqrt{1 + C^2} |w|_{H^1(\Omega_h)} |w|_* \leq C_T \sqrt{1 + C^2} \left( \epsilon_1 |w|_{H^1(\Omega_h)}^2 + \frac{|w|_*^2}{\epsilon_1} \right). \end{aligned} \quad (47)$$

Now, we bound the terms with  $r(w)$ . In order to achieve that, we start by noticing that

$$\left\| r^k(w) \right\|_{L^2(T^k)} \leq \sum_{i \in I^{kB}} \left| w(\mathbf{x}_i^k) \right| \left\| \ell_i^k \right\|_{L^2(T^k)}, \quad (48)$$

$$\left\| \nabla r^k(w) \right\|_{L^2(T^k)} \leq \sum_{i \in I^{kB}} \left| w(\mathbf{x}_i^k) \right| \left\| \nabla \ell_i^k \right\|_{L^2(T^k)}. \quad (49)$$

From standard results, it holds for mesh independent constants  $C_1$  and  $C_2$

$$\left\| \ell_i^k \right\|_{L^2(T^k)} \leq C_1 h_k \quad \text{and} \quad \left\| \nabla \ell_i^k \right\|_{L^2(T^k)} \leq C_2.$$

Using Proposition (A.1), Lemma (A.1) and following similar arguments as in [28], we may prove that

$$|w(\mathbf{x}_j^k)| \leq C_{\partial\Omega} C_\infty C_J h_k \|\nabla w\|_{L^2(T^k)}.$$

Then, we get

$$\left\| r^k(w) \right\|_{L^2(T^k)} \leq \tilde{C}_1 h^2 \|\nabla w\|_{L^2(T^k)}, \quad (50)$$

where  $\tilde{C}_1 = (N - 1)C_1 C_{\partial\Omega_h} C_\infty C_J$ , and

$$\left\| \nabla r^k(w) \right\|_{L^2(T^k)} \leq \tilde{C}_2 h \|\nabla w\|_{L^2(T^k)}, \quad (51)$$

where  $\tilde{C}_2 = (N - 1)C_2 C_{\partial\Omega_h} C_\infty C_J$ .

Thus, using the inequality (51), we get for the fifth term in (44)

$$\sum_{k \in I^B} (\nabla w, \nabla r^k(w))_{L^2(T^k)} \leq \sum_{k \in I^B} \|\nabla w\|_{L^2(T^k)} \left\| \nabla r^k(w) \right\|_{L^2(T^k)} \leq \tilde{C}_2 h \|\nabla_h w\|_{L^2(\Omega_h)}^2. \quad (52)$$



Using the inequality (50) and considering  $\epsilon_2 > 0$ , we obtain for the sixth term in (44)

$$\begin{aligned} \sum_{k \in I^B} (cw, r^k(w))_{L^2(T^k)} &\leq \sum_{k \in I^B} \|cw\|_{L^2(T^k)} \|r^k(w)\|_{L^2(T^k)} \leq \sum_{k \in I^B} \|cw\|_{L^2(T^k)} \tilde{C}_1 h^2 \|\nabla w\|_{L^2(T^k)} \\ &\leq \tilde{C}_1 h^2 \|cw\|_{L^2(\Omega_h)} \|\nabla_h w\|_{L^2(\Omega_h)} \\ &\leq \frac{\tilde{C}_1}{2} h^2 \left( \epsilon_2 \|cw\|_{L^2(\Omega_h)}^2 + \frac{\|\nabla_h w\|_{L^2(\Omega_h)}^2}{\epsilon_2} \right). \end{aligned} \quad (53)$$

Let us now estimate the last three terms in (44). Note that, for  $\epsilon_3 > 0$ ,

$$\begin{aligned} \sum_{e \in \mathcal{E}_0^B} \int_e \frac{\eta}{h_e} \llbracket r(w) \rrbracket \cdot \llbracket w \rrbracket ds &\leq \sum_{e \in \mathcal{E}_h^B} \int_e \frac{\eta}{h_e} |\llbracket r(w) \rrbracket \cdot \llbracket w \rrbracket| ds \leq \eta \sum_{e \in \mathcal{E}_h^B} \left\| \frac{\llbracket r(w) \rrbracket}{h_e^{1/2}} \right\|_{L^2(e)} \left\| \frac{\llbracket w \rrbracket}{h_e^{1/2}} \right\|_{L^2(e)} \\ &\leq \eta |r(w)|_* |w|_* \leq \frac{\eta}{2} \left( \epsilon_3 |r(w)|_*^2 + \frac{|w|_*^2}{\epsilon_3} \right). \end{aligned} \quad (54)$$

Using a trace inequality ([2], equation (2.4)) for  $e \in \partial T^k$ ,  $T^k \in \mathcal{T}_h$

$$\|v\|_{L^2(e)}^2 \leq \tilde{C}_T^2 \left( h_e^{-1} \|v\|_{L^2(T^k)}^2 + h_e |v|_{H^1(T^k)}^2 \right), \quad v \in H^1(T^k), \quad (55)$$

for  $v = r(w)$ , and employing the inequality (5) we obtain

$$\begin{aligned} |r(w)|_*^2 &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e^{-1} \|r^k(w)\|_{L^2(e)}^2 \right) \\ &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e^{-1} \tilde{C}_T^2 h_e^{-1} \left( \|r^k(w)\|_{L^2(T^k)}^2 + h_e^2 \|\nabla r^k(w)\|_{L^2(T^k)}^2 \right) \right) \\ &\leq \sum_{k=1}^K \tilde{C}_T^2 \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} \frac{1}{\mu^2 h_k^2} \left( \tilde{C}_1^2 h_k^4 \|\nabla w\|_{L^2(T^k)}^2 + \tilde{C}_2^2 h_k^4 \|\nabla w\|_{L^2(T^k)}^2 \right) \right) \\ &\leq C_*^2 h^2 \|\nabla_h w\|_{L^2(\Omega_h)}^2, \end{aligned}$$

where  $C_*^2 = \tilde{C}_T^2 (\tilde{C}_1^2 + \tilde{C}_2^2) / \mu^2$ . Then, we can write (54) as

$$\sum_{e \in \mathcal{E}_0^B} \int_e \frac{\eta}{h_e} \llbracket r(w) \rrbracket \cdot \llbracket w \rrbracket ds \leq \frac{\eta}{2} \left( \epsilon_3 C_*^2 h^2 \|\nabla_h w\|_{L^2(\Omega_h)}^2 + \frac{|w|_*^2}{\epsilon_3} \right). \quad (56)$$

For the eighth term in (44) we have

$$\sum_{e \in \mathcal{E}_0^B} \int_e \llbracket w \rrbracket \cdot \{ \{ \nabla_h r(w) \} \} ds \geq \frac{-1}{2} \sum_{e \in \mathcal{E}_0^B} \int_e \left( \left| h_e^{-1/2} \llbracket w \mathbf{n} \rrbracket \right|^2 + \left| h_e^{1/2} \left\{ \left\{ \frac{\partial r(w)}{\partial n} \right\} \right\} \right|^2 \right) ds, \quad (57)$$

and then we need to bound each of the terms on the right-hand side of the inequality above. Note that

$$\sum_{e \in \mathcal{E}_0^B} \int_e \left| h_e^{-1/2} \llbracket w \mathbf{n} \rrbracket \right|^2 ds \leq \sum_{e \in \mathcal{E}_h} \int_e h_e^{-1} \llbracket w \rrbracket^2 ds = |w|_*^2. \quad (58)$$

On the other hand, using the inequality (39), we get

$$\begin{aligned} \sum_{e \in \mathcal{E}_0^B} \int_e \left| h_e^{1/2} \left\{ \left\{ \frac{\partial r(w)}{\partial n} \right\} \right\} \right|^2 ds &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e \left\| \frac{\partial r(w)}{\partial n} \right\|_{L^2(e)}^2 \right) \\ &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e C_T^2 h_e^{-1} \left( \left\| \nabla r^k(w) \right\|_{L^2(T^k)}^2 + h_e^2 \left| r^k(w) \right|_{H^2(T^k)}^2 \right) \right). \end{aligned}$$

Using (51) and recalling that  $h_e^2 |r^k(w)|_{H^2(T^k)}^2 \leq C^2 |r^k(w)|_{H^1(T^k)}^2$ , (see [7]), we obtain

$$\begin{aligned} \sum_{e \in \mathcal{E}_0^B} \int_e \left| h_e^{1/2} \left\{ \left\{ \frac{\partial r(w)}{\partial n} \right\} \right\} \right|^2 ds &\leq \sum_{k=1}^K C_T^2 (1 + C^2) \left\| \nabla r^k(w) \right\|_{L^2(T^k)}^2 \\ &\leq C_T^2 (1 + C^2) \tilde{C}_2^2 h^2 \left\| \nabla_h w \right\|_{L^2(\Omega_h)}^2. \end{aligned} \quad (59)$$

Finally, replacing (58) and (59) in (57), we get

$$\sum_{e \in \mathcal{E}_0^B} \int_e \llbracket w \rrbracket \cdot \{ \{ \nabla_h r(w) \} \} ds \geq \frac{-1}{2} \left( |w|_*^2 + C_T^2 (1 + C^2) \tilde{C}_2^2 h^2 \left\| \nabla_h w \right\|_{L^2(\Omega_h)}^2 \right). \quad (60)$$

Applying a similar argument, we can estimate a lower bound for last term in (44). Note that

$$\sum_{e \in \mathcal{E}_0^B} \int_e \llbracket r(w) \rrbracket \cdot \{ \{ \nabla_h w \} \} ds \geq \frac{-1}{2} \sum_{e \in \mathcal{E}_0^B} \int_e \left( \left| h_e^{-1} \llbracket r(w) \mathbf{n} \rrbracket \right|^2 + \left| h_e \left\{ \left\{ \frac{\partial w}{\partial n} \right\} \right\} \right|^2 \right) ds, \quad (61)$$

and so we need to bound each of the terms on the right-hand side of the inequality above. Recalling assumption (6) we may write

$$\sum_{e \in \mathcal{E}_0^B} \int_e \left| h_e^{-1} \llbracket r(w) \mathbf{n} \rrbracket \right|^2 ds \leq \frac{1}{h_{\min}} |r(w)|_*^2 \leq \frac{1}{h_{\min}} C_*^2 h^2 |w|_{H^1(\mathcal{T}_h)}^2 \leq \tilde{\rho} C_*^2 h |w|_{H^1(\mathcal{T}_h)}^2. \quad (62)$$

On the other hand,

$$\begin{aligned} \sum_{e \in \mathcal{E}_0^B} \int_e \left| h_e \left\{ \left\{ \frac{\partial w}{\partial \mathbf{n}} \right\} \right\} \right|^2 ds &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e^2 \left\| \frac{\partial w}{\partial \mathbf{n}} \right\|_{L^2(e)}^2 \right) \\ &\leq \sum_{k=1}^K \left( \sum_{e \in \partial T^k \cap \mathcal{E}_h} h_e^2 C_T^2 h_e^{-1} \left( \|\nabla w\|_{L^2(T^k)}^2 + h_e^2 |w|_{H^2(T^k)}^2 \right) \right) \\ &\leq C_T^2 (1 + C^2) \sum_{k=1}^K h_k \|\nabla w\|_{L^2(T^k)}^2 \leq C_T^2 (1 + C^2) h |w|_{H^1(\mathcal{T}_h)}^2. \end{aligned} \quad (63)$$

Thus, replacing (62) and (63) in (61), we obtain

$$\sum_{e \in \mathcal{E}_0^B} \int_e \llbracket r(w) \rrbracket \cdot \{ \{ \nabla_h w \} \} ds \geq \frac{-1}{2} \left( \tilde{\rho} C_*^2 + C_T^2 (1 + C^2) \right) h |w|_{H^1(\mathcal{T}_h)}^2. \quad (64)$$

Let us finally consider the terms containing the function  $c$  in (44). Taking (53) into account and considering

$$h^2 < h_1^2 = 2 / \left( \tilde{C}_1 \epsilon_2 \|c\|_{L^\infty(\Omega_h)} \right), \quad (65)$$

we get

$$\begin{aligned} \sum_{k=1}^K (cw, w)_{L^2(T^k)} - \sum_{k \in I^B} (cw, r^k(w))_{L^2(T^k)} &\geq \sum_{k=1}^K \left( (cw, w)_{L^2(T^k)} - \frac{\tilde{C}_1}{2} h^2 \epsilon_2 \|cw\|_{L^2(T^k)}^2 \right) \\ &= \sum_{k=1}^K \left( (cw, w)_{L^2(T^k)} + (cw, -\frac{\tilde{C}_1}{2} h^2 \epsilon_2 cw)_{L^2(T^k)} \right) \\ &= \sum_{k=1}^K \left( cw, \left( 1 - \frac{\tilde{C}_1}{2} h^2 \epsilon_2 c \right) w \right)_{L^2(T^k)} \geq c_{\min} c' \|w\|_{L^2(\Omega_h)}^2, \end{aligned}$$

with  $c_{\min} = \min_{\mathbf{x} \in \Omega_h} c(\mathbf{x})$  and

$$c' = 1 - \frac{\tilde{C}_1}{2} h^2 \epsilon_2 \|c\|_{L^\infty(\Omega_h)}. \quad (66)$$

Using (10), we may write

$$\|w\|^2 \leq C_{aux}^2 \left( \|w\|_{L^2(\Omega_h)}^2 + \|\nabla_h w\|_{L^2(\Omega_h)}^2 + |w|_*^2 \right). \quad (67)$$

Note that using an inequality of Poincaré-Friedrichs type valid for  $w \in H^1(\mathcal{T}_h)$  (see [2], Lemma 2.1), we get

$$\begin{aligned} \|w\|^2 &\leq C_{aux}^2 \left( C_P \left( \|\nabla_h w\|_{L^2(\Omega_h)}^2 + |w|_*^2 \right) + \|\nabla_h w\|_{L^2(\Omega_h)}^2 + |w|_*^2 \right) \\ &\leq C_{aux}^2 (1 + C_P) \left( \|\nabla_h w\|_{L^2(\Omega_h)}^2 + |w|_*^2 \right). \end{aligned} \quad (68)$$

Then, combining the bounds for each term of (44), namely (45)–(47), (52), (53), (56), (60), (64), we may write (44) as

$$a_h(w, v) \geq C_\alpha \|w\|^2, \quad (69)$$

where  $C_\alpha = \min\{\widehat{C}_1, \widehat{C}_2, \widehat{C}_3\}/C_{aux}^2$ , if  $c_{min} > 0$ , and  $C_\alpha = \min\{\widehat{C}_2, \widehat{C}_3\}/(C_{aux}^2(1 + C_P))$ , if  $c_{min} = 0$ , and  $\widehat{C}_1 = c_{min}c'$ , with  $c'$  is given by (66),

$$\begin{aligned} \widehat{C}_2 &= \left(1 - C_T \sqrt{1 + C^2} \epsilon_1\right) - h \left( \tilde{C}_2 + \frac{1}{2} \left( C_*^2 \tilde{\rho} + C_T^2 (1 + C^2) \right) \right) \\ &\quad - h^2 \left( \frac{\tilde{C}_1}{2\epsilon_2} + \frac{1}{2} C_T^2 (1 + C^2) \tilde{C}_2^2 + \eta \epsilon_3 \frac{C_*^2}{2} \right) \end{aligned}$$

and  $\widehat{C}_3 = -C_T \sqrt{1 + C^2} / \epsilon_1 - 1/2 + \eta C_w - \eta / (2\epsilon_3)$ .

Note that, considering (65) and  $c_{min} > 0$ ,  $\widehat{C}_1 > 0$  (see (66)). For  $\widehat{C}_2$ , if we take  $\epsilon_1 < 1/(C_T \sqrt{1 + C^2})$  and  $h < h_2$ , where  $h_2$  is the positive root of the equation (in  $h$ )  $\widehat{C}_2 = 0$ , we have  $\widehat{C}_2 > 0$ . For  $\widehat{C}_3$ , considering  $\epsilon_3 > 1/(2C_w)$  and taking

$$\eta = 2 \left( \frac{1}{C_w - \frac{1}{2\epsilon_3}} \right) \left( \frac{1}{2} + \frac{C_T \sqrt{1 + C^2}}{\epsilon_1} \right),$$

we get

$$\widehat{C}_3 = \frac{1}{2} + \frac{C_T \sqrt{1 + C^2}}{\epsilon_1} > 0.$$

Thus, considering  $h$  sufficiently small in the sense that  $h < h_0 = \min\{h_1, h_2\}$ , we have (69) with  $C_\alpha > 0$ .

Note that using (50)

$$\|v\|_{L^2(T^k)} \leq \|w\|_{L^2(T^k)} + \tilde{C}_1 h^2 \|\nabla w\|_{L^2(T^k)} \leq \sqrt{2}(1 + \tilde{C}_1 h^2) \|w\|_{H^1(T^k)}. \quad (70)$$

Similarly, considering (51), we obtain

$$\|\nabla v\|_{L^2(T^k)} \leq (1 + \tilde{C}_2 h) \|\nabla w\|_{L^2(T^k)}. \quad (71)$$

Note that  $|v|_*^2 \leq 2C_*^2 h^2 \|\nabla_h w\|_{L^2(\Omega_h)}^2 + 2|w|_*^2$  and then

$$\|v\|^2 \leq \widehat{C}_{aux}^2 \left( \|v\|_{L^2(\Omega_h)}^2 + \|\nabla_h v\|_{L^2(\Omega_h)}^2 + |v|_*^2 \right) \leq C_v^2 \|w\|^2,$$

where  $C_v^2 = \widehat{C}_{aux}^2 \max\{2, 2(1 + \tilde{C}_1 h_0^2)^2 + (1 + \tilde{C}_2 h_0)^2 + 2C_*^2 h_0^2\}$ . We get

$$a_h(w, v) \geq \frac{C_\alpha}{C_v} \|w\| \|v\| \Rightarrow \frac{a_h(w, v)}{\|w\| \|v\|} \geq \frac{C_\alpha}{C_v} = \alpha.$$

□

Using condition (i) of Corollary 3.1, since  $\dim(\mathcal{W}_h) = \dim(\mathcal{V}_h)$ , the fact that the inf-sup condition (43) holds, implies that (33)–(34) is uniquely solvable.

## 4 Error Estimates

In this section, we derive error estimates for the DG-ROD method for convex and non-convex domains. We first analyse whether the Galerking orthogonality holds in each case, and if it does not hold, we derive estimates for the resulting residual. To estimate the error, we use the inf-sup condition (43) and classical interpolation results. Let  $I_h(w) \in \mathcal{W}_h$  be the  $\mathcal{P}_N$ -interpolate of  $w$  at the nodes associated with  $\mathcal{W}_h$ . First, we note that if  $k \notin I^B$ , then  $I_h(w)$  is the standard interpolate of  $w$  at the mesh nodes  $\mathbf{x}_i^k$  (see left panel of Figure 2). If  $k \in I^B$ , then  $I_h(w)$  is the interpolate of  $w$  at the set of  $m_N + 2$  mesh nodes  $\mathbf{x}_i^k$  that do not lie in the interior of  $e^{kB}$ , together with the  $N - 1$  points lying on  $\partial\Omega$  associated with the mesh nodes of  $T^k$  lying in the interior of  $e^{kB}$  (see right panel of Figure 2).

### 4.1 Convex case

We start by discussing the consistency of the method. Let  $\Omega$  be a convex domain,  $u \in H^2(\Omega)$  the exact solution of the boundary value problem (1)–(2) and  $v \in \mathcal{V}_h$ . Attending that  $\Omega_h \subset \Omega$ ,  $\llbracket u \rrbracket = 0$ ,  $\llbracket \nabla u \rrbracket = 0$  and using the estimate (21), we get

$$\begin{aligned} a_h(u, v) &= (\nabla_h u, \nabla_h v)_{L^2(\Omega_h)} + (cu, v)_{L^2(\Omega_h)} \\ &\quad - \int_{\Gamma_0} \llbracket u \rrbracket \cdot \{\{\nabla_h v\}\} ds - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{\{\nabla_h u\}\} ds + \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket v \rrbracket \cdot \llbracket u \rrbracket ds \\ &= (\nabla_h u, \nabla_h v)_{L^2(\Omega_h)} + (cu, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{\{\nabla_h u\}\} ds \\ &= -(\nabla_h \cdot \nabla_h u, v)_{L^2(\Omega_h)} + \int_{\Gamma} \llbracket v \rrbracket \cdot \{\{\nabla_h u\}\} ds + \int_{\Gamma_0} \{\{v\}\} \llbracket \nabla_h u \rrbracket ds \end{aligned}$$

$$+(cu, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{\{\nabla_h u\}\} ds = (f, v)_{L^2(\Omega_h)}.$$

Thus, the method is consistent and the Galerkin orthogonality holds

$$a_h(u - u_h, v) = 0, \quad \forall v \in \mathcal{V}_h. \quad (72)$$

**Theorem 4.1.** *Let  $\Omega$  be convex and  $u \in H^{N+1}(\Omega)$  be the solution of (1)–(2). Then, for  $h$  sufficiently small and for a suitable constant  $\mathcal{C}$  independent of  $h$  and  $u$ , the solution  $u_h$  of (33)–(34) satisfies*

$$\| \|u - u_h\| \| \leq \mathcal{C} h^N |u|_{H^{N+1}(\Omega)}. \quad (73)$$

*Proof.* First, note that  $\| \|u - u_h\| \| \leq \| \|u - I_h(u)\| \| + \| \|u_h - I_h(u)\| \|$ , with  $I_h(u)$  the  $\mathcal{P}_N$ -interpolate of  $u$  at the nodes associated with  $\mathcal{W}_h$ . Using the inf-sup inequality (Theorem 3.2), we get

$$\| \|u_h - I_h(u)\| \| \leq \frac{1}{\alpha} \sup_{v \in \mathcal{V}_h \setminus \{0\}} \frac{a_h(u_h - I_h(u), v)}{\| \|v\| \|}. \quad (74)$$

Adding and subtracting  $u$  in the first argument of  $a_h$ , using the Galerkin orthogonality and the boundedness inequality (35) yields

$$\| \|u - u_h\| \| \leq \left(1 + \frac{C_b}{\alpha}\right) \| \|u - I_h(u)\| \|.$$

Recall that, as the interpolant  $I_h(u)$  is discontinuous across the inter-elements boundaries, the jumps  $u - I_h(u)$  will not be zero. Therefore

$$\begin{aligned} \| \|u - I_h(u)\| \|^2 &= \sum_{k=1}^K \left( \| \|u - I_h(u)\| \|_{L^2(T^k)}^2 + |u - I_h(u)|_{H^1(T^k)}^2 + h_k^2 |u - I_h(u)|_{H^2(T^k)}^2 \right) \\ &\quad + \sum_{e \in \mathcal{E}_h} h_e^{-1} \| \|u - I_h(u)\| \|_{L^2(e)}^2 \end{aligned} \quad (75)$$

and, using (75) and (55), we obtain

$$\begin{aligned} \| \|u - I_h(u)\| \|^2 &\leq C \sum_{k=1}^K \left( \| \|u - I_h(u)\| \|_{L^2(T^k)}^2 + |u - I_h(u)|_{H^1(T^k)}^2 + h_k^2 |u - I_h(u)|_{H^2(T^k)}^2 \right. \\ &\quad \left. + h_k^{-2} \| \|u - I_h(u)\| \|_{L^2(T^k)}^2 \right). \end{aligned} \quad (76)$$

From Lemma A.2, considering  $p = 2$ ,  $m = N + 1$ ,  $j = 0, 1, 2$  and  $h < 1$ , we establish

$$\| \|u - I_h(u)\| \| \leq \mathcal{C}_a h^N |u|_{H^{N+1}(\Omega)}.$$

Thus, (73) holds with  $\mathcal{C} = C_a (1 + C_b/\alpha)$ .  $\square$

The next theorem establishes that the DG-ROD solution exhibits an optimal  $\mathcal{O}(h^{N+1})$  convergence rate in the  $L^2$ -norm when  $N$ -degree piecewise polynomials are used, under certain regularity conditions on the solution.

**Theorem 4.2.** *Let  $\Omega$  be convex and  $u$  be the solution of (1)–(2) belonging to  $H^{N+1+r}(\Omega)$ , with  $r = 1/2 + \epsilon$ , for  $\epsilon > 0$  arbitrary small. Then, given  $h$  sufficiently small, the solution  $u_h$  of (33)–(34) satisfies for  $N > 1$  and a suitable constant  $\mathcal{C}_0$  independent of  $h$  and  $u$*

$$\|u - u_h\|_{L^2(\Omega_h)} \leq \mathcal{C}_0 h^{N+1} \|u\|_{H^{N+1+r}(\Omega)}. \quad (77)$$

*Proof.* Recall that every function in  $\mathcal{W}_h$  is defined in  $\bar{\Omega} \setminus \Omega_h$  and let  $z \in H_0^1(\Omega)$  be the solution of

$$\begin{aligned} -\Delta z(\mathbf{x}) + c(\mathbf{x})z(\mathbf{x}) &= u(\mathbf{x}) - u_h(\mathbf{x}), & \mathbf{x} \in \Omega, \\ z(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\Omega. \end{aligned}$$

We know that  $u - u_h \in L^2(\Omega)$ ,  $z \in H^2(\Omega)$ . It is well-known that if  $\partial\Omega$  is  $C^2$  (see, for example, [17], Theorem 6.3.4) or  $\Omega$  is a convex polygonal bounded domain with a Lipschitz boundary (see [20], Theorem 4.3.1.4) there exists a constant  $C(\Omega)$  such that

$$\|z\|_{H^2(\Omega)} \leq C(\Omega) \|u - u_h\|_{L^2(\Omega)}. \quad (78)$$

Then,

$$\|u - u_h\|_{L^2(\Omega)} \leq C(\Omega) \frac{\|u - u_h\|_{L^2(\Omega)}^2}{\|z\|_{H^2(\Omega)}} = C(\Omega) \frac{(u - u_h, -\Delta z + cz)_{L^2(\Omega)}}{\|z\|_{H^2(\Omega)}}. \quad (79)$$

Considering  $\Delta h = \Omega \setminus \Omega_h$ , we have

$$\begin{aligned} (u - u_h, -\Delta z + cz)_{L^2(\Omega)} &= (u - u_h, -\Delta z + cz)_{L^2(\Omega_h)} + (u - u_h, -\Delta z + cz)_{L^2(\Delta h)} \\ &= (\nabla_h(u - u_h), \nabla z)_{L^2(\Omega_h)} + (u - u_h, cz)_{L^2(\Omega_h)} \\ &\quad - \sum_{k=1}^K \int_{\partial T^k} (u - u_h) \frac{\partial z}{\partial n} \, ds + \widehat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z), \end{aligned}$$

with

$$\widehat{a}_{\Delta h}(u - u_h, z) = \int_{\Delta h} \nabla_h(u - u_h) \cdot \nabla z + (u - u_h)cz \, d\mathbf{x}, \quad (80)$$

and

$$b_{1h}(u - u_h, z) = - \int_{\partial\Omega} (u - u_h) \frac{\partial z}{\partial n} \, ds. \quad (81)$$

Using the equality (20) and attending that  $[[z]] = 0$  and  $[[\nabla z]] = 0$ , we may write

$$(u - u_h, -\Delta z + cz)_{L^2(\Omega)} = (\nabla_h(u - u_h), \nabla z)_{L^2(\Omega_h)} + (u - u_h, cz)_{L^2(\Omega_h)}$$

$$\begin{aligned}
& - \int_{\Gamma} \llbracket u - u_h \rrbracket \cdot \{\{\nabla z\}\} \, ds - \int_{\Gamma_0} \{\{u - u_h\}\} \llbracket \nabla z \rrbracket \, ds \\
& - \int_{\Gamma_0} \llbracket z \rrbracket \cdot \{\{\nabla_h(u - u_h)\}\} \, ds + \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket u - u_h \rrbracket \cdot \llbracket z \rrbracket \, ds \\
& + \widehat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z) \\
& = a_h(u - u_h, z) + \widehat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z) \\
& - \int_{\partial\Omega_h} (u - u_h) \frac{\partial z}{\partial n} \, ds. \tag{82}
\end{aligned}$$

In order to estimate the bilinear forms, consider  $\Pi_h(z)$  a continuous piecewise linear interpolate of  $z$  in  $\Omega$  at the vertexes of the mesh. Then, setting  $z_h = \Pi_h(z)$ , in  $\Omega_h$ , we have  $z_h \in \mathcal{V}_h$ . Therefore, since  $\Omega_h \subset \Omega$

$$a_h(u, z_h) = (f, z_h)_{L^2(\Omega_h)} = a_h(u_h, z_h). \tag{83}$$

Now, observe that

$$\begin{aligned}
\widehat{a}_{\Delta h}(u - u_h, z) &= \widehat{a}_{\Delta h}(u - u_h, z - \Pi_h(z)) + \widehat{a}_{\Delta h}(u - u_h, \Pi_h(z)) \\
&= \widehat{a}_{\Delta h}(u - u_h, z - \Pi_h(z)) + \sum_{k \in I^B} \int_{\Delta_k} -\Delta(u - u_h) \Pi_h(z) \, d\mathbf{x} \\
&\quad + c(u - u_h) \Pi_h(z) \, d\mathbf{x} + \sum_{k \in I^B} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} \frac{\partial(u - u_h)}{\partial n} \Pi_h(z) \, ds \\
&= \widehat{a}_{\Delta h}(u - u_h, z - \Pi_h(z)) + b_{2h}(u - u_h, \Pi_h(z)) + b_{3h}(u - u_h, \Pi_h(z)),
\end{aligned}$$

where

$$b_{2h}(u - u_h, \Pi_h(z)) = \sum_{k \in I^B} \int_{\Delta_k} -\Delta(u - u_h) \Pi_h(z) + c(u - u_h) \Pi_h(z) \, d\mathbf{x}, \tag{84}$$

$$b_{3h}(u - u_h, \Pi_h(z)) = \sum_{k \in I^B} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} \frac{\partial(u - u_h)}{\partial n} \Pi_h(z) \, ds. \tag{85}$$

Setting  $e_h(z) = z - \Pi_h(z)$ , we get

$$b_{4h}(u - u_h, e_h(z)) = \widehat{a}_{\Delta h}(u - u_h, z - \Pi_h(z)). \tag{86}$$

Then, we may write

$$\widehat{a}_{\Delta h}(u - u_h, z) = b_{2h}(u - u_h, \Pi_h(z)) + b_{3h}(u - u_h, \Pi_h(z)) + b_{4h}(u - u_h, e_h(z)). \tag{87}$$

Now, using the Galerkin orthogonality (72), since  $\Pi_h(z) \in \mathcal{V}_h$ , we may write

$$a_h(u - u_h, z) = a_h(u - u_h, z - \Pi_h(z) + \Pi_h(z)) = a_h(u - u_h, e_h(z)). \tag{88}$$



Considering

$$b_{5h}(u - u_h, z) = - \int_{\partial\Omega_h} (u - u_h) \frac{\partial z}{\partial n} ds, \quad (89)$$

and combining (87) and (88) into (82), we get

$$\begin{aligned} (u - u_h, -\Delta z + cz)_{L^2(\Omega)} &= a_h(u - u_h, e_h(z)) + b_{1h}(u - u_h, z) + b_{2h}(u - u_h, \Pi_h(z)) \\ &\quad + b_{3h}(u - u_h, \Pi_h(z)) + b_{4h}(u - u_h, e_h(z)) + b_{5h}(u - u_h, z). \end{aligned} \quad (90)$$

Thus, we obtain

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega)} &\leq C(\Omega) \frac{a_h(u - u_h, e_h(z)) + b_{1h}(u - u_h, z) + b_{2h}(u - u_h, \Pi_h(z))}{\|z\|_{H^2(\Omega)}} \\ &\quad + C(\Omega) \frac{b_{3h}(u - u_h, \Pi_h(z)) + b_{4h}(u - u_h, e_h(z)) + b_{5h}(u - u_h, z)}{\|z\|_{H^2(\Omega)}}. \end{aligned} \quad (91)$$

Using the boundedness inequality (35) and applying Theorem 4.1, we note that

$$a_h(u - u_h, e_h(z)) \leq C_b \|u - u_h\| \|e_h(z)\| \leq C_b \mathcal{C} h^N \|u\|_{H^{N+1}(\Omega)} \|e_h(z)\|. \quad (92)$$

From (76) and applying Lemma A.2 with  $j = 0, 1, 2$  and since  $h < 1$ , we establish

$$\|e_h(z)\| \leq C_{\Omega, z} h \|z\|_{H^2(\Omega)}, \quad (93)$$

with  $C_{\Omega, z}$  a mesh-independent constant. Thus, we rewrite (92) as

$$a_h(u - u_h, e_h(z)) \leq C_a h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (94)$$

where  $C_a = C_b \mathcal{C} C_{\Omega, z}$ .

Estimates for  $b_{ih}$ , with  $i = 1, 2, 3, 4, 5$ , can be established as follows (see Appendix B):

$$b_{1h}(u - u_h, z) \leq C_{b1} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (95)$$

$$b_{2h}(u - u_h, \Pi_h(z)) \leq C_{b2} h^{N+2} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (96)$$

$$b_{3h}(u - u_h, \Pi_h(z)) \leq C_{b3} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (97)$$

$$b_{4h}(u_h - u, e_h(z)) \leq C_{b4} h^{N+3/2} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (98)$$

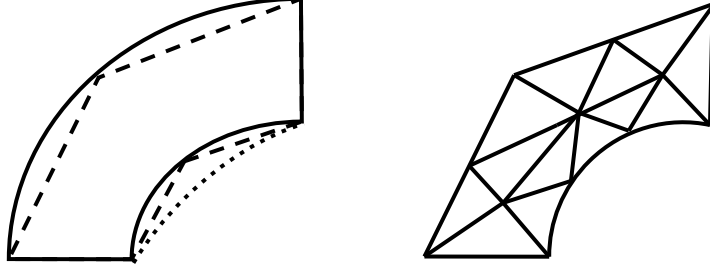
$$b_{5h}(u - u_h, z) \leq C_{b5} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}. \quad (99)$$

Finally, combining (94) with the estimates for  $b_{ih}$ , with  $i = 1, 2, 3, 4, 5$ , (95)–(99) into (91), owing to the fact  $h < 1$ , we obtain (77) with  $\mathcal{C}_0 = C(\Omega) (C_a + C_{b1} + C_{b2} + C_{b3} + C_{b4} + C_{b5})$ .  $\square$

## 4.2 Non-convex case

Now, we consider a non-convex domain  $\Omega$ . In this case, as the Galerkin orthogonality does not hold, there will be a non-zero residual  $a_h(u, v) - (f, v)_{L^2(\Omega_h)}$ ,  $v \in \mathcal{V}_h$ . We introduce a smooth domain  $\tilde{\Omega}$  close to  $\Omega$  such that  $\tilde{\Omega}_h = \Omega \cup \Omega_h \subset \tilde{\Omega}$  and  $\text{length}(\partial\tilde{\Omega}) - \text{length}(\partial\Omega) \leq \epsilon$ , for  $\epsilon$  sufficiently small (see Figure 3). Similarly to the norm defined in (7), we consider the following norm in  $\Omega \cap \Omega_h$  for  $u \in H^2(\mathcal{T}_h)$ ,

$$\left( \| \| u \| \|' \right)^2 = \sum_{k=1}^K \left( \| u \|_{H^1(T^k \cap \Omega)}^2 + h_k^2 |u|_{H^2(T^k \cap \Omega)}^2 \right) + \sum_{e \in \mathcal{E}_h} h_e^{-1} \| \llbracket u \rrbracket \|_{L^2(e \cap \Omega)}^2. \quad (100)$$



**Fig. 3: Left panel:** Example of a non-convex domain  $\Omega$  delimited by the solid lines, a polygonal mesh  $\Omega_h$  delimited by the dashed lines and  $\tilde{\Omega}$  delimited by the dotted lines. **Right panel:** example of  $\Omega \cap \Omega_h$ .

Consider  $f$  extended to  $\tilde{\Omega} \setminus \Omega$  such that  $f \in H^{N-1}(\tilde{\Omega})$  and we still denote the extended function by  $f$ . Assume a continuous extension of  $c$  to  $\tilde{\Omega} \setminus \Omega$ . Then, the following theorem holds:

**Theorem 4.3.** *Assume that there exists a function  $\tilde{u}$  defined in  $\tilde{\Omega}$  such that  $\tilde{u} \in H^{N+1}(\tilde{\Omega})$ ,  $\tilde{u}$  coincide with  $u$  in  $\Omega$ ,  $\tilde{u}$  satisfies (1) in  $\tilde{\Omega}$  and  $\tilde{u}$  vanishes on  $\partial\Omega$  in the sense of trace. Then, for  $h$  sufficiently small there exists a mesh-independent constant  $\tilde{C}$  such that*

$$\| \| u - u_h \| \|' \leq \tilde{C} h^N |\tilde{u}|_{H^{N+1}(\tilde{\Omega})}, \quad (101)$$

where  $\| \| \cdot \| \|'$  denotes the norm defined in (100).

*Proof.* We extend every  $v \in \mathcal{V}_h$  by zero on  $\tilde{\Omega} \setminus \Omega_h$ . Thanks to the properties of  $\tilde{u}$ , note that the proof of this theorem is based on the same arguments of the proof of Theorem 4.1. Since  $\| \| u - u_h \| \|' \leq \| \| \tilde{u} - u_h \| \|$ , we obtain (101).  $\square$

Note that, given a regular  $f$  in  $\tilde{\Omega}$ , the existence of an associated  $\tilde{u}$  satisfying the above assumptions is not ensured. Thus, let us consider that  $f$  vanishes in  $\Omega_h \setminus \Omega$ . Denoting by  $\tilde{u}$  the regular extension of  $u$  to  $\tilde{\Omega}$  such that  $\tilde{u} \in H^{N+1}(\tilde{\Omega})$  and  $\tilde{u}|_{\Omega} = u$ , in the following theorems we estimate the non-zero residual  $a_h(\tilde{u}, v) - (f, v)_{L^2(\Omega_h)}$  considering two different approaches.

**Theorem 4.4.** *Let  $u \in H^{N+1}(\Omega)$  be the solution of (1)–(2). Provided  $h$  sufficiently small, there exists mesh independent constants  $\tilde{C}_1$  and  $\widehat{C}_0$  such that the numerical solution  $u_h$  satisfies*

$$\| \|u - u_h\| \|' \leq \tilde{C}_1 h^N |\tilde{u}|_{H^{N+1}(\tilde{\Omega})} + \widehat{C}_0 h^{5/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})}, \quad (102)$$

where  $\tilde{u}$  is the regular extension of  $u$  to  $\tilde{\Omega}$  such that  $\tilde{u} \in H^{N+1}(\tilde{\Omega})$ .

*Proof.* Note that

$$\| \|u - u_h\| \|' \leq \| \tilde{u} - u_h \| \| \leq \| \tilde{u} - I_h(\tilde{u}) \| \| + \| u_h - I_h(\tilde{u}) \| \|, \quad (103)$$

with  $I_h(\tilde{u})$  the  $\mathcal{P}_N$ -interpolate of  $\tilde{u}$  at the nodes associated with  $\mathcal{W}_h$ . From the inf-sup condition (Theorem 3.2), we get

$$\| \|u_h - I_h(\tilde{u})\| \| \leq \frac{1}{\alpha} \sup_{v \in \mathcal{V}_h \setminus \{0\}} \frac{a_h(u_h - I_h(\tilde{u}), v)}{\| \|v\| \|}. \quad (104)$$

Adding and subtracting  $\tilde{u}$  in the first argument of  $a_h$  yields

$$a_h(u_h - I_h(\tilde{u}), v) \leq |a_h(\tilde{u} - I_h(\tilde{u}), v)| + |a_h(u_h - \tilde{u}, v)|. \quad (105)$$

Following the same argument as in the proof of Theorem 4.1, using the boundedness inequality (35) we obtain

$$|a_h(\tilde{u} - I_h(\tilde{u}), v)| \leq \tilde{C}_b \| \tilde{u} - I_h(\tilde{u}) \| \| \|v\| \|.$$

As the Galerkin orthogonality does not hold, we need to estimate  $a_h(u_h - \tilde{u}, v)$ . First, note that

$$\begin{aligned} a_h(\tilde{u}, v) &= (\nabla_h \tilde{u}, \nabla_h v)_{L^2(\Omega_h)} + (c\tilde{u}, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{ \{ \nabla_h \tilde{u} \} \} \, ds \\ &= -(\nabla_h \cdot \nabla_h \tilde{u}, v)_{L^2(\Omega_h)} + \int_{\Gamma} \llbracket v \rrbracket \cdot \{ \{ \nabla_h \tilde{u} \} \} \, ds + \int_{\Gamma_0} \{ \{ v \} \} \llbracket \nabla_h \tilde{u} \rrbracket \, ds \\ &\quad + (c\tilde{u}, v)_{L^2(\Omega_h)} - \int_{\Gamma_0} \llbracket v \rrbracket \cdot \{ \{ \nabla_h \tilde{u} \} \} \, ds \\ &= -(\nabla \cdot \nabla \tilde{u}, v)_{L^2(\Omega_h)} + (c\tilde{u}, v)_{L^2(\Omega_h)} \\ &= \sum_{k \in \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(T^k)} + \sum_{k \notin \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(T^k)}. \end{aligned}$$

Then,

$$a_h(u_h - \tilde{u}, v) = (f, v)_{L^2(\Omega_h)} - a_h(\tilde{u}, v)$$

$$\begin{aligned}
&= \sum_{k \in \mathcal{Q}^B} (f, v)_{L^2(\Delta_k)} - \sum_{k \in \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(\Delta_k)} \\
&= \sum_{k \in \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(\Delta_k)},
\end{aligned}$$

since  $f = 0$  in  $\Omega_h \setminus \Omega$ . Now, following the same argument as in Appendix B for the estimate for  $b_{2h}$ , we obtain

$$\begin{aligned}
\left| \sum_{k \in \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(\Delta_k)} \right| &\leq \sum_{k \in \mathcal{Q}^B} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\Delta_k)} \|v\|_{L^2(\Delta_k)} \\
&\leq \sum_{k \in \mathcal{Q}^B} \sqrt{C_{\partial\Omega} h_k^{3/2}} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\Delta_k)} \|v\|_{L^\infty(\Delta_k)}. \quad (106)
\end{aligned}$$

Recalling that  $v = 0$  on  $\partial\Omega_h$ , by the Mean Value Theorem and Proposition A.1, we get

$$|v(P)| \leq C_{\partial\Omega} h_k^2 \|\nabla v\|_{L^\infty(T^k \cup \Delta_k)}, \forall P \in \Delta_k, T^k, k \in I^B. \quad (107)$$

Considering the inequality (107) and Lemma A.1, we may write

$$\|v\|_{L^\infty(\Delta_k)} \leq C_{\partial\Omega} C_J h_k \|\nabla v\|_{L^2(T^k)}. \quad (108)$$

Then, replacing (108) in (106), we get

$$\begin{aligned}
\left| \sum_{k \in \mathcal{Q}^B} (-\Delta \tilde{u} + c\tilde{u}, v)_{L^2(\Delta_k)} \right| &\leq \sum_{k \in \mathcal{Q}^B} C_{\partial\Omega}^{3/2} C_J h_k^{5/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\Delta_k)} \|\nabla v\|_{L^2(T^k)} \\
&\leq C_{\partial\Omega}^{3/2} C_J h^{5/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \|\nabla v\|_{L^2(\Omega_h)} \\
&\leq C_{\partial\Omega}^{3/2} C_J h^{5/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \|v\| \quad (109)
\end{aligned}$$

and so the inequality (104) may be written as

$$\|u_h - I_h(\tilde{u})\| \leq \frac{\tilde{C}_b}{\alpha} \|\tilde{u} - I_h(\tilde{u})\| + \frac{C_{\partial\Omega}^{3/2} C_J h^{5/2}}{\alpha} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})}.$$

Finally, replacing the previous inequality in (103), we establish

$$\|u - u_h\|' \leq \left(1 + \frac{\tilde{C}_b}{\alpha}\right) \|\tilde{u} - I_h(\tilde{u})\| + \frac{C_{\partial\Omega}^{3/2} C_J h^{5/2}}{\alpha} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})}.$$

Now, following similar arguments as in (75) and using Lemma A.2, considering  $p = 2$ ,  $m = N + 1$ ,  $j = 0, 1, 2$  and  $h < 1$ , we establish

$$\|\tilde{u} - I_h(\tilde{u})\| \leq \tilde{C}_a h^N |\tilde{u}|_{H^{N+1}(\tilde{\Omega})}.$$

Thus, (102) holds with  $\tilde{C}_1 = \tilde{C}_a \left(1 + \tilde{C}_b/\alpha\right)$  and  $\hat{C}_0 = C_{\partial\Omega}^{3/2} C_J/\alpha$ .  $\square$

We can estimate a similar result by considering  $\|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})}$ .

**Theorem 4.5.** *Let  $u \in H^{N+1}(\Omega)$  be the solution of (1)–(2). Provided  $h$  sufficiently small, there exists a mesh independent constants  $\tilde{C}_1$  and  $C'_0$  such that the numerical solution  $u_h$  satisfies*

$$\| \|u - u_h\| \|' \leq \tilde{C}_1 h^N |\tilde{u}|_{H^{N+1}(\tilde{\Omega})} + C'_0 h^{7/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})}, \quad (110)$$

where  $\tilde{u}$  is the regular extension of  $u$  to  $\tilde{\Omega}$  such that  $\tilde{u} \in H^{N+1}(\tilde{\Omega})$ .

*Proof.* According to the Sobolev embedding Theorem [1], since  $\tilde{u} \in H^{N+1}(\tilde{\Omega})$ , then  $\Delta\tilde{u} \in L^\infty(\tilde{\Omega})$ . Now, following the same steps as in the proof of the Theorem 4.4 up to equation (106), and applying the Cauchy-Schwarz inequality, we get

$$\begin{aligned} \left| \sum_{k \in \mathcal{Q}^B} (-\Delta\tilde{u} + c\tilde{u}, v)_{L^2(\Delta_k)} \right| &\leq \sum_{k \in \mathcal{Q}^B} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^2(\Delta_k)} \|v\|_{L^2(\Delta_k)} \\ &\leq \sum_{k \in \mathcal{Q}^B} \sqrt{C_{\partial\Omega} h_k^{3/2}} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\Delta_k)} \sqrt{C_{\partial\Omega} h_k^{3/2}} \|v\|_{L^\infty(\Delta_k)} \\ &\leq C_{\partial\Omega}^2 C_J h^{7/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})} \sum_{k \in \mathcal{Q}^B} h_k^{1/2} \|\nabla v\|_{L^2(T^k)} \\ &\leq C_{\partial\Omega}^2 C_J h^{7/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})} \left( \sum_{k \in \mathcal{Q}^B} h_k \right)^{1/2} \|v\| \\ &\leq C_{\partial\Omega}^2 C_J C(\partial\Omega) h^{7/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})} \|v\|, \end{aligned} \quad (111)$$

assuming that exists a mesh-independent constant  $C(\partial\Omega)$  such that  $\sum_{k \in \mathcal{Q}^B} h_k \leq C^2(\partial\Omega)$ . Thus, (110) holds with  $\tilde{C}_1 = \tilde{C}_a \left(1 + \tilde{C}_b/\alpha\right)$  and  $C'_0 = C_{\partial\Omega}^2 C_J C(\partial\Omega)/\alpha$ .  $\square$

Note that, using the Theorem 4.4 for  $N = 2$  and a suitable constant  $\tilde{C}_2$ , we get

$$\| \|u - u_h\| \|' \leq \tilde{C}_2 h^2 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \quad (112)$$

and considering the Theorem 4.5 for  $N = 3$  and a suitable constant  $\tilde{C}_3$ , we obtain

$$\| \|u - u_h\| \|' \leq \tilde{C}_3 h^3 \left( |\tilde{u}|_{H^4(\tilde{\Omega})} + h^{1/2} \|-\Delta\tilde{u} + c\tilde{u}\|_{L^\infty(\tilde{\Omega})} \right). \quad (113)$$

We now establish error estimates in the  $L^2$ - norm in the case of a non-convex domain  $\Omega$ , by requiring more regularity from the solution  $u$ . The optimal convergence can be achieved not only when  $u$  is more regular but also when the computational domain  $\Omega_h$  approximates better the physical domain  $\Omega$ , i.e., when  $\Omega_h \setminus \Omega$  is of order

$h^q$ , with  $q > 2$  [23]. However, unless the assumptions of the Theorem 4.3 hold, with our definition of  $\Omega_h$ , optimally is not attained for  $N > 2$ , see Proposition A.2.

**Theorem 4.6.** *Let  $N = 2$ . Assume that  $\Omega$  is not convex and  $u \in H^{3+r}(\Omega)$  is the solution of (1)–(2), for  $r = 1/2 + \epsilon$ , with  $\epsilon > 0$  arbitrarily small. Then for  $h$  sufficiently small, the following error estimate holds:*

$$\|u - u_h\|_{L^2(\Omega \cap \Omega_h)} \leq \tilde{C}_0 h^3 \left( \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right), \quad (114)$$

where  $\tilde{C}_0$  is a mesh independent constant and  $\tilde{u}$  is the regular extension of  $u$  to  $\tilde{\Omega}$  such that  $\tilde{u} \in H^{3+r}(\tilde{\Omega})$ , for  $r = 1/2 + \epsilon$ , with  $\epsilon > 0$ .

*Proof.* Recalling the proof of Theorem 4.2, let  $z \in H_0^1(\Omega)$  be the solution of

$$\begin{aligned} -\Delta z(\mathbf{x}) + c(\mathbf{x})z(\mathbf{x}) &= u(\mathbf{x}) - u_h(\mathbf{x}), & \mathbf{x} \in \Omega, \\ z(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\Omega. \end{aligned}$$

Then, considering  $\Omega = (\Omega \cap \Omega_h) \cup \Delta_h$  and using integration by parts we obtain

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega)} &\leq C(\Omega) \frac{(u - u_h, -\Delta z + cz)_{L^2(\Omega)}}{\|z\|_{H^2(\Omega)}} \\ &\leq C(\Omega) \frac{a'_h(u - u_h, z) + \hat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z) - \int_{\partial\Omega_h \cap \Omega} (u - u_h) \frac{\partial z}{\partial n} ds}{\|z\|_{H^2(\Omega)}}, \end{aligned} \quad (115)$$

where  $\hat{a}_{\Delta h}$  and  $b_{1h}$  are defined in (80) and (81), respectively and

$$\begin{aligned} a'_h(z, u - u_h) &= (\nabla_h(u - u_h), \nabla z)_{L^2(\Omega \cap \Omega_h)} + (u - u_h, cz)_{L^2(\Omega \cap \Omega_h)} \\ &\quad - \int_{\Gamma_0} \llbracket z \rrbracket \cdot \{ \{ \nabla_h(u - u_h) \} \} ds - \int_{\Gamma_0} \llbracket u - u_h \rrbracket \cdot \{ \{ \nabla z \} \} ds \\ &\quad + \int_{\Gamma_0} \frac{\eta}{h_e} \llbracket u - u_h \rrbracket \cdot \llbracket z \rrbracket ds. \end{aligned}$$

Thus, since  $\|u - u_h\|_{L^2(\Omega \cap \Omega_h)} \leq \|u - u_h\|_{L^2(\Omega)}$ , we have

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega \cap \Omega_h)} &\leq C(\Omega) \frac{a'_h(u - u_h, z) + \hat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z)}{\|z\|_{H^2(\Omega)}} \\ &\quad - C(\Omega) \frac{\int_{\partial\Omega_h \cap \Omega} (u - u_h) \frac{\partial z}{\partial n} ds}{\|z\|_{H^2(\Omega)}}. \end{aligned} \quad (116)$$

Since  $f = 0$  in  $\Omega_h \setminus \Omega$ ,  $\forall v \in \mathcal{V}_h$  we get

$$a_h(u_h, v_h) = \int_{\Omega \cap \Omega_h} -(\Delta u + cu)v_h d\mathbf{x} = - \int_{\partial\Omega \cap \Omega_h} \frac{\partial u}{\partial n} v_h ds + a'_h(u, v_h). \quad (117)$$

Thus, using (117), for  $v \in \mathcal{V}_h$ , we may write

$$-a'_h(u - u_h, v_h) + b_{6h}(u_h, v_h) + b_{7h}(u - u_h, v_h) = 0, \quad (118)$$

where

$$b_{6h}(u_h, v_h) = - \sum_{k \in \mathcal{Q}^B} \int_{\Delta_k} (-\Delta u_h + cu_h) v_h \, d\mathbf{x}, \quad (119)$$

and

$$b_{7h}(u - u_h, v_h) = \sum_{k \in \mathcal{Q}^B} \int_{\partial\Omega \cap T^k} \frac{\partial(u - u_h)}{\partial n} v_h \, ds. \quad (120)$$

Now, considering  $v_h = \Pi_h(z)$  in (118) and

$$b_{8h}(u - u_h, z) = - \int_{\partial\Omega_h \cap \Omega} (u - u_h) \frac{\partial z}{\partial n} \, ds, \quad (121)$$

we may write (116) as

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega \cap \Omega_h)} &\leq C(\Omega) \frac{a'_h(u - u_h, e_h(z)) + \widehat{a}_{\Delta h}(u - u_h, z) + b_{1h}(u - u_h, z)}{\|z\|_{H^2(\Omega)}} \\ &\quad + C(\Omega) \frac{b_{6h}(u_h, \Pi_h(z)) + b_{7h}(u - u_h, \Pi_h(z)) + b_{8h}(u - u_h, z)}{\|z\|_{H^2(\Omega)}}, \end{aligned}$$

where  $e_h(z) = z - \Pi_h(z)$ . Recalling  $b_{2h}$ ,  $b_{3h}$  and  $b_{4h}$  given by (84), (85) and (86), respectively, and adding and subtracting  $\Pi_h(z)$  in the second argument of the bilinear form  $\widehat{a}_{\Delta h}$ , note that

$$\begin{aligned} \widehat{a}_{\Delta h}(u - u_h, z) + b_{7h}(u - u_h, \Pi_h(z)) &= \widehat{a}_{\Delta h}(u - u_h, e_h(z)) + \widehat{a}_{\Delta h}(u - u_h, \Pi_h(z)) \\ &\quad + b_{7h}(u - u_h, \Pi_h(z)) \\ &= b_{4h}(u - u_h, e_h(z)) + b_{2h}(u - u_h, \Pi_h(z)) \\ &\quad + b_{3h}(u - u_h, \Pi_h(z)). \end{aligned}$$

Thus, we may write

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega \cap \Omega_h)} &\leq C(\Omega) \frac{b_{1h}(u - u_h, z) + b_{2h}(u - u_h, \Pi_h(z)) + b_{3h}(u - u_h, \Pi_h(z))}{\|z\|_{H^2(\Omega)}} \\ &\quad + C(\Omega) \frac{b_{4h}(u - u_h, e_h(z)) + b_{6h}(u_h, \Pi_h(z)) + b_{8h}(u - u_h, z)}{\|z\|_{H^2(\Omega)}} \\ &\quad + C(\Omega) \frac{a'_h(u - u_h, e_h(z))}{\|z\|_{H^2(\Omega)}}. \end{aligned} \quad (122)$$

We are left to estimate upper bounds for the bilinear forms  $a'_h$ ,  $b_{6h}$  and  $b_{8h}$ . Using the boundedness of the bilinear form and applying the Theorem 4.4 with  $N = 2$  (see

inequality (112)), we first note that

$$a'_h(u - u_h, e_h(z)) \leq C_b \tilde{C}_2 h^2 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \|e_h(z)\|'. \quad (123)$$

On the other hand, similiar to (93), we get

$$\|e_h(z)\|' \leq C'_{\Omega, z} h |z|_{H^2(\Omega)},$$

with  $C'_{\Omega, z}$  a mesh-independent constant. Thus, we rewrite (123) as

$$a'_h(u - u_h, e_h(z)) \leq C'_a h^3 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \|z\|_{H^2(\Omega)}, \quad (124)$$

where  $C'_a = C_b \tilde{C}_2 C'_{\Omega, z}$ . An estimate for an upper bound for  $b_{6h}$  can be established as follows (see Appendix B):

$$b_{6h}(u_h, \Pi_h(z)) \leq \tilde{C}_{b6} h^{7/2} \left( \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + h^{1/2} \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \right) \|z\|_{H^2(\Omega)}. \quad (125)$$

Following similar arguments as in the estimates for  $b_{5h}$  and applying the error estimate (112), we get

$$\begin{aligned} b_{8h}(u - u_h, z) &\leq \tilde{C}_{b8} h^3 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \|z\|_{H^2(\Omega)} \\ &\leq \tilde{C}_{b8} h^3 \left( \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \|z\|_{H^2(\Omega)}. \end{aligned} \quad (126)$$

Estimates for  $b_{1h}$ ,  $b_{2h}$ ,  $b_{3h}$  and  $b_{4h}$  can be obtained by following the same arguments as in the proof of Theorem 4.2, taking  $N = 2$  and noticing that in this case, we apply Theorem 4.4 instead of Theorem 4.1 (see Appendix B). Thus,  $|u|_{H^3(\Omega)}$  is replaced by  $|\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})}$  in the estimates for  $b_{ih}$ ,  $i = 1, 2, 3, 4$ .

Finally, combining (124), (125) and (126) with the estimates for  $b_{ih}$ , with  $i = 1, 2, 3, 4$ , into (122), owing to the fact  $h < 1$ , we obtain (114) with  $\tilde{C}_0 = C(\Omega) \left( C'_a + \tilde{C}_{b1} + \tilde{C}_{b2} + \tilde{C}_{b3} + \tilde{C}_{b4} + \tilde{C}_{b6} + \tilde{C}_{b8} \right)$ , where  $\tilde{C}_{bi}$  is the constant in the estimate for  $b_{ih}$ .  $\square$

## 5 Numerical Results

Let denote by  $u_h$  an approximation of the solution  $u$  for a given mesh  $\mathcal{T}_h$  and  $\|u - u_h\|$  the norm of the error. The method is of convergence order  $p$  if one has asymptotically

$$\|u - u_h\| \leq Ch^p,$$



with  $C$  a real constant independent of  $h$ . The errors are assessed at the node points of the elements,  $T^k \in \mathcal{T}_h$ ,  $k = 1, \dots, K$ . We compute the  $L^2$ -errors

$$E_2(\mathcal{T}_h) = \|u - u_h\|_{L^2(\Omega_h)} = \sqrt{\sum_{k=1}^K \|u - u_h\|_{L^2(T^k)}^2}.$$

Recall that

$$\|u_h\|_{L^2(T^k)}^2 = \left( u_h^k, u_h^k \right)_{L^2(T^k)} = \int_{T^k} \sum_{i=1}^{N_p} \sum_{j=1}^{N_p} u_i^k u_j^k \ell_i^k(\mathbf{x}) \ell_j^k(\mathbf{x}) \, d\mathbf{x} = (\mathbf{u}^k)^T M^k \mathbf{u}^k,$$

where  $M_{ij}^k = \left( \ell_i^k, \ell_j^k \right)_{L^2(T^k)}$ .

Consider two different meshes, denoted as  $\mathcal{T}_{h_1}$  and  $\mathcal{T}_{h_2}$ , whose the corresponding numerical solutions are denoted as  $u_{h_1}$  and  $u_{h_2}$ , respectively. Then, the convergence order between two successively finer meshes is determined as

$$\mathcal{O}_2(\mathcal{T}_{h_1}, \mathcal{T}_{h_2}) = \frac{\log(E_2(\mathcal{T}_{h_1})/E_2(\mathcal{T}_{h_2}))}{\log(h_1/h_2)}.$$

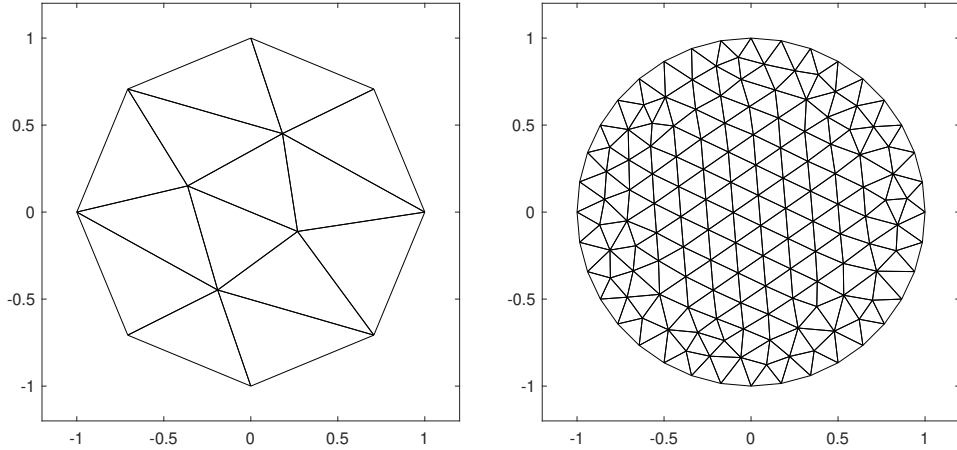
For the numerical tests, we consider  $c(\mathbf{x}) = 1$ .

## 5.1 Disk domain

Consider the reaction-diffusion equation on a disk of radius  $R = 1$  with a homogeneous Dirichlet boundary condition. An analytical solution is manufactured for problem (1)–(2) and is given as  $u(x, y) = x \sin(1 - x^2 - y^2)$ , from which the corresponding source term is deduced. Simulations are carried out with successively finer meshes generated by Gmsh (version 4.6.0) [18] (see Figure 4).

Simulations are first performed for the classical DG method prescribing the homogeneous Dirichlet boundary condition at the nodes of the computational boundary (the edges of the mesh). More precisely, each node of the computational boundary has a corresponding node on the real boundary where the Dirichlet boundary condition is prescribed. For the classical DG method, the value evaluated at the physical boundary point is used at the corresponding node on the computational boundary. The results, reported in Table 1, demonstrate the accuracy deterioration from such a geometrical mismatch without any specific treatment for curved boundaries, and the error convergence is limited to the second-order.

Table 2 reports the  $L^2$ -errors and convergence orders for the DG-ROD method taking  $N = 2, 3, 4$ . As observed, the quality of the approximations obtained with the DG-ROD method is in good agreement with the theoretical results.



**Fig. 4:** Unstructured meshes generated for the disk domain. Mesh with  $K = 14$  and  $h = 9.34\text{E}-01$  (left panel) and mesh with  $K = 262$  and  $h = 2.34\text{E}-01$  (right panel).

**Table 1:** Errors and convergence orders for the classical DG method in the disk domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
14	9.34E-01	9.21E-02	—	8.98E-02	—	8.76E-02	—
64	4.70E-01	2.47E-02	1.9	2.43E-02	1.9	2.40E-02	1.9
262	2.34E-01	4.89E-03	2.3	4.81E-03	2.3	4.78E-03	2.3
1096	1.13E-01	1.08E-03	2.1	1.07E-03	2.1	1.07E-03	2.1
4136	5.69E-02	2.67E-04	2.0	2.66E-04	2.0	2.65E-04	2.0

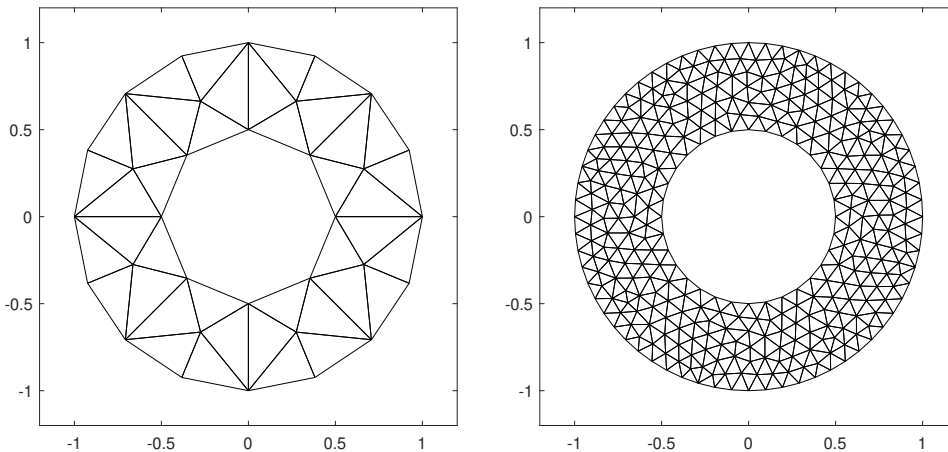
**Table 2:** Errors and convergence orders for the DG-ROD method in the disk domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
14	9.34E-01	1.79E-02	—	9.69E-04	—	8.93E-04	—
64	4.70E-01	7.36E-04	4.6	7.00E-05	3.8	1.65E-05	5.8
262	2.34E-01	5.64E-05	3.7	4.59E-06	3.9	3.55E-07	5.5
1096	1.13E-01	3.47E-06	3.8	2.53E-07	4.0	8.16E-09	5.2
4136	5.69E-02	2.52E-07	3.8	1.67E-08	4.0	2.43E-10	5.1

## 5.2 Annulus domain

Now, we consider an annulus domain with inner radius  $R_I = 0.5$  and outer radius  $R_E = 1$  meshed with triangular elements (see Figure 5). The analytic solution corresponds

to the manufactured solution  $u(x, y) = \log(x^2 + y^2)$  and the boundaries are prescribed with constant Dirichlet boundary conditions. The numerical simulations are carried out with successively finer meshes generated by Gmsh. As for the previous test case, simulations are firstly performed for the classical DG method and the results are reported in Table 3. On the other hand, Table 4 reports the errors and convergence rate for the DG-ROD method where the optimal convergence orders are recovered due to the polynomial reconstruction of the boundary conditions. In this case, the solution satisfies the conditions of Theorem 4.3 for non-convex domains. Thus, the method recovers the optimal convergence orders for  $N > 2$ .



**Fig. 5:** Unstructured mesh generated for the annulus domain. Mesh with  $K = 40$  and  $h = 5.00\text{E}-01$  (left panel) and mesh with  $K = 608$  and  $h = 1.31\text{E}-01$  (right panel).

**Table 3:** Errors and convergence orders for the classical DG method in the annulus domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
40	5.00E-01	8.90E-02	—	9.02E-02	—	9.09E-02	—
144	2.57E-01	2.25E-02	2.1	2.27E-02	2.1	2.28E-02	2.1
608	1.31E-01	5.71E-03	2.0	5.75E-03	2.0	5.76E-03	2.0
2576	6.45E-02	1.29E-03	2.1	1.30E-03	2.1	1.30E-03	2.1
10226	3.18E-02	3.23E-04	2.0	3.24E-04	2.0	3.24E-04	2.0

### 5.3 Rose-shaped domain

Consider a geometry generated by applying a diffeomorphic transformation to an annular domain, denoted as  $\Omega'$ , with interior and exterior physical boundaries with

**Table 4:** Errors and convergence orders for the DG-ROD method in the annulus domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
40	5.00E-01	4.48E-03	—	4.45E-04	—	8.17E-05	—
144	2.57E-01	5.16E-04	3.2	3.67E-05	3.8	2.74E-06	5.1
608	1.31E-01	3.96E-05	3.8	1.51E-06	4.7	6.01E-08	5.6
2576	6.45E-02	2.64E-06	3.8	4.60E-07	4.9	9.62E-10	5.9
10226	3.18E-02	2.03E-07	3.6	1.83E-09	4.6	2.26E-11	5.3

radius  $r_I$  and  $r_E$ , respectively. The diffeomorphic transformation corresponds to the mapping  $\Omega' \rightarrow \Omega$ , where  $\Omega$  is the rose-shaped domain, given in polar coordinates

$$\Omega' \rightarrow \Omega : \begin{bmatrix} r' \\ \theta' \end{bmatrix} \rightarrow \begin{bmatrix} r \\ \theta \end{bmatrix} = \begin{bmatrix} R(r', \theta'; \alpha, \beta) \\ \theta \end{bmatrix}, \quad (127)$$

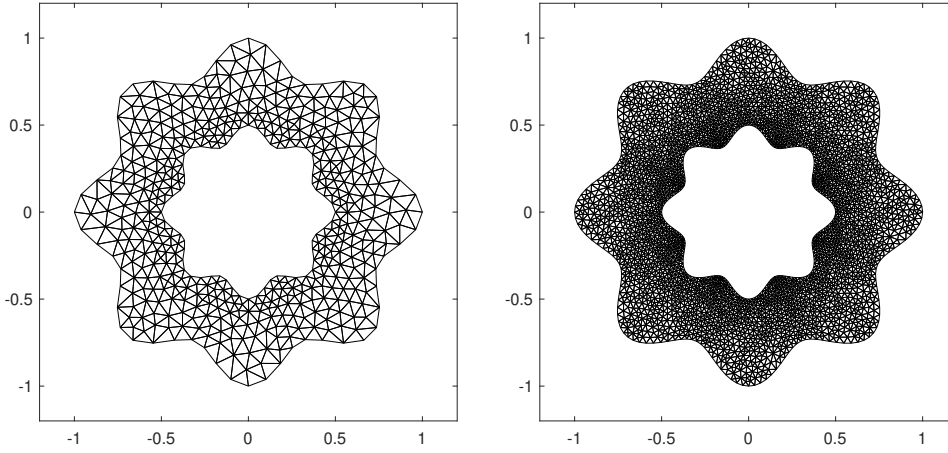
where  $\alpha$  is the number of petals and function  $R(r', \theta') := R(r', \theta'; \alpha, \beta)$  corresponds to a periodic radius perturbation of magnitude in  $[-\beta, \beta]$ , with  $\beta \in \mathbb{R}$ , given as

$$R(r', \theta'; \alpha, \beta) = r' (1 - \beta + \beta \cos(\alpha\theta')). \quad (128)$$

Thus, the interior and exterior physical boundaries parametrisation are given as  $R_I := R(r_I, \theta)$  and  $R_E := R(r_E, \theta)$ , respectively. The analytic solution corresponds to the manufactured solution  $u(x, y) = \log(x^2 + y^2)$ . In this test case, the interior and exterior boundaries are prescribed with a non-constant Dirichlet boundary condition.

We consider  $r_I = 0.5$ ,  $r_E = 1$ , the number of petals is  $\alpha = 8$ , and the perturbation magnitude is  $\beta = 0.1$ . The rose-shaped domain is meshed with triangular elements (see Figure 6). The reaction-diffusion equation is solved and the approximate solution is compared with the exact solution. The numerical simulations are carried out with successively finer meshes generated by Gmsh. Table 5 reports the errors and converge rates for the classical DG method. The results confirm the accuracy deterioration due to the lack of specific treatment for curved boundaries, where the error convergence is limited to the second order. On the other hand, the results for the DG-ROD method are reported in Table 6, where the optimal convergence orders are recovered. The method behaves similarly to the previous test case, where the optimal convergence orders can be achieved for  $N > 2$ .

For further numerical results, the authors refer to [30]. Recall that in [30] the overall DG-ROD method was obtained by considering an iterative procedure of the DG method and the polynomial reconstruction and for the numerical results the authors computed the  $L^\infty$ -errors and  $L^1$ -errors. In the mentioned paper, the authors considered different approaches to obtain the nodes associated with  $\mathcal{W}_h$  located on  $\partial\Omega$ , namely the intuitive construction of the nodes lying on normals to edges of  $\partial\Omega_h$ . In this work, the construction of such nodes for the numerical results is described in item (4).



**Fig. 6:** Unstructured mesh generated for the rose-shaped domain. Mesh with  $K = 888$  and  $h = 1.41\text{E}-01$  (left panel) and mesh with  $K = 6912$  and  $h = 5.25\text{E}-02$  (right panel).

**Table 5:** Errors and convergence orders for the classical DG method in the rose-shaped domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
3072	7.50E-02	1.85E-03	—	1.82E-03	—	1.81E-03	—
4792	6.33E-02	1.18E-03	2.7	1.16E-03	2.6	1.16E-03	2.6
6912	5.25E-02	8.16E-04	2.0	8.08E-04	2.0	8.07E-04	2.0
10478	4.27E-02	5.26E-04	2.1	5.22E-04	2.1	5.21E-04	2.1
15346	3.52E-02	3.56E-04	2.0	3.54E-04	2.0	3.53E-04	2.0

**Table 6:** Errors and convergence orders for the DG-ROD method in the rose-shaped domain with the Dirichlet boundary conditions.

$K$	$h$	$N = 2$		$N = 3$		$N = 4$	
		$E_2$	$O_2$	$E_2$	$O_2$	$E_2$	$O_2$
3072	7.50E-02	2.05E-06	—	2.84E-08	—	4.58E-10	—
4792	6.33E-02	9.64E-07	4.5	1.07E-08	5.8	1.43E-10	6.9
6912	5.25E-02	5.90E-07	2.6	5.25E-09	3.8	5.83E-11	4.8
10478	4.27E-02	2.98E-07	3.3	2.14E-09	4.4	1.94E-11	5.3
15346	3.52E-02	1.78E-07	2.7	1.03E-09	3.8	8.48E-12	4.3

## 6 Conclusions

We have discussed a modified DG scheme, defined on a polygonal mesh  $\Omega_h$  for solving boundary value problems on a two-dimensional curved boundary domain  $\Omega$ , where piecewise linear elements approximate the curved boundaries. The DG-ROD method corrects the error resulting from the approximation of the curved boundary  $\partial\Omega$  by the computational boundary  $\partial\Omega_h$  by means of polynomial reconstructions of the boundary conditions. This correction is reflected in the variational formulation of the problem.

We present a study on the existence and uniqueness of the solution for the reaction-diffusion equation with homogeneous Dirichlet boundary conditions. We provided a complete mathematical analysis of the convergence in the natural norm of the DG method, as well as  $L^2$ -error estimates, considering convex and non-convex domains. For the convex domains, we prove that the DG-ROD solution exhibits an optimal  $\mathcal{O}(h^{N+1})$  convergence rate in the  $L^2$ -norm when  $N$ -degree piecewise polynomials are used, under certain regularity conditions on the solution. For non-convex domains, unless the solution satisfies certain regularity conditions and the computational domain  $\Omega_h$  approximates better the physical domain  $\Omega$ , i.e., when  $\Omega_h \setminus \Omega$  is of order  $h^q$ , with  $q > 2$ , optimally is not attained for  $N > 2$ . In other words, the error is affected by the geometrical mismatch of order  $\mathcal{O}(h^2)$  between the curved physical boundaries and the associated piecewise linear representation. Then, for non-convex domains, we prove that the DG-ROD solution exhibits an optimal  $\mathcal{O}(h^3)$  convergence rate in the  $L^2$ -norm when piecewise polynomials of degree  $N = 2$  are used, under certain regularity conditions on the solution.

The sharpness of the theoretical results is confirmed by a series of numerical experiments in convex and non-convex domains. It is important to highlight that the assumption on the mesh size  $h$  is just a sufficient condition for the formal analysis given in this work. Good numerical results can be obtained even for coarse meshes. For example, in the test case 5.1, the polynomial reconstructions of the boundary condition correct the error from approximating the curved boundary with a polygonal boundary even for a mesh with  $K = 14$ , where the corresponding mesh size  $h$  has the same order of magnitude as the radius of the disk.

Extensions of this work considering nonlinear equations and time-dependent problems are challenging and this will be carried out in the future. For future work, we also plan to extend this approach to other boundary conditions (Neumann, Robin) and derive error estimates for the respective problems.

**Acknowledgements.** The authors were financially supported by the Fundação para a Ciência e a Tecnologia (Portuguese Foundation for Science and Technology) under the scope of the projects UIDB/00324/2020 (<https://doi.org/10.54499/UIDB/00324/2020>) and UIDP/00324/2020 (<https://doi.org/10.54499/UIDP/00324/2020>) (Centre for Mathematics of the University of Coimbra). M. Santos acknowledges FCT for the support under the Ph.D. scholarship UI/BD/153816/2022.

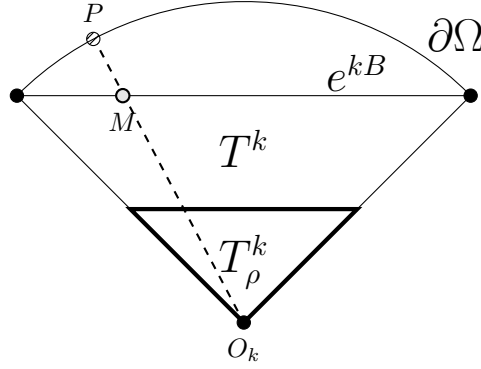
## Declarations

**Conflict of interest.** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Appendix A Technical tools

The following assumptions and results can be found (except Lemma A.4) in [28].

**Assumption A.1.** Let  $T_\rho^k$  be the homothetic transformation of  $T^k$  with center  $O_k$  (the vertex of  $T^k$  not located on  $\partial\Omega_h$ ) and ratio  $\rho < 1$ . Consider  $h$  small enough for the intersection  $P$  with  $\partial\Omega$  of a straight line joining any point of  $T_\rho^k$  to a point  $M \in e^{kB}$  is uniquely defined for all  $T^k$ ,  $k \in I^B$  (see Figure A1).



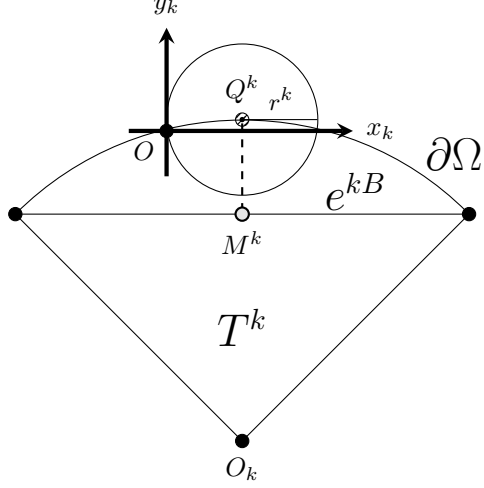
**Fig. A1:** Example of a homothetic transformation of  $T^k$  with center  $O_k$  and  $\rho = 1/2$ .

Let  $Q^k$  be the closest intersection with  $\partial\Omega$  of the perpendicular to  $e^{kB}$  passing through its mid-point  $M^k$ . We know that exists a ball  $B(Q^k, r^k)$  and a straight line  $\Pi^k$  swept by the coordinate  $x_k$  of an orthogonal coordinate system  $(O, x_k, y_k)$  with a suitably chosen origin  $O$ , such that a function  $f_k(x_k)$  uniquely expresses the coordinate  $y_k$  of points located on  $\partial\Omega$ , as long as they lie in  $B(Q^k, r^k)$  [17] (see Figure A2).

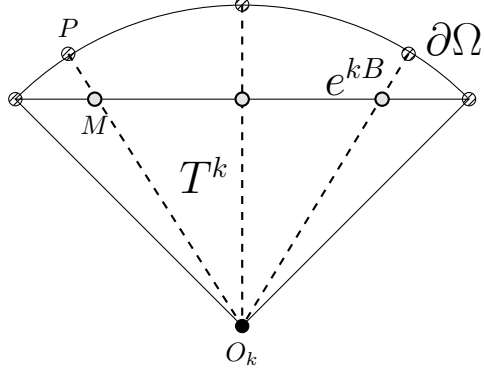
**Assumption A.2.** Consider  $h$  small enough such that  $\Pi^k$  is aligned with  $e^{kB}$  and the ball  $B(Q^k, r^k)$  contains  $e^{kB}$ ,  $\forall T^k, k \in I^B$ .

**Proposition A.1** ([28], Proposition 2.1). If Assumption A.1 and Assumption A.2 hold there exists a constant  $C_{\partial\Omega}$  depending only on  $\partial\Omega$  such that  $\forall M \in e^{kB}$  the length of the segment joining  $M$  and  $P \in \partial\Omega$  aligned with  $O_k$  and  $M$  is bounded above by  $C_{\partial\Omega} h_k^2$  (see Figure A3).

**Proposition A.2** ([28], Proposition 2.2). Assume that  $\partial\Omega$  is of the piecewise  $C^{N+1}$ -class for  $N > 1$ . Let  $v^{(j)}$  denote the derivative of order  $j$  with respect to  $x$  of a sufficiently differentiable function  $v(x)$ ,  $0 \leq j \leq N + 1$ , with  $v^{(0)} = v$ ,  $v^{(1)} = v'$ ,  $v^{(2)} = v''$ . If assumption A.2 holds, there exists a constant  $C_{\partial\Omega}^j$  depending only of  $\partial\Omega$



**Fig. A2:** Each point  $P \in \partial\Omega \cap B(Q^k, r^k)$  has coordinates  $(x_k, f(x_k))$  in the cartesian coordinate system  $(O, x_k, y_k)$ .



**Fig. A3:** Intersection of  $\partial\Omega$  with the straight line joining  $M$  and  $O_k$ .

such that

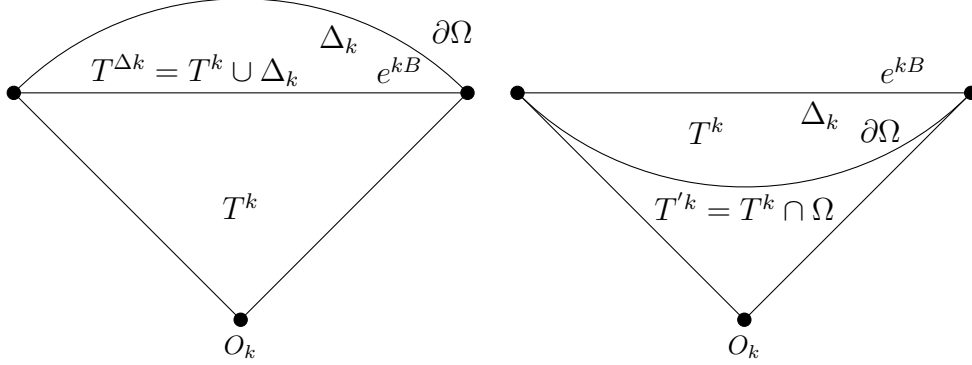
$$\left| f_k^{(j)}(M) \right| \leq C_{\partial\Omega}^j h_k^{\max\{2-j, 0\}}, \quad \forall M \in e^{kB} \quad \text{for } j = 0, 1, \dots, N+1. \quad (\text{A1})$$

**Lemma A.1** ([28], Lemma 3.1). *Consider  $h$  small enough such that Assumption A.1 and Assumption A.2 hold. Then, there exists two mesh-independent constants  $C_\infty$  and  $C_J$  depending only on  $\partial\Omega$  and the shape regularity of  $\mathcal{T}_h$  such that  $\forall w \in \mathcal{P}_N(T^k \cup \Delta_k)$  and  $\forall T^k, k \in I^B$  (see Figure A4), it holds*

$$\|w\|_{L^\infty(T^k \cup \Delta_k)} \leq C_\infty \|w\|_{L^\infty(T^k \cap \Omega)} \quad (\text{A2})$$



$$\|w\|_{L^\infty(T^k \cup \Delta_k)} \leq C J h_k^{-1} \|w\|_{L^2(T^k \cap \Omega)}. \quad (\text{A3})$$



**Fig. A4:** Example for the convex case, where  $T^k \subset \Omega$ , (left panel) and for the concave case, where  $T^k \not\subset \Omega$  (right panel).

Let  $D^j w$  be the  $j$ -th order tensor whose components are the  $j$ -th order partial derivatives with respect to the space of variables of a function  $w$ . In the following, we introduce some technical lemmas that are useful in proving the error estimates.

**Lemma A.2** ([28], Lemma 4.1). *Let  $m$  be an integer,  $m > 1$ , and  $w \in H^m(\Omega)$  such that  $w|_{\partial\Omega} = 0$ , for  $j = 0, 1, \dots, m$ . Let  $I_h(w)$  be the  $\mathcal{P}_N$ -interpolate of  $w$ . Then, for  $p \in [1, \infty]$ , there exists a mesh-independent constant  $C_\Omega$  such that:*

$$\left\| D^j (w - I_h(w)) \right\|_{L^p(\Omega)} \leq C_\Omega h^{m-j} |w|_{W_p^m(\Omega)}. \quad (\text{A4})$$

**Lemma A.3** ([28], Lemma 4.2). *Let  $r = 1/2 + \epsilon$  for a certain  $\epsilon \in (0, 1/2)$  and  $w \in H^{N+1+r}(\Omega \cup \Omega_h)$  be such that  $w|_{\partial\Omega} = 0$ . Let  $\tilde{T}$  be a closed set fulfilling  $(T^k \cap \Omega) \subseteq \tilde{T} \subseteq (T^k \cup \Delta_k)$  and  $w_h$  be a function in  $\mathcal{W}_h$  extended to  $\Delta_k$ ,  $k \in I^B$ . Then there exist constants  $C_j$  independent of  $T^k$  and  $h_k$  such that for  $j = 1, 2, \dots, N$  it holds*

$$\left\| D^j (w - w_h) \right\|_{L^\infty(\tilde{T})} \leq C_j h_k^{-j} \left( \left\| \nabla (w - w_h) \right\|_{L^2(T^k \cap \Omega)} + h_k^N |w|_{H^{N+1}(T^k \cap \Omega)} + h_k^{N+r} \|w\|_{H^{N+1+r}(T^k \cup \Delta_k)} \right). \quad (\text{A5})$$

Now, we deduce a similar result to (A5) for  $j = 0$ .

**Lemma A.4.** *Let  $r = 1/2 + \epsilon$  for a certain  $\epsilon \in (0, 1/2)$  and  $w \in H^{N+1+r}(\Omega \cup \Omega_h)$  be such that  $w|_{\partial\Omega} = 0$ . Let  $\tilde{T}$  be a closed set fulfilling  $(T^k \cap \Omega) \subseteq \tilde{T} \subseteq (T^k \cup \Delta_k)$  and  $w_h$  be a function in  $\mathcal{W}_h$  extended to  $\Delta_k$ ,  $k \in I^B$ . Then there exist a constant  $C_0$*

independent of  $T^k$  and  $h_k$  such that

$$\begin{aligned} \|w - w_h\|_{L^\infty(\tilde{T})} &\leq \frac{C_0}{h_k} \left( \|w - w_h\|_{L^2(T^k \cap \Omega)} + h_k^{N+1} |w|_{H^{N+1}(T^k \cap \Omega)} \right. \\ &\quad \left. + h_k^{N+r+1} \|w\|_{H^{N+1+r}(T^k \cup \Delta_k)} \right). \end{aligned} \quad (\text{A6})$$

*Proof.* Recall that  $I_h(w)$  interpolate  $w$  in  $\Omega \cup \Omega_h$  and  $w - w_h = (w - I_h(w)) + (I_h(w) - w_h)$ . Now, for all  $T^k$ , with  $k \in I^B$

$$\begin{aligned} \|w - w_h\|_{L^\infty(\tilde{T})} &\leq \|w - I_h(w)\|_{L^\infty(T^k \cup \Delta_k)} + \|I_h(w) - w_h\|_{L^\infty(T^k \cup \Delta_k)} \\ &\leq \|w - I_h(w)\|_{L^\infty(T^k \cup \Delta_k)} + \frac{C_J}{h_k} \|I_h(w) - w_h\|_{L^2(T^k \cap \Omega)}, \end{aligned} \quad (\text{A7})$$

using the inequality (A3). Consider the mapping  $\Theta_k$  from  $T^k \cup \Delta_k$  to a unit element  $\widehat{T^k \cup \Delta_k}$ , where  $\Theta_k(x, y) = (x, y)/h_k$ . Setting  $\Theta_k(T^k \cup \Delta_k) = \widehat{T^k \cup \Delta_k}$ , we note that  $H^{N+1+r}(\widehat{T^k \cup \Delta_k})$  is continuously embedded in  $W_\infty^N(\widehat{T^k \cup \Delta_k})$ , i.e., there exists a constant  $C_e$  depending only on  $\widehat{T^k \cup \Delta_k}$  such that [1]

$$\|\hat{v}\|_{W_\infty^N(\widehat{T^k \cup \Delta_k})} \leq C_e \|\hat{v}\|_{H^{N+1+r}(\widehat{T^k \cup \Delta_k})}, \quad \forall \hat{v} \in H^{N+1+r}(\widehat{T^k \cup \Delta_k}). \quad (\text{A8})$$

On the other hand, consider  $\hat{w}$  and  $\hat{I}(\hat{w})$  the transformations under  $\Theta_k$  in  $\widehat{T^k \cup \Delta_k}$  of  $w$  and  $I_h(w)$ , respectively. Notice that  $\hat{I}(\hat{w})$  is a  $\mathcal{P}_N$ -interpolate of  $\hat{w}$  in  $\widehat{T^k \cup \Delta_k}$ . Then, there exists a constant  $\hat{C}_{\mathcal{T}_h}$  depending on  $\widehat{T^k \cup \Delta_k}$  such that

$$\|w - I_h(w)\|_{L^\infty(T^k \cup \Delta_k)} = \|\hat{w} - \hat{I}(\hat{w})\|_{L^\infty(\widehat{T^k \cup \Delta_k})} \leq \hat{C}_{\mathcal{T}_h} C_e \|\hat{w}\|_{H^{N+1+r}(\widehat{T^k \cup \Delta_k})},$$

using (A8). Thus, applying standard transformations to functions in fractional Sobolev spaces, we may write for suitable mesh-independent constants  $C_0^\Delta$  [29]

$$\|w - I_h(w)\|_{L^\infty(T^k \cup \Delta_k)} \leq C_0^\Delta h_k^{N+r} \|w\|_{H^{N+1+r}(T^k \cup \Delta_k)}. \quad (\text{A9})$$

Note that

$$\|I_h(w) - w_h\|_{L^2(T^k \cap \Omega)} \leq \|I_h(w) - w\|_{L^2(T^k \cap \Omega)} + \|w - w_h\|_{L^2(T^k \cap \Omega)}. \quad (\text{A10})$$

Now, following the proof of Theorem 4.4.4 in [7], we get

$$\|I_h(w) - w\|_{L^2(T^k \cap \Omega)} \leq C_L h_k^{N+1} |w|_{H^{N+1}(T^k \cap \Omega)}, \quad (\text{A11})$$

with  $C_L$  a mesh-independent constant.

Finally, combining (A7), (A9), (A10) and (A11), we get

$$\begin{aligned} \|w - w_h\|_{L^\infty(\tilde{T})} &\leq C_0^\Delta h_k^{N+r} \|w\|_{H^{N+1+r}(T^k \cup \Delta_k)} + C_J h_k^{-1} \left( C_L h_k^{N+1} |w|_{H^{N+1}(T^k \cap \Omega)} \right. \\ &\quad \left. + \|w - w_h\|_{L^2(T^k \cap \Omega)} \right) \end{aligned}$$

and (A6) follows with  $C_0 = \max\{C_J, C_J C_L, C_0^\Delta\}$ .  $\square$

## Appendix B Upper bounds estimates

In what follows we derive upper bounds for  $b_{ih}$ , with  $i = 1, 2, 3, 4, 5$ , given by (81), (84), (85), (86), (89), respectively, used in Theorem 4.2 and for  $b_{6h}$  (119), used in Theorem 4.6.

### Estimate for $b_{1h}(u - u_h, z)$ defined by (81)

Using the Cauchy-Schwarz inequality and applying the trace theorem, there exists a constant  $C_t$  depending only on  $\Omega$  ([1], Theorem 1) such that

$$b_{1h}(u - u_h, z) \leq \|u - u_h\|_{L^2(\partial\Omega)} \|\nabla z\|_{L^2(\partial\Omega)} \leq C_t \|u - u_h\|_{L^2(\partial\Omega)} \|z\|_{H^2(\Omega)}.$$

In order to estimate  $\|u - u_h\|_{L^2(\partial\Omega)}$ , consider for each element  $T^k$ ,  $k \in I^B$ , a local orthogonal frame  $(O; x, y)$  whose origin  $O$  is a vertex of  $T^k$  in  $\partial\Omega$ ,  $x$  is the abscissa along the edge  $e^{kB}$  and  $y$  increases from  $e^{kB}$  to  $\partial\Omega$ . Let  $s$  be the curvilinear abscissa along  $(T^k \cup \Delta_k) \cap \partial\Omega$  with origin at  $O$ . Note that  $s$  can be uniquely expressed in terms of  $x$ , for  $x \in [0, l_{kB}]$ , where  $l_{kB}$  is the length of  $e^{kB}$ . Considering  $f_k(x)$  the  $y$ -abscissa of the points in  $(T^k \cup \Delta_k) \cap \partial\Omega$ , let  $\tilde{u}_h$  be the function of  $x$  defined by  $\tilde{u}_h(x) = [u - u_h](x, f_k(x)) = [u - u_h](s(x))$ . Since  $ds = \sqrt{1 + (f_k')^2} dx$ , we have  $\text{length}((T^k \cup \Delta_k) \cap \partial\Omega) \leq C_q l_{kB}$ , with  $C_q = \sqrt{1 + (h_0 \|f_k''\|_{L^\infty(0, l_{kB})})^2}$ .<sup>1</sup> Thus,

$$\begin{aligned} \|u - u_h\|_{L^2(\partial\Omega)}^2 &= \sum_{k \in I^B} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} |(u - u_h)(s(x))|^2 ds \\ &= \sum_{k \in I^B} \int_0^{l_{kB}} |\tilde{u}_h(x)|^2 \sqrt{1 + (f_k')^2} dx \leq C_q \sum_{k \in I^B} \int_0^{l_{kB}} |\tilde{u}_h(x)|^2 dx. \end{aligned} \tag{B12}$$

Since  $\tilde{u}_h$  vanishes at  $N + 1$  different points in  $[0, l_{kB}]$  and  $u_h|_{T^k} \in \mathcal{P}_N(T^k)$ , from standard results for one-dimensional interpolation [25] there exists a mesh independent

---

<sup>1</sup>Consider  $h_0$  as defined in Theorem 3.2.

constant  $C_L$  such that

$$\|\tilde{u}_h\|_{L^2(0, l_{kB})} \leq C_L h_k^{N+1} \left( \int_0^{l_{kB}} |\tilde{u}_h^{(N+1)}(x)|^2 dx \right)^{1/2}. \quad (\text{B13})$$

On the other hand, using Proposition A.2, there is mesh-independent constants  $c_{j, \partial\Omega}$  such that

$$\max_{x \in [0, l_{kB}]} |f_k^{(j)}(x)| \leq c_{j, \partial\Omega} h_k^{2-j}, \quad \text{for } j = 1, \dots, N+1, \forall T^k, k \in I^{kB}. \quad (\text{B14})$$

Thus, taking into account that the derivatives of  $u_h$  of order greater than  $N$  vanish in  $T^k \cup \Delta_k$ , using the chain rule yields for suitable mesh-independent constants  $c_j$ ,  $j = 0, 1, \dots, N$ :

$$|\tilde{u}_h^{(N+1)}| \leq c_0 |D^{N+1}(u)| + \sum_{j=1}^N c_j h_k^{1-j} |D^{N+1-j}(u - u_h)|. \quad (\text{B15})$$

Then, combining (B12), (B13) and (B15), and using the Cauchy-Schwarz inequality, for a suitable mesh-independent constant  $C_{N,0}$ , we get

$$\begin{aligned} \|u - u_h\|_{L^2(\partial\Omega)}^2 &\leq C_q C_L^2 \sum_{k \in I^B} h_k^{2(N+1)} \int_0^{l_{kB}} \left( c_0 |D^{N+1}(u)| + \sum_{j=1}^N c_j h_k^{1-j} |D^{N+1-j}(u - u_h)| \right)^2 dx \\ &\leq C_{N,0} \left( h^{2(N+1)} \int_{\partial\Omega} |D^{N+1}(u)|^2 ds \right. \\ &\quad \left. + \sum_{k \in I^B} h_k^{2(N+1)} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} \sum_{j=1}^N h_k^{2(1-j)} |D^{N+1-j}(u - u_h)|^2 ds \right). \quad (\text{B16}) \end{aligned}$$

From the trace theorem ([1], Theorem 1), we know that there exists a constant  $C_r$  such that

$$\|D^{N+1}(u)\|_{L^2(\partial\Omega)}^2 = \int_{\partial\Omega} |D^{N+1}(u)|^2 ds \leq C_r^2 \|u\|_{H^{N+1+r}(\Omega)}^2. \quad (\text{B17})$$

On the other hand, using the curved triangle  $T^k \cup \Delta_k$  and considering  $i = N+1-j$

$$\begin{aligned} &\int_{(T^k \cup \Delta_k) \cap \partial\Omega} \sum_{j=1}^N h_k^{2(1-j)} |D^{N+1-j}(u - u_h)|^2 ds \\ &\leq C_q h_k \sum_{j=1}^N h_k^{2(1-j)} \left\| D^{N+1-j}(u - u_h) \right\|_{L^\infty(T^k \cup \Delta_k)}^2 \end{aligned}$$

$$= C_q h_k \sum_{i=1}^N h_k^{2(i-N)} \left\| D^i(u - u_h) \right\|_{L^\infty(T^k \cup \Delta_k)}^2. \quad (\text{B18})$$

Now, using Lemmas A.3 with  $w = u$  and  $w_h = u_h$ , we get

$$\begin{aligned} \left\| D^i(u - u_h) \right\|_{L^\infty(T^k \cup \Delta_k)} &\leq C_i h_k^{-i} \left( \left\| \nabla(u - u_h) \right\|_{L^2(T^k \cap \Omega)} + h_k^N |u|_{H^{N+1}(T^k \cap \Omega)} \right. \\ &\quad \left. + h_k^{N+r} \|u\|_{H^{N+r+1}(T^k \cup \Delta_k)} \right). \end{aligned} \quad (\text{B19})$$

From (B16), combining (B18) and (B19), and applying the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned} &\sum_{k \in I^B} h_k^{2(N+1)} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} \sum_{j=1}^N h_k^{2(1-j)} \left| D^{N+1-j}(u - u_h) \right|^2 ds \\ &\leq \sum_{k \in I^B} h_k^{2(N+1)} C_q h_k \sum_{i=1}^N h_k^{2(i-N)} \left\| D^i(u - u_h) \right\|_{L^\infty(T^k \cup \Delta_k)}^2 \\ &\leq \sum_{k \in I^B} h_k^{2(N+1)} C_q h_k \sum_{i=1}^N h_k^{2(i-N)} C_i^2 h_k^{-2i} \left( \left\| \nabla(u - u_h) \right\|_{L^2(T^k \cap \Omega)} + h_k^N |u|_{H^{N+1}(T^k \cap \Omega)} \right. \\ &\quad \left. + h_k^{N+r} \|u\|_{H^{N+r+1}(T^k \cup \Delta_k)} \right)^2 \\ &\leq \sum_{k \in I^B} h_k^3 3C_q \sum_{i=1}^N C_i^2 \left( \left\| \nabla(u - u_h) \right\|_{L^2(T^k \cap \Omega)}^2 + h_k^{2N} |u|_{H^{N+1}(T^k \cap \Omega)}^2 \right. \\ &\quad \left. + h_k^{2N+2r} \|u\|_{H^{N+r+1}(T^k \cup \Delta_k)}^2 \right). \end{aligned}$$

Recalling that  $r = 1/2 + \epsilon$ ,  $h < 1$  and taking into account Theorem 4.1, the definition of the semi norm  $|\cdot|_{H^{N+1}(\Omega)}$  and the definition of the norm  $\|\cdot\|_{H^{N+1+r}(\Omega)}$ , we infer that for a suitable mesh independent constant  $C_{N,1}$  it holds

$$\begin{aligned} &\sum_{k \in I^B} h_k^{2(N+1)} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} \sum_{j=1}^N h_k^{2(1-j)} \left| D^{N+1-j}(u - u_h) \right|^2 ds \\ &\leq h^3 3C_q \sum_{i=1}^N C_i^2 \left( C^2 h^{2N} |u|_{H^{N+1}(\Omega)}^2 + h_k^{2N} |u|_{H^{N+1}(\Omega)}^2 + h_k^{2N+2r} \|u\|_{H^{N+r+1}(\Omega)}^2 \right) \\ &\leq C_{N,1} h^{2(N+1)} \|u\|_{H^{N+r+1}(\Omega)}^2. \end{aligned} \quad (\text{B20})$$

Finally, combining (B17) and (B20) in (B16), we get

$$\|u - u_h\|_{L^2(\partial\Omega)}^2 \leq C_{N,0} \left( h^{2(N+1)} C_r^2 \|u\|_{H^{N+1+r}(\Omega)}^2 + C_{N,1} h^{2(N+1)} \|u\|_{H^{N+r+1}(\Omega)}^2 \right)$$

$$\leq C_{N,2}^2 h^{2(N+1)} \|u\|_{H^{N+1+r}(\Omega)}^2,$$

where  $C_{N,2}^2 = C_{N,0}(C_r^2 + C_{N,1})$ . Thus

$$\|u - u_h\|_{L^2(\partial\Omega)} \leq C_{N,2} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)}. \quad (\text{B21})$$

Then,

$$b_{1h}(u - u_h, z) \leq C_{b1} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)},$$

where  $C_{b1} = C_t C_{N,2}$ .

### Estimate for $b_{2h}(u - u_h, \Pi_h(z))$ defined by (84)

Since  $\Pi_h(z)$  is piecewise linear,  $\nabla_h \Pi_h(z)$  is constant in  $T^k \cup \Delta_k$ . Recalling that  $\Pi_h(z) = 0$  on  $\partial\Omega_h$ , by the Mean Value Theorem and Proposition A.1, we get for  $P_1 \in \partial\Omega_h$ ,

$$\begin{aligned} |\Pi_h(z)(P)| &\leq \text{length}(\overline{PP_1}) \left| \nabla(\Pi_h(z)|_{T^k}) \right| \\ &\leq C_{\partial\Omega} h_k^2 \left| \nabla(\Pi_h(z)|_{T^k}) \right|, \forall P \in \Delta_k, T^k, k \in I^B. \end{aligned} \quad (\text{B22})$$

Using the Cauchy-Schwarz inequality, the inequality above and noticing that  $\text{area}(\Delta_k) \leq C_{\partial\Omega} h_k^3$ , we may write

$$\begin{aligned} b_{2h}(u - u_h, \Pi_h(z)) &= \sum_{k \in I^B} \int_{\Delta_k} (-\Delta(u - u_h) + c(u - u_h)) \Pi_h(z) \, d\mathbf{x} \\ &\leq \left(1 + \|c\|_{L^\infty(\Omega \setminus \Omega_h)}\right) \sum_{k \in I^B} \left( \|\Delta(u - u_h)\|_{L^2(\Delta_k)} + \|u - u_h\|_{L^2(\Delta_k)} \right) \\ &\quad \times \|\Pi_h(z)\|_{L^2(\Delta_k)} \\ &\leq C_c \sum_{k \in I^B} C_{\partial\Omega}^{1/2} h_k^{3/2} \left( \|\Delta(u - u_h)\|_{L^\infty(\Delta_k)} + \|u - u_h\|_{L^\infty(\Delta_k)} \right) \\ &\quad \times C_{\partial\Omega}^{1/2} h_k^{3/2} \|\Pi_h(z)\|_{L^\infty(\Delta_k)} \\ &\leq C_c \sum_{k \in I^B} C_{\partial\Omega}^{1/2} h_k^{3/2} \left( \|\Delta(u - u_h)\|_{L^\infty(T^k \cup \Delta_k)} + \|u - u_h\|_{L^\infty(T^k \cup \Delta_k)} \right) \\ &\quad \times C_{\partial\Omega}^{3/2} h_k^{7/2} \|\nabla \Pi_h(z)\|_{L^\infty(T^k \cup \Delta_k)}, \end{aligned}$$

where  $C_c = 1 + \|c\|_{L^\infty(\Omega \setminus \Omega_h)}$ . From Lemma A.1, we obtain

$$\begin{aligned} b_{2h}(u - u_h, \Pi_h(z)) &\leq C_c C_{\partial\Omega}^2 C_J \sum_{k \in I^B} h_k^4 \left( \|\Delta(u - u_h)\|_{L^\infty(T^k \cup \Delta_k)} + \|u - u_h\|_{L^\infty(T^k \cup \Delta_k)} \right) \\ &\quad \times \|\nabla \Pi_h(z)\|_{L^2(T^k \cap \Delta_k)}. \end{aligned}$$

Now, considering Lemma A.3 with  $j = 2$  and Lemma A.4, applying the Cauchy-Schwarz inequality and recalling that  $h_k \leq h$ , we get

$$\begin{aligned}
b_{2h}(u - u_h, \Pi_h(z)) &\leq C_c C_{\partial\Omega}^2 C_J \sum_{k \in I^B} \left( h_k^2 C_2 \left( \|\nabla(u - u_h)\|_{L^2(T^k \cap \Omega)} + h_k^N |u|_{H^{N+1}(T^k \cap \Omega)} \right. \right. \\
&\quad \left. \left. + h_k^{N+r} \|u\|_{H^{N+1+r}(T^k \cup \Delta_k)} \right) + h_k^3 C_0 \left( \|u - u_h\|_{L^2(T^k \cap \Omega)} \right. \right. \\
&\quad \left. \left. + h_k^{N+1} |u|_{H^{N+1}(T^k \cap \Omega)} + h_k^{N+1+r} \|u\|_{H^{N+1+r}(T^k \cup \Delta_k)} \right) \right) \|\nabla \Pi_h(z)\|_{L^2(T^k)} \\
&\leq C_c C_{\partial\Omega}^2 C_J (C_0 + C_2) h^2 \left( \sqrt{2} \|u - u_h\|_{H^1(\mathcal{T}_h)} + 2h^N |u|_{H^{N+1}(\Omega_h)} \right. \\
&\quad \left. + 2h^{N+r} \|u\|_{H^{N+1+r}(\Omega)} \right) \|\nabla \Pi_h(z)\|_{L^2(\Omega_h)}.
\end{aligned}$$

Applying Theorem 4.1 and since  $h < 1$ , we may write

$$\begin{aligned}
b_{2h}(u - u_h, \Pi_h(z)) &\leq C_c C_{\partial\Omega}^2 C_J (C_0 + C_2) (\sqrt{2}\mathcal{C} + 4) h^{N+2} \|u\|_{H^{N+1+r}(\Omega)} \|\nabla \Pi_h(z)\|_{L^2(\Omega_h)} \\
&= \tilde{C}_{21} h^{N+2} \|u\|_{H^{N+1+r}(\Omega)} \|\nabla \Pi_h(z)\|_{L^2(\Omega_h)},
\end{aligned}$$

where  $\tilde{C}_{21} = C_c C_{\partial\Omega}^2 C_J (C_0 + C_2) (\sqrt{2}\mathcal{C} + 4)$ . On the other hand, from Lemma A.2 we have

$$\|\nabla(z - \Pi_h(z))\|_{L^2(\Omega_h)} \leq \|\nabla(z - \Pi_h(z))\|_{L^2(\Omega)} \leq C_\Omega h |z|_{H^2(\Omega)}.$$

Then, using the Cauchy-Schwarz inequality, we get

$$\begin{aligned}
\|\nabla \Pi_h(z)\|_{L^2(\Omega)}^2 &= \|\nabla(z - z + \Pi_h(z))\|_{L^2(\Omega)}^2 \\
&\leq 2 \|\nabla(z - \Pi_h(z))\|_{L^2(\Omega)}^2 + 2 \|\nabla z\|_{L^2(\Omega)}^2 \\
&\leq 2C_\Omega^2 h^2 |z|_{H^2(\Omega)}^2 + 2 \|z\|_{H^2(\Omega)}^2 \\
&\leq (2 + 2C_\Omega^2 h_0^2) \|z\|_{H^2(\Omega)}^2 = \tilde{C}_\Omega^2 \|z\|_{H^2(\Omega)}^2, \tag{B23}
\end{aligned}$$

with  $\tilde{C}_\Omega = \sqrt{2 + 2C_\Omega^2 h_0^2}$ . Thus, we derive for  $C_{b2} = \tilde{C}_{21} \tilde{C}_\Omega$

$$b_{2h}(u - u_h, \Pi_h(z)) \leq C_{b2} h^{N+2} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}.$$

### Estimate for $b_{3h}(u - u_h, \Pi_h(z))$ defined by (85)

We know that  $\nabla \Pi_h(z)$  is constant on each element. Then,  $\|\nabla \Pi_h(z)\|_{L^\infty(T^k \cup \Delta_k)} = \|\nabla \Pi_h(z)\|_{L^\infty(T^k)}$ . Let  $w|_{T^k} = \left| [\nabla_h(u - u_h)]|_{T^k} \right|$ ,  $\forall T^k, k \in I^B$ . Thus, recalling (B22),

we get

$$\begin{aligned} b_{3h}(u - u_h, \Pi_h(z)) &\leq \sum_{k \in I^B} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} w|_{T^k} |\Pi_h(z)| \, ds \\ &\leq \sum_{k \in I^B} C_{\partial\Omega} h_k^2 \|\nabla \Pi_h(z)\|_{L^\infty(T^k)} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} w|_{T^k} \, ds. \end{aligned}$$

Applying Lemma A.1, we obtain

$$b_{3h}(u - u_h, \Pi_h(z)) \leq \sum_{k \in I^B} C_{\partial\Omega} C_J h_k \|\nabla \Pi_h(z)\|_{L^2(T^k)} \int_{(T^k \cup \Delta_k) \cap \partial\Omega} w|_{T^k} \, ds.$$

Now, consider the master triangle  $\hat{T}$  with vertices  $(0, 0)$ ,  $(0, 1)$ ,  $(1, 0)$  in the reference frame  $(\hat{O}, \hat{x}, \hat{y})$ , where the origin  $\hat{O}$  is one of the vertices of  $T^k$  belonging to  $\partial\Omega$ . Denote the transformation of  $(T^k \cup \Delta_k) \cap \partial\Omega$  under the affine mapping  $\mathcal{F}_k$  from  $T^k$  to  $\hat{T}$  by  $\partial\hat{T}$ . Taking into account that  $\text{length}(T^k \cup \Delta_k) \cap \partial\Omega = \int_0^{l_{kB}} \sqrt{1 + (f'_k)^2} \, d\mathbf{x} \leq C_q h_k$ , where  $C_q = \sqrt{1 + \left(h_0 \|f''_k\|_{L^\infty(0, l_{kB})}\right)^2}$ , and  $\text{length}(\partial\hat{T}) = 1$ , we have

$$b_{3h}(u - u_h, \Pi_h(z)) \leq \sum_{k \in I^B} C_{\partial\Omega} C_J C_q h_k^2 \|\nabla \Pi_h(z)\|_{L^2(T^k)} \int_{\partial\hat{T}} \hat{w} \, d\hat{s},$$

where  $\hat{w}$  is the transformation of  $w|_{T^k}$  under the mapping  $\mathcal{F}_k$ .

We apply the Trace Theorem to the transformation  $\widehat{T^k \cup \Delta_k}$  of  $T^k \cup \Delta_k$  under  $\mathcal{F}_k$ . Since  $\partial\Omega$  is smooth and  $h$  is sufficiently small, there exists a constant  $\hat{C}_t$  independent of  $T^k$  such that

$$\int_{\partial\hat{T}} \hat{w} \, d\hat{s} \leq \hat{C}_t \left( \int_{\widehat{T^k \cup \Delta_k}} \left( \hat{w}^2 + |\widehat{\nabla} \hat{w}|^2 \right) d\hat{s} \right)^{1/2},$$

where  $\widehat{\nabla}$  is the gradient operator for functions defined in  $\widehat{T^k \cup \Delta_k}$ .

Now, moving back to  $T^k \cup \Delta_k$ , we get for a suitable mesh-independent constant  $\tilde{C}_3$

$$\begin{aligned} b_{3h}(u - u_h, \Pi_h(z)) &\leq C_{\partial\Omega} C_J C_q \hat{C}_t \sum_{k \in I^B} h_k^2 \|\nabla \Pi_h(z)\|_{L^2(T^k)} \left( \int_{\widehat{T^k \cup \Delta_k}} \left( \hat{w}^2 + |\widehat{\nabla} \hat{w}|^2 \right) d\hat{s} \right)^{1/2} \\ &\leq \tilde{C}_3 \sum_{k \in I^B} h_k \|\nabla \Pi_h(z)\|_{L^2(T^k)} \left( \int_{T^k \cup \Delta_k} \left( w|_{T^k}^2 + h_k^2 |\nabla w|_{T^k}|^2 \right) ds \right)^{1/2}. \end{aligned}$$



Applying the Cauchy-Schwarz inequality, we have

$$b_{3h}(u - u_h, \Pi_h(z)) \leq \tilde{C}_3 h \|\nabla_h \Pi_h(z)\|_{L^2(\Omega_h)} \left( \sum_{k \in I^B} \|\nabla(u - u_h)\|_{L^2(T^k \cup \Delta_k)}^2 + h_k^2 \|D^2(u - u_h)\|_{L^2(T^k \cup \Delta_k)}^2 \right)^{1/2}.$$

Note that

$$\|D^2(u - u_h)\|_{L^2(T^k \cup \Delta_k)} \leq \sqrt{\text{area}(T^k \cup \Delta_k)} \|D^2(u - u_h)\|_{L^\infty(T^k \cup \Delta_k)}.$$

Now, observe that  $\text{area}(T^k \cup \Delta_k) \leq \text{area}(T^k) + C_{\partial\Omega} h_k^3 \leq h_k^2(1/2 + C_{\partial\Omega} h_0)$ . Applying Lemma A.3, the Cauchy-Schwarz inequality and Theorem 4.1, we get

$$\begin{aligned} & \sum_{k \in I^B} \left( \|\nabla(u - u_h)\|_{L^2(T^k \cup \Delta_k)}^2 + h_k^2 \|D^2(u - u_h)\|_{L^2(T^k \cup \Delta_k)}^2 \right) \\ & \leq \sum_{k \in I^B} \left( h_k^2(1/2 + C_{\partial\Omega} h_0) \|\nabla(u - u_h)\|_{L^\infty(T^k \cup \Delta_k)}^2 + h_k^2 \left( h_k^2(1/2 + C_{\partial\Omega} h_0) \right) \right. \\ & \quad \left. \times \|D^2(u_h - u)\|_{L^\infty(T^k \cup \Delta_k)}^2 \right) \\ & \leq \sum_{k \in I^B} 3(C_1^2 + C_2^2)(1/2 + C_{\partial\Omega} h_0) \left( \|\nabla(u - u_h)\|_{L^2(T^k \cap \Omega)}^2 + h_k^{2N} |u|_{H^{N+1}(T^k \cap \Omega)}^2 \right. \\ & \quad \left. + h_k^{2N+2r} \|u\|_{H^{N+r+1}(T^k \cup \Delta_k)}^2 \right) \leq \bar{C}_3^2 h^{2N} \|u\|_{H^{N+r+1}(\Omega)}^2, \end{aligned}$$

with  $\bar{C}_3^2 = 3(C_1^2 + C_2^2)(1/2 + C_{\partial\Omega} h_0)(C^2 + 2)$ . Thus, using (B23)

$$b_{3h}(u - u_h, \Pi_h(z)) \leq C_{b3} h^{N+1} \|u\|_{H^{N+r+1}(\Omega)} \|z\|_{H^2(\Omega)},$$

with  $C_{b3} = \tilde{C}_3 \tilde{C}_\Omega \bar{C}_3$ .

### Estimate for $b_{4h}(u - u_h, e_h(z))$ defined by (86)

Using the Cauchy-Schwarz inequality and attending the fact that  $\text{area}(\Delta_k) \leq C_\Omega h_k^3$ , we obtain

$$\begin{aligned} b_{4h}(u - u_h, e_h(z)) &= \widehat{a}_{\Delta h}(u - u_h, z - \Pi_h(z)) \\ &= \sum_{k \in I^B} \int_{\Delta_k} \nabla(u - u_h) \cdot \nabla(z - \Pi_h(z)) + c(u - u_h)(z - \Pi_h(z)) \, d\mathbf{x} \\ &\leq C_c \sum_{k \in I^B} C_{\partial\Omega}^{1/2} h_k^{3/2} \|\nabla(u - u_h)\|_{L^\infty(T^k \cup \Delta_k)} \|\nabla(z - \Pi_h(z))\|_{L^2(T^k \cup \Delta_k)} \end{aligned}$$

$$+ C_{\partial\Omega}^{1/2} h_k^{3/2} \|u - u_h\|_{L^\infty(T^k \cup \Delta_k)} \|z - \Pi_h(z)\|_{L^2(T^k \cup \Delta_k)},$$

with  $C_c = 1 + \|c\|_{L^\infty(\Omega \setminus \Omega_h)}$ . Now, applying Lemma A.3 with  $j = 1$  and Lemma A.4, the Cauchy-Schwarz inequality and using Theorem 4.1, we arrive at

$$\begin{aligned} b_{4h}(u - u_h, e_h(z)) &\leq C_c C_{\partial\Omega}^{1/2} (C_1 + C_0) h^{1/2} \left( \|z - \Pi_h(z)\|_{L^2(\Omega)} + \|\nabla(z - \Pi_h(z))\|_{L^2(\Omega)} \right) \\ &\quad \times \left( \|u - u_h\|_{L^2(\Omega_h)} + \|\nabla_h(u - u_h)\|_{L^2(\Omega_h)} + 2h^N |u|_{H^{N+1}(\Omega_h)} \right. \\ &\quad \left. + 2h^{N+r} \|u\|_{H^{N+r+1}(\Omega)} \right) \\ &\leq C_c C_{\partial\Omega}^{1/2} (C_1 + C_0) h^{1/2} \left( \|z - \Pi_h(z)\|_{L^2(\Omega)} + \|\nabla(z - \Pi_h(z))\|_{L^2(\Omega)} \right) \\ &\quad \times \left( (2 + \sqrt{2}\mathcal{C}) h^N |u|_{H^{N+1}(\Omega)} + 2h^{N+r} \|u\|_{H^{N+r+1}(\Omega)} \right) \\ &\leq C_c C_{\partial\Omega}^{1/2} (C_1 + C_0) (4 + \sqrt{2}\mathcal{C}) h^{N+1/2} \left( \|z - \Pi_h(z)\|_{L^2(\Omega)} \right. \\ &\quad \left. + \|\nabla(z - \Pi_h(z))\|_{L^2(\Omega)} \right) \|u\|_{H^{N+r+1}(\Omega)}. \end{aligned}$$

Considering Lemma A.2 with  $j = 0, 1$  we get

$$b_{4h}(u - u_h, e_h(z)) \leq C_{b4} h^{N+3/2} \|u\|_{H^{N+r+1}(\Omega)} \|z\|_{H^2(\Omega)},$$

where  $C_{b4} = 2C_c C_{\partial\Omega}^{1/2} (C_1 + C_0) (4 + \sqrt{2}\mathcal{C}) C_\Omega$ .

### Estimate for $b_{5h}(u - u_h, z)$ defined by (89)

Attending to the definition of jump and average along a boundary edge, we may write

$$\begin{aligned} b_{5h}(u - u_h, z) &= - \int_{\partial\Omega_h} (u - u_h) \frac{\partial z}{\partial n} \, ds \\ &= - \sum_{e \in \partial\Omega_h} \int_e \llbracket u - u_h \rrbracket \cdot \nabla(z - \Pi_h(z) + \Pi_h(z)) \, ds \\ &\leq \left| \sum_{e \in \partial\Omega_h} \int_e \llbracket u - u_h \rrbracket \cdot \nabla(z - \Pi_h(z) + \Pi_h(z)) \, ds \right|. \end{aligned}$$

Thus, applying the Cauchy-Schwarz inequality we get

$$\begin{aligned} b_{5h}(u - u_h, z) &\leq \sum_{e \in \partial\Omega_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla(z - \Pi_h(z) + \Pi_h(z)) \right\|_{L^2(e)} \\ &\leq \sum_{e \in \partial\Omega_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla(z - \Pi_h(z)) \right\|_{L^2(e)} \quad (\text{B24}) \end{aligned}$$

$$+ \sum_{e \in \partial\Omega_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla \Pi_h(z) \right\|_{L^2(e)}. \quad (\text{B25})$$

Now, following similar arguments as in Subsection 3.1, using Theorem 4.1 and Lemma A.2 , we may write (B24) as follows:

$$\begin{aligned}
& \sum_{e \in \partial\Omega_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla(z - \Pi_h(z)) \right\|_{L^2(e)} \\
& \leq \sum_{e \in \mathcal{E}_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla(z - \Pi_h(z)) \right\|_{L^2(e)} \\
& \leq \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \left\| \llbracket u - u_h \rrbracket \right\|_{L^2(e)}^2 \right)^{1/2} \left( \sum_{e \in \mathcal{E}_h} h_e \left\| \frac{\partial e_h(z)}{\partial n} \right\|_{L^2(e)}^2 \right)^{1/2} \\
& \leq |u - u_h|_* \left( \sum_{k=1}^K C_T^2 \left( |e_h(z)|_{H^1(T^k)}^2 + h_k^2 |e_h(z)|_{H^2(T^k)}^2 \right) \right)^{1/2} \\
& \leq \|u - u_h\| C_T \sqrt{2} C_\Omega h |z|_{H^2(\Omega)} \leq C'_{b5} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \quad (\text{B26})
\end{aligned}$$

with  $C'_{b5} = \sqrt{2} C C_T C_\Omega$ .

Considering  $\partial\Delta_k$  the boundary of  $\Delta_k$ , for  $k \in I^B$ , applying Lemma A.1 and inequality (B23), we can rewrite (B25) in the following way

$$\begin{aligned}
& \sum_{e \in \partial\Omega_h} \left\| h_e^{-1/2} \llbracket u - u_h \rrbracket \right\|_{L^2(e)} \left\| h_e^{1/2} \nabla \Pi_h(z) \right\|_{L^2(e)} \\
& \leq \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \left\| \llbracket u - u_h \rrbracket \right\|_{L^2(e)}^2 \right)^{1/2} \left( \sum_{e \in \partial\Omega_h} h_e \left\| \nabla \Pi_h(z) \right\|_{L^2(e)}^2 \right)^{1/2} \\
& \leq \left( \sum_{e \in \mathcal{E}_h} h_e^{-1} \left\| \llbracket u - u_h \rrbracket \right\|_{L^2(e)}^2 \right)^{1/2} \left( \sum_{k \in I^B} h_k \left\| \nabla \Pi_h(z) \right\|_{L^2(\partial\Delta_k)}^2 \right)^{1/2} \\
& \leq |u - u_h|_* \left( \sum_{k \in I^B} h_k C_t^2 \left\| \nabla \Pi_h(z) \right\|_{H^1(\Delta_k)}^2 \right)^{1/2} \\
& = |u - u_h|_* \left( \sum_{k \in I^B} h_k C_t^2 \left\| \nabla \Pi_h(z) \right\|_{L^2(\Delta_k)}^2 \right)^{1/2} \\
& \leq |u - u_h|_* \left( \sum_{k \in I^B} h_k^4 C_t^2 C_{\partial\Omega} \left\| \nabla \Pi_h(z) \right\|_{L^\infty(\Delta_k)}^2 \right)^{1/2} \\
& = |u - u_h|_* \left( \sum_{k \in I^B} h_k^4 C_t^2 C_{\partial\Omega} \left\| \nabla \Pi_h(z) \right\|_{L^\infty(T^k \cup \Delta_k)}^2 \right)^{1/2}
\end{aligned}$$

$$\begin{aligned}
&\leq |u - u_h|_* \left( \sum_{k \in I^B} h_k^2 C_t^2 C_{\partial\Omega} C_J^2 \|\nabla \Pi_h(z)\|_{L^2(T^k \cap \Delta_k)}^2 \right)^{1/2} \\
&\leq \mathcal{C} C_t C_{\partial\Omega}^{1/2} C_J h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|\nabla \Pi_h(z)\|_{L^2(\Omega)} \\
&\leq C_{b5}'' h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)}, \tag{B27}
\end{aligned}$$

with  $C_{b5}'' = \mathcal{C} C_t C_{\partial\Omega}^{1/2} C_J \tilde{C}_\Omega$ .

Thus, combining (B26) and (B27), we get

$$b_{5h}(u - u_h, z) \leq C_{b5} h^{N+1} \|u\|_{H^{N+1+r}(\Omega)} \|z\|_{H^2(\Omega)},$$

where  $C_{b5} = C_{b5}' + C_{b5}''$ .

### Estimate for $b_{6h}(u_h, \Pi_h(z))$ defined by (119)

Following similar arguments as in (B22), we get for  $P_1 \in \partial\Omega_h$ ,

$$|\Pi_h(z)(P)| \leq C_{\partial\Omega} h_k^2 \left| \nabla (\Pi_h(z)|_{T^k}) \right|, \forall P \in \Delta_k, T^k, k \in \mathcal{Q}^B. \tag{B28}$$

Using the Cauchy-Schwarz inequality, the inequality above and noticing that  $\text{area}(\Delta_k) \leq C_{\partial\Omega} h_k^3$ , we may write

$$\begin{aligned}
b_{6h}(u_h, \Pi_h(z)) &= \sum_{k \in \mathcal{Q}^B} \int_{\Delta_k} (-\Delta u_h + c u_h) \Pi_h(z) \, d\mathbf{x} \\
&\leq C_c' \sum_{k \in \mathcal{Q}^B} C_{\partial\Omega}^{1/2} h_k^{3/2} \left( \|\Delta u_h\|_{L^\infty(T^k)} + \|u_h\|_{L^\infty(T^k)} \right) C_{\partial\Omega}^{3/2} h_k^{7/2} \|\nabla \Pi_h(z)\|_{L^\infty(T^k)},
\end{aligned}$$

where  $C_c' = 1 + \|c\|_{L^\infty(\Omega_h \setminus \Omega)}$ . From Lemma A.1, we obtain

$$b_{6h}(u_h, \Pi_h(z)) \leq C_c' C_{\partial\Omega}^2 C_J \sum_{k \in \mathcal{Q}^B} h_k^4 \left( \|\Delta u_h\|_{L^\infty(T^k)} + \|u_h\|_{L^\infty(T^k)} \right) \|\nabla \Pi_h(z)\|_{L^2(T^k \cap \Omega)}.$$

Note that by adding and subtracting the exact solution  $u$ , we get

$$\begin{aligned}
\|\Delta u_h\|_{L^\infty(T^k)} + \|u_h\|_{L^\infty(T^k)} &\leq \|\Delta(u - u_h)\|_{L^\infty(T^k)} + \|u - u_h\|_{L^\infty(T^k)} + \|\Delta u\|_{L^\infty(T^k)} \\
&\quad + \|u\|_{L^\infty(T^k)}.
\end{aligned}$$

Thus, considering the previous inequality, Lemma A.3 with  $j = 2$  and Lemma A.4, applying the Cauchy-Schwarz inequality, and recalling that  $h_k \leq h$  and that  $\tilde{u}$  is the regular extension of  $u$  to  $\tilde{\Omega}$  such that  $\tilde{u}|_\Omega = u$ , we conclude that for a suitable

mesh-independent constant  $C'_5$  it holds

$$\begin{aligned}
b_{6h}(u_h, \Pi_h(z)) &\leq C'_c C_{\partial\Omega}^2 C_J \sum_{k \in \mathcal{Q}^B} h_k^4 \|\nabla \Pi_h(z)\|_{L^2(T^k \cap \Omega)} \left( C \left( \|\Delta u\|_{L^\infty(\Omega)} + \|u\|_{L^\infty(\Omega)} \right) \right. \\
&\quad + \frac{C_2}{h_k^2} \left( \|\nabla(u - u_h)\|_{L^2(T^k \cap \Omega)} + h_k^2 |u|_{H^3(T^k \cap \Omega)} + h_k^{2+r} \|\tilde{u}\|_{H^{3+r}(T^k \cup \Delta_k)} \right) \\
&\quad \left. + \frac{C_0}{h_k} \left( \|u - u_h\|_{L^2(T^k \cap \Omega)} + h_k^3 |u|_{H^3(T^k \cap \Omega)} + h_k^{3+r} \|\tilde{u}\|_{H^{3+r}(T^k \cup \Delta_k)} \right) \right) \\
&\leq C'_5 \|\nabla \Pi_h(z)\|_{L^2(\Omega \cap \Omega_h)} \left( h^{7/2} \|u\|_{W_\infty^2(\Omega)} + h^2 \left( \|u - u_h\|_{L^2(\Omega \cap \Omega_h)} \right. \right. \\
&\quad \left. \left. + \|\nabla(u - u_h)\|_{L^2(\Omega \cap \Omega_h)} \right) + h^4 |u|_{H^3(\Omega \cap \Omega_h)} + h^{4+r} \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} \right). \tag{B29}
\end{aligned}$$

Now, we note that by the Sobolev embedding Theorem, there exists a constant  $C_s$  such that

$$\|u\|_{W_\infty^2(\Omega)} \leq C_s \|u\|_{H^{3+r}(\Omega)} \leq C_s \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})}. \tag{B30}$$

Considering the previous inequality and applying (112), for a suitable mesh-independent constant, we may rewrite (B29) as follows:

$$\begin{aligned}
b_{6h}(u_h, \Pi_h(z)) &\leq \bar{C}_{51} \|\nabla \Pi_h(z)\|_{L^2(\Omega \cap \Omega_h)} \left( h^{7/2} \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + h^4 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} \right. \right. \\
&\quad \left. \left. + h^{1/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) + h^4 |u|_{H^3(\Omega \cap \Omega_h)} + h^{4+r} \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} \right) \\
&\leq \bar{C}_{52} \|\nabla \Pi_h(z)\|_{L^2(\Omega \cap \Omega_h)} \left( h^{7/2} \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + h^4 \left( |\tilde{u}|_{H^3(\tilde{\Omega})} \right. \right. \\
&\quad \left. \left. + h^{1/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \right). \tag{B31}
\end{aligned}$$

On the other hand, considering similar arguments as in inequality (B23), we obtain

$$\|\nabla \Pi_h(z)\|_{L^2(\Omega \cap \Omega_h)}^2 \leq \tilde{C}_\Pi^2 \|z\|_{H^2(\Omega)}^2, \tag{B32}$$

with  $\tilde{C}_\Pi = \sqrt{2 + 2C_\Omega'^2 h_0^2}$ . Thus, we derive for  $\tilde{C}_{b6} = \bar{C}_{52} \tilde{C}_\Pi$

$$\begin{aligned}
b_{6h}(u_h, \Pi_h(z)) &\leq \tilde{C}_{b6} h^{7/2} \left( \|\tilde{u}\|_{H^{3+r}(\tilde{\Omega})} + h^{1/2} \left( |\tilde{u}|_{H^3(\tilde{\Omega})} + h^{1/2} \|-\Delta \tilde{u} + c\tilde{u}\|_{L^2(\tilde{\Omega})} \right) \right) \\
&\quad \times \|z\|_{H^2(\Omega)}. \tag{B33}
\end{aligned}$$

## References

- [1] Adams RA (1975) Sobolev spaces. Academic Press, New York

- [2] Arnold DN (1982) An interior penalty finite element method with discontinuous elements. *SIAM J Numer Anal* 19(4):742–760. <https://doi.org/10.1137/0719052>
- [3] Arnold DN, Brezzi F, Cockburn B, et al (2002) Unified analysis of discontinuous Galerkin methods for elliptic problems. *SIAM J Numer Anal* 39(5):1749–1779. <https://doi.org/10.1137/S0036142901384162>
- [4] Atallah NM, Canuto C, Scovazzi G (2022) The high-order shifted boundary method and its analysis. *Comput Methods Appl Mech Eng* 394:114885. <https://doi.org/10.1016/j.cma.2022.114885>
- [5] Bassi F, Rebay S (1995) Accurate 2D Euler computations by means of a high order discontinuous finite element method. In: Deshpande SM, Desai SS, Narasimha R (eds) *Fourteenth International Conference on Numerical Methods in Fluid Dynamics*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp 234–240. [https://doi.org/10.1007/3-540-59280-6\\_128](https://doi.org/10.1007/3-540-59280-6_128)
- [6] Bernstein DS (2009) *Matrix Mathematics: Theory, Facts, and Formulas* (Second Edition). Princeton University Press, Princeton, <https://doi.org/10.1515/9781400833344>
- [7] Brenner SC, Scott LR (2008) *The Mathematical Theory of Finite Element Methods*. In: *Texts in Applied Mathematics*, vol 15. Springer New York, New York, <https://doi.org/10.1007/978-0-387-75934-0>
- [8] Brezzi F (1974) On the existence, uniqueness and approximation of saddle-point problems arising from Lagrangian multipliers. *RAIRO Analyse Numérique* 8(R2):129–151. <https://doi.org/10.1051/m2an/197408R201291>
- [9] Clain S, Lopes D, Pereira RMS (2021) Very high-order cartesian-grid finite difference method on arbitrary geometries. *J Comput Phys* 434:110217. <https://doi.org/10.1016/j.jcp.2021.110217>
- [10] Costa R, Clain S, Loubère R, et al (2018) Very high-order accurate finite volume scheme on curved boundaries for the two-dimensional steady-state convection–diffusion equation with Dirichlet condition. *Appl Math Model* 54:752–767. <https://doi.org/10.1016/j.apm.2017.10.016>
- [11] Costa R, Nóbrega JM, Clain S, et al (2019) Very high-order accurate polygonal mesh finite volume scheme for conjugate heat transfer problems with curved interfaces and imperfect contacts. *Comput Methods Appl Mech Eng* 357:112560. <https://doi.org/10.1016/j.cma.2019.07.029>
- [12] Costa R, Nóbrega JM, Clain S, et al (2021) Efficient very high-order accurate polyhedral mesh finite volume scheme for 3D conjugate heat transfer problems in curved domains. *J Comput Phys* 445:110604. <https://doi.org/10.1016/j.jcp.2021.110604>

- [13] Costa R, Clain S, Machado GJ, et al (2022) Very high-order accurate finite volume scheme for the steady-state incompressible Navier–Stokes equations with polygonal meshes on arbitrary curved boundaries. *Comput Methods Appl Mech Eng* 396:115064. <https://doi.org/10.1016/j.cma.2022.115064>
- [14] Costa R, Clain S, Machado GJ, et al (2023) Imposing slip conditions on curved boundaries for 3D incompressible flows with a very high-order accurate finite volume scheme on polygonal meshes. *Comput Methods Appl Mech Eng* 415:116274. <https://doi.org/10.1016/j.cma.2023.116274>
- [15] Cuminato JA, Ruas V (2015) Unification of distance inequalities for linear variational problems. *Comp Appl Math* 34(3):1009–1033. <https://doi.org/10.1007/s40314-014-0163-6>
- [16] Duvigneau R (2020) CAD-consistent adaptive refinement using a NURBS-based discontinuous Galerkin method. *Int J Numer Methods Fluids* 92(9):1096–1117. <https://doi.org/10.1002/fld.4819>
- [17] Evans LC (1998) *Partial differential equations*. American Mathematical Society, University of California, Berkeley, CA, <https://doi.org/10.1090/gsm/019>
- [18] Geuzaine C, Remacle JF (2009) Gmsh: a three-dimensional finite element mesh generator with built-in pre- and post-processing facilities. *Int J Numer Meth Engng* 79(11):1309–1331. <https://doi.org/10.1002/nme.2579>
- [19] Ghalati MK (2017) Numerical analysis and simulation of discontinuous Galerkin method for time-domain maxwell’s equations. PhD thesis, Universidade de Coimbra, Departamento de Matemática da Faculdade de Ciências e Tecnologia
- [20] Grisvard P (2011) *Elliptic Problems in Nonsmooth Domains*. Society for Industrial and Applied Mathematics, Philadelphia, <https://doi.org/10.1137/1.9781611972030>
- [21] Hesthaven J, Warburton T (2008) Nodal discontinuous galerkin methods: Algorithms, analysis, an applications. In: *Texts in Applied Mathematics*, vol 54. Springer-Verlag, New York, <https://doi.org/10.1007/978-0-387-72067-8>
- [22] Krivodonova L, Berger M (2006) High-order accurate implementation of solid wall boundary conditions in curved geometries. *J Comput Phys* 211:492–512. <https://doi.org/10.1016/j.jcp.2005.05.029>
- [23] Lew AJ, Negri M (2011) Optimal convergence of a discontinuous-Galerkin-based immersed boundary method. *ESAIM: M2AN*, 45(4):651–674. <https://doi.org/10.1051/m2an/2010069>
- [24] Mu L, Wang J, Wang Y, et al (2014) Interior penalty discontinuous Galerkin method on very general polygonal and polyhedral meshes. *J Comput Appl Math*

255:432–440. <https://doi.org/10.1016/j.cam.2013.06.003>

- [25] Quarteroni A, Sacco R, Saleri F (2017) Numerical mathematics. In: Texts in Applied Mathematics, vol 37. Springer, New York, <https://doi.org/10.1007/b98885>
- [26] Ruas V (2017) Variational formulations yielding high-order finite-element solutions in smooth domains without curved elements. *J Appl Math Phys* 5:2127–2139. <https://doi.org/10.4236/jamp.2017.511174>
- [27] Ruas V (2019) Accuracy enhancement for non-isoparametric finite-element simulations in curved domains; application to fluid flow. *Comput Math Appl* 77(6):1756–1769. <https://doi.org/10.1016/j.camwa.2018.05.029>
- [28] Ruas V (2020) Optimal Lagrange and Hermite finite elements for Dirichlet problems in curved domains with straight-edged triangles. *Z Angew Math Mech* 100. <https://doi.org/10.1002/zamm.201900296>
- [29] Sanchez A, Arcangeli R (1984) Estimations des erreurs de meilleure approximation polynomiale et d'interpolation de Lagrange dans les espaces de Sobolev d'ordre non entier. *Numer Math* 45:301–321. <https://doi.org/10.1007/BF01389473>
- [30] Santos M, Araújo A, Barbeiro S, et al (2024) Very high-order accurate discontinuous Galerkin method for curved boundaries with polygonal meshes. *J Sci Comput* 100(3):66. <https://doi.org/10.1007/s10915-024-02613-2>
- [31] Strang G, Berger AE (1973) The change in solution due to change in domain. In: Spencer DC (ed) Proceedings of symposia in pure mathematics, Partial Differential Equations, vol 23. American Mathematical Society, Providence, Rhode Island, pp 199–205, <https://doi.org/10.1090/pspum/023>
- [32] Yin J, Xu L, Xie P, et al (2020) A curved boundary treatment for discontinuous Galerkin method applied to Euler equations on triangular and tetrahedral grids. *Comput Phys Commun* 258:107549. <https://doi.org/10.1016/j.cpc.2020.107549>
- [33] Zhang X (2016) A curved boundary treatment for discontinuous Galerkin schemes solving time dependent problems. *J Comput Phys* 308:153–170. <https://doi.org/10.1016/j.jcp.2015.12.036>