



**HAL**  
open science

# HDG and HDG+ methods for harmonic wave problems with convection

Hélène Barucq, Nathan Rouxelin, Sébastien Tordeux

► **To cite this version:**

Hélène Barucq, Nathan Rouxelin, Sébastien Tordeux. HDG and HDG+ methods for harmonic wave problems with convection. [Research Report] RR-9410, Inria Bordeaux - Sud-Ouest; LMAP UMR CNRS 5142; Université de Pau et des Pays de l'Adour (UPPA), Pau, FRA. 2021. hal-03253415

**HAL Id: hal-03253415**

**<https://inria.hal.science/hal-03253415>**

Submitted on 8 Jun 2021

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

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



# HDG and HDG+ methods for harmonic wave problems with convection

Hélène Barucq, Nathan Rouxelin, Sébastien Tordeux

**RESEARCH  
REPORT**

**N° 9410**

June 2021

Project-Team Makutu





# HDG and HDG+ methods for harmonic wave problems with convection

Hélène Barucq<sup>\*</sup>, Nathan Rouxelin<sup>\*</sup>, Sébastien  
Tordeux<sup>\*</sup>

Project-Team Makutu

Research Report n° 9410 — June 2021 — 88 pages

**Abstract:** In this report, we introduce three variants of the HDG method based on two weak formulations of the convected Helmholtz equation. Two of them are standard HDG methods with the same interpolation degree for all the unknowns and one of them uses a higher interpolation degree for the volumetric scalar unknown. For those three numerical methods, a detailed analysis including local and global well-posedness, as well as convergence estimates is carried out. We then provide implementation details and numerical experiments to illustrate our theoretical results.

**Key-words:** Hybridizable Discontinuous Galerkin Method (HDG), aeroacoustics, convected Helmholtz equation, harmonic regime

---

<sup>\*</sup> Makutu, Inria, e2s-UPPA. Emails : {helene.barucq, nathan.rouxelin, sebastien.tordeux}@inria.fr

**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

200 avenue de la Vieille Tour  
33405 Talence Cedex

# Méthodes HDG et HDG+ pour des problèmes d'ondes convectées en régime harmonique

**Résumé :** Dans ce rapport, nous construisons trois variantes de la méthode HDG, basées sur deux formulations faibles de l'équation d'Helmholtz convectée. Deux de ces méthodes sont des méthodes HDG standard qui utilisent le même degré d'interpolation polynomiale pour toutes les inconnues. La troisième méthode, quant à elle, utilise un degré d'interpolation plus élevé pour l'inconnue scalaire volumique, à l'instar des méthodes HDG+. Pour toutes ces méthodes, une analyse détaillée a été effectuée, elle inclut des résultats d'existence et unicité locale et globale ainsi qu'une étude de convergence. Pour finir, nous présentons les détails de l'implémentation de ces méthodes et des expériences numériques qui illustrent nos résultats théoriques.

**Mots-clés :** Méthode de Galerkin discontinue hybridisable, aéroacoustique, équation d'Helmholtz convectée, domaine fréquentiel

# Contents

<b>Introduction</b>	<b>4</b>
<b>1 Model problem</b>	<b>5</b>
1.1 First-order formulations	6
<b>2 Notations</b>	<b>8</b>
2.1 Approximation spaces	8
2.2 Hermitian products and norms	9
2.3 Faces, jumps and averages	10
<b>3 HDG method for the total flux formulation</b>	<b>11</b>
3.1 Constructing the formulation	11
3.2 Choice of penalization parameter	14
3.3 Local solvability	21
3.4 Error analysis	24
3.5 Global solvability	26
<b>4 HDG(+) methods for the diffusive flux formulation</b>	<b>27</b>
4.1 Construction of the method	28
4.2 Local solvability	31
4.3 Error analysis of the HDG+ method	34
4.4 Error analysis of the HDG method with diffusive flux	47
<b>5 Implementation</b>	<b>48</b>
5.1 Framework and notations	49
5.2 Implementation of the diffusive flux HDG method	50
5.3 Implementation of the total flux HDG method	53
5.4 Implementation of the HDG+ method	54
5.5 Comparison of the cost of the HDG and HDG+ methods	60
<b>6 Numerical experiments</b>	<b>61</b>
6.1 Convergence rate	61
6.2 A posteriori error estimate	72
6.3 Is the upwinding mechanism necessary ?	76
6.4 Point-sources in a uniform flow	78
6.5 Gaussian jet	81
<b>Conclusion</b>	<b>82</b>
<b>A Intermediate results for the error analysis of the HDG method with diffusive flux</b>	<b>83</b>
<b>References</b>	<b>88</b>

## Introduction

Nowadays the solar interior is studied by considering the propagation of aeroacoustic waves in time-harmonic domain. Realistic models of solar oscillations require to approximate non-standard Hilbert settings leading to non-conforming methods, such as *Discontinuous Galerkin Methods*. As those methods have a very important numerical cost, we consider the so-called *Hybridizable Discontinuous Galerkin Methods* (HDG), which relies on a static condensation process to reduce the number of degrees of freedom.

As a first step towards the use of HDG in helioseismology, we construct and study HDG for the simplest aeroacoustic model : the *convected Helmholtz equation*.

HDG have been used and validated by numerous authors for various problems such as elliptic equations in [CGL09, CDG<sup>+</sup>09, CC12, CC14], acoustic wave propagation in [GM11, GSV18, NPRC15], elastic wave propagation in [HPS17, BDMP21, CS13, FCS15, BCDL15], Maxwell equations in [CQSS17, CQS18, CLOS20]. These methods have also been used to implement the forward propagator in the context of quantitative inverse problems in [FS20] where a specific formulation of the adjoint method is developed. In this paper, we will consider the HDG+ variant of HDG, introduced in [Leh10], where different polynomial degrees are used for the different unknowns. This HDG+ has been considered for various applications in [CQSS17, Oik14, Oik16, Oik18, QSS16, QS16a, QS16b, Hun19] and to the best of our knowledge, the case of the convected Helmholtz equation has not been addressed yet.

Theory for HDGs is rather similar to the one for mixed finite elements and the actual connection was first established by Cockburn and his coworkers in [CGS10]. For a self-contained introduction to the theory of HDG, we refer to [DS19]. For a historical perspective on HDG, we refer to [Coc14].

For a comparison between HDG and Continuous Galerkin methods, we refer to [KSC12, YMKS16]. The relationship between HDG and HHO (Hybrid High-Order, another new generation of high-order face-based finite element method) has been studied in [CDPE16].

**Main results:** We construct three variants of the HDG method for convected acoustics in the frequency domain. Our main results include a detailed analysis of those methods where the most important properties of the method are proved including local and global solvability, convergence rate for regular solutions. The choice of the penalization parameter is also discussed. Finally, those three methods were implemented in **hawen** (see [Fau21]) and we also provide numerical experiments. Using those numerical experiments, we can conclude that the HDG+ and HDG- $\sigma_h$  methods should be preferred to the HDG- $q_h$  method as they seem more robust.

**Organization of this paper:** This work is organized as follows:

- in [SECTION 1](#): we present the convected Helmholtz equation and recall some results on this equation, we also present two ways to reach a first-order in space formulation;
- in [SECTION 2](#): we introduce some notations and the approximation spaces needed to construct HDGs developed in this paper;
- in [SECTION 3](#): we construct the HDG- $\sigma_h$  method based on the *total flux formulation* of the convected Helmholtz equation, we also provide theoretical results and discuss the optimal choice of penalization parameter for this method;
- in [SECTION 4](#): we construct the HDG- $q_h$  and HDG+ methods based on the *diffusive flux formulation* of the convected Helmholtz equation, we also provide detailed analysis of those methods;

- in [SECTION 5](#): we give details on how those methods can be implemented in a nodal settings;
- in [SECTION 6](#): we present numerical experiments to illustrate our theoretical results, as well as some illustrative examples.

## 1 Model problem

As a model problem we consider the so-called *convected Helmholtz equation*

$$\rho_0 \left( -\omega^2 p - 2i\omega \mathbf{v}_0 \cdot \nabla p + \mathbf{v}_0 \cdot \nabla (\mathbf{v}_0 \cdot \nabla p) \right) - \operatorname{div} \left( \rho_0 c_0^2 \nabla p \right) = s \quad (1)$$

where  $\omega$  is the angular frequency,  $\rho_0$  is the density of the fluid,  $\mathbf{v}_0$  is the velocity of the fluid,  $c_0$  is the adiabatic sound speed, and  $s$  is the acoustic source.

**Validity of this equation:** Equation (1) is the simplest aeroacoustic models and therefore has a limited validity. This equation can be used for

- a *uniform background flow*, in this case the unknown  $p$  can be interpreted as a pressure perturbation,
- a *potential background flow*, in this case the unknown  $p$  should be interpreted as an *acoustic potential* and the physical quantities can be retrieved using the following identities

$$\begin{aligned} \text{Pressure perturbation:} & \quad p' = -\rho_0 c_0 (-i\omega + \mathbf{v}_0 \cdot \nabla) p, \\ \text{Velocity perturbation:} & \quad \mathbf{v}' = -c_0 \nabla p, \end{aligned}$$

see [[Pie90](#), Sec. II.].

**Combining the second-order differential operators:** We will assume that the background flow is incompressible which leads to the following local mass conservation equation

$$\operatorname{div}(\rho_0 \mathbf{v}_0) = 0.$$

With this assumption, we have

$$\begin{aligned} \rho_0 \mathbf{v}_0 \cdot \nabla (\mathbf{v}_0 \cdot \nabla p) &= \operatorname{div}(\rho_0 (\mathbf{v}_0 \cdot \nabla p) \mathbf{v}_0) - (\mathbf{v}_0 \cdot \nabla p) \operatorname{div}(\rho_0 \mathbf{v}_0) \\ &= \operatorname{div}(\rho_0 (\mathbf{v}_0 \cdot \nabla p) \mathbf{v}_0) \\ &= \operatorname{div}(\rho_0 \mathbf{v}_0 \mathbf{v}_0^T \nabla p) \end{aligned}$$

Leading to

$$\rho_0 \left( -\omega^2 p - 2i\omega \mathbf{v}_0 \cdot \nabla p \right) - \operatorname{div}(\mathbf{K}_0 \nabla p) = s \quad (2)$$

where  $\mathbf{K}_0 = \rho_0 (c_0^2 \mathbf{Id} - \mathbf{v}_0 \mathbf{v}_0^T)$ .

It is easy to prove that

**Lemma 1.1:**

$\mathbf{K}_0$  is symmetric positive-definite and

$$\operatorname{Sp}(\mathbf{K}_0) = \left\{ \rho_0 c_0^2, \rho_0 (c_0^2 - |\mathbf{v}_0|^2) \right\}$$

**Proof:**  $\mathbf{K}_0 \mathbf{v}_0 = \rho_0 (c_0^2 - |\mathbf{v}_0|^2) \mathbf{v}_0$  and  $\mathbf{K}_0 \mathbf{u} = \rho_0 c_0^2 \mathbf{u}$  for all  $\mathbf{u} \in \mathbf{v}_0^\perp$ .



**Fredholm type :** If the background flow is subsonic, ie.

$$\inf_{\mathcal{O}} (c_0^2 - |\mathbf{v}_0|^2) > 0, \quad (3)$$

then (2) leads to a problem of Fredholm type. Indeed, by using [LEMMA 1.1](#) we can conclude that  $-\operatorname{div}(\mathbf{K}_0 \nabla p)$  is a coercive operator, and that the convected Helmholtz equation therefore has a *coercive + compact* structure.

**Boundary conditions:** Let  $\Gamma$  be the boundary of the domain  $\mathcal{O}$  and let  $\mathbf{n}$  be the outward-facing normal vector.

We will use the following boundary conditions

$$\text{Neumann:} \quad (\mathbf{K}_0 \nabla p) \cdot \mathbf{n} + 2i\omega(\rho_0 \mathbf{v}_0 \cdot \mathbf{n})p = g_N \quad \text{on } \Gamma_N \quad (4a)$$

$$\text{Dirichlet:} \quad p = g_D \quad \text{on } \Gamma_D \quad (4b)$$

$$\text{Impedance:} \quad (\mathbf{K}_0 \nabla p) \cdot \mathbf{n} + \mathcal{Z}p = g_I \quad \text{on } \Gamma_I \quad (4c)$$

and

$$\Gamma = \Gamma_N \cup \Gamma_D \quad \Gamma_D \cap \Gamma_N = \emptyset.$$

**Remark 1.1:** In this report, we will only consider Dirichlet (4b) and Neumann (4a) boundary conditions. Impedance boundary condition (4c) is useful to consider local *absorbing boundary conditions* which will be considered in a future work.

## 1.1 First-order formulations

As it is usually done in the framework of HDG methods, we will rewrite (2) as a first-order in space system. Notice that we have chosen to keep a second-order dependance in frequency. Adaptation of our method to a first-order in frequency formulation is straightforward.

We will compare two different ways to reach a first-order in space formulation.

To lighten the notations in the remaining of this paper, we introduce the following vector field

$$\mathbf{b}_0 := \rho_0 \mathbf{v}_0,$$

that satisfies the following mass conservation equation

$$\operatorname{div}(\mathbf{b}_0) = 0.$$

### 1.1.1 Diffusive flux formulation:

We begin by introducing the *diffusive flux*

$$\mathbf{q} := -\mathbf{K}_0 \nabla p$$

as a new unknown, leading to the following first-order in space system

$$\mathbf{W}_0 \mathbf{q} + \nabla p = 0 \quad (5a)$$

$$-\rho_0 \omega^2 p - 2i\omega \mathbf{b}_0 \cdot \nabla p + \operatorname{div}(\mathbf{q}) = s \quad (5b)$$

where

$$\mathbf{W}_0 := \mathbf{K}_0^{-1} = \frac{1}{\rho_0 c_0^2} \left[ \mathbf{Id} + \frac{\mathbf{v}_0 \mathbf{v}_0^T}{c_0^2 - |\mathbf{v}_0|^2} \right]. \quad (6)$$

Note that  $\mathbf{K}_0$  is always invertible thanks to (3), indeed we have

$$\det \mathbf{K}_0 = \rho_0 c_0^2 (c_0^2 - |\mathbf{v}_0|^2) \neq 0.$$

The second equality in (6) comes from the *Sherman-Morrison formula*, see [\[SM50\]](#) :

**Lemma 1.2:**

If  $\mathbf{A} \in GL_n(\mathbb{R})$  and  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ , then  $\mathbf{A} + \mathbf{u}\mathbf{v}^T$  is invertible if and only if  $1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq 0$  and

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u}\mathbf{v}^T \mathbf{A}^{-1}}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}}.$$

With this formulation, the Neumann boundary condition (4a) becomes

$$\mathbf{q} \cdot \mathbf{n} - 2i\omega(\mathbf{b}_0 \cdot \mathbf{n})p = -g_N.$$

**Variational formulation:** We can now write a variational formulation for (5a)–(5b) : Seek  $(\mathbf{q}, p) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$  such that for all  $(\mathbf{r}, w) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$

$$\int_{\mathcal{O}} \mathbf{W}_0 \mathbf{q} \cdot \mathbf{r}^* \, d\mathbf{x} - \int_{\mathcal{O}} p \operatorname{div}(\mathbf{r}^*) \, d\mathbf{x} + \int_{\partial\mathcal{O}} p \mathbf{r}^* \cdot \mathbf{n} \, d\sigma = 0 \quad (7a)$$

$$-\omega^2 \int_{\mathcal{O}} \rho_0 p w^* \, d\mathbf{x} + 2i\omega \int_{\mathcal{O}} p \mathbf{b}_0 \cdot \nabla w^* \, d\mathbf{x} - \int_{\mathcal{O}} \mathbf{q} \cdot \nabla w^* \, d\mathbf{x} + \int_{\partial\mathcal{O}} w^* \mathbf{q} \cdot \mathbf{n} - 2i\omega p w^* \mathbf{b}_0 \cdot \mathbf{n} \, d\sigma = \int_{\mathcal{O}} s w^* \, d\mathbf{x} \quad (7b)$$

where the boundary integrals should formally be interpreted as the duality bracket  $\langle \cdot, \cdot \rangle_{H^{-\frac{1}{2}}(\partial\mathcal{O}), H^{\frac{1}{2}}(\partial\mathcal{O})}$  between  $H^{-\frac{1}{2}}(\partial\mathcal{O})$  and  $H^{\frac{1}{2}}(\partial\mathcal{O})$ .

### 1.1.2 Total flux formulation:

As  $\operatorname{div}(\mathbf{b}_0) = 0$ , we notice that

$$2i\omega \mathbf{b}_0 \cdot \nabla p = \operatorname{div}(2i\omega p \mathbf{b}_0),$$

and we can therefore rewrite (2) as

$$-\rho_0 \omega^2 p - \operatorname{div}(\mathbf{K}_0 \nabla p + 2i\omega p \mathbf{b}_0) = s.$$

This leads to another possible first-order in space formulation. We introduce the *total flux*

$$\boldsymbol{\sigma} := -\mathbf{K}_0 \nabla p - 2i\omega p \mathbf{b}_0,$$

leading to the following system

$$\mathbf{W}_0 \boldsymbol{\sigma} + \nabla p + 2i\omega p \mathbf{W}_0 \mathbf{b}_0 = 0, \quad (8a)$$

$$-\rho_0 \omega^2 p + \operatorname{div}(\boldsymbol{\sigma}) = s. \quad (8b)$$

With this formulation, the Neumann boundary condition (4a) becomes

$$\boldsymbol{\sigma} \cdot \mathbf{n} = -g_N.$$

**Variational formulation:** We can now write a variational formulation for (8a)–(8b) : Seek  $(\boldsymbol{\sigma}, p) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$  such that for all  $(\mathbf{r}, w) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$

$$\int_{\mathcal{O}} \mathbf{W}_0 \boldsymbol{\sigma} \cdot \mathbf{r}^* \, d\mathbf{x} - \int_{\mathcal{O}} p \operatorname{div}(\mathbf{r}^*) \, d\mathbf{x} + 2i\omega \int_{\mathcal{O}} p \mathbf{W}_0 \mathbf{b}_0 \cdot \mathbf{r}^* \, d\mathbf{x} + \int_{\partial\mathcal{O}} p \mathbf{r}^* \cdot \mathbf{n} \, d\sigma = 0 \quad (9a)$$

$$-\omega^2 \int_{\mathcal{O}} \rho_0 p w^* \, d\mathbf{x} - \int_{\mathcal{O}} \boldsymbol{\sigma} \cdot \nabla w^* \, d\mathbf{x} + \int_{\partial\mathcal{O}} w^* \boldsymbol{\sigma} \cdot \mathbf{n} \, d\sigma = \int_{\mathcal{O}} s w^* \, d\mathbf{x} \quad (9b)$$

where the boundary integrals should formally be interpreted as the duality bracket  $\langle \cdot, \cdot \rangle_{H^{-\frac{1}{2}}(\partial\mathcal{O}), H^{\frac{1}{2}}(\partial\mathcal{O})}$  between  $H^{-\frac{1}{2}}(\partial\mathcal{O})$  and  $H^{\frac{1}{2}}(\partial\mathcal{O})$ .

## 2 Notations

In this section, we introduce the notations and approximation spaces that will be used to construct the HDG methods considered in this paper.

### 2.1 Approximation spaces

We consider a mesh  $\mathcal{T}_h$  of the domain  $\mathcal{O}$  of dimension  $n$ . For an element  $K \in \mathcal{T}_h$ , we denote by  $\mathcal{E}(K)$  the set of its edges. We also consider

$$\begin{aligned} \text{The set of boundary edges:} & \quad \mathcal{E}_h^b := \{e = \partial K \cap \Gamma \mid K \in \mathcal{T}_h\}, \\ \text{The set of interior edges:} & \quad \mathcal{E}_h^i := \{e = \partial K_+ \cap \partial K_- \mid K_+, K_- \in \mathcal{T}_h\}, \\ \text{The set of all edges:} & \quad \mathcal{E}_h := \mathcal{E}_h^b \cup \mathcal{E}_h^i. \end{aligned}$$

To study the convergence of the methods, we will assume that the mesg has the usual *shape-regularity* property, see [EG04, Def. 1.107].

For  $K \in \mathcal{T}_h$ , we denote by  $\mathcal{P}_k(K)$  the space of polynomials of total degree at most  $k$  defined on  $K$ . We will also use the space of vectorial polynomials  $\mathcal{P}_k(K) = \mathcal{P}_k(K)^n$ . Even if those spaces can be defined for  $k > 0$ , in this paper we will usually assume that  $k \geq 2$  as HDG method of lower order have no interest from a computational point of view.

On each element  $K \in \mathcal{T}_h$ , we introduce the following approximation spaces for the pressure and the flux

$$\begin{aligned} \mathbf{V}_h(K) &:= \left\{ \mathbf{q} \in \mathbf{L}^2(\mathcal{O}) \mid \mathbf{q}|_K \in \mathcal{P}_k(K) \right\} & \text{for the flux } \mathbf{q}_h \text{ or } \boldsymbol{\sigma}_h, \\ W_h(K) &:= \left\{ p \in L^2(\mathcal{O}) \mid p \in \mathcal{P}_\ell(K) \right\} & \text{for the pressure } p_h, \end{aligned}$$

where  $\ell$  can be equal to  $k$  or  $k + 1$  depending on the formulation.

To construct HDG formulations, we will need to add a surfacic unknown, called the *numerical trace* and denoted by  $\widehat{p}_h$ , to the problem. This unknown will be the main unknown of the method as the *static condensation* process will allow to eliminate the volumetric unknowns to obtain a so-called *global problem*. To approximate this new unknown we introduce the following space for  $e \in \mathcal{E}(K)$

$$M_h(e) := \left\{ \mu \in L^2(\mathcal{E}_h) \mid \mu|_e \in \mathcal{P}_k(e) \right\}.$$

As those approximation spaces are discontinuous, we can construct the *global approximation spaces* as the cartesian product of the local ones

$$\begin{aligned} \mathbf{V}_h &:= \prod_{K \in \mathcal{T}_h} \mathbf{V}_h(K) & \text{for the flux } \mathbf{q}_h \text{ or } \boldsymbol{\sigma}_h, \\ W_h &:= \prod_{K \in \mathcal{T}_h} W_h(K) & \text{for the pressure } p_h, \\ M_h &:= \prod_{e \in \mathcal{E}_h} M_h(e) & \text{for the trace } \widehat{p}_h. \end{aligned}$$

In [FIGURE 1](#), we have depicted the differences in the *degrees of freedom* for the continuous (CG), discontinuous (DG) and hybridizable discontinuous (HDG) Galerkin methods. The degrees of freedom of the HDG methods are the ones associated with the numerical trace  $\widehat{p}_h$ . As the numerical cost of the method is directly linked to the number of degrees of freedom, we can clearly see that the HDG method is less expensive than the DG method.

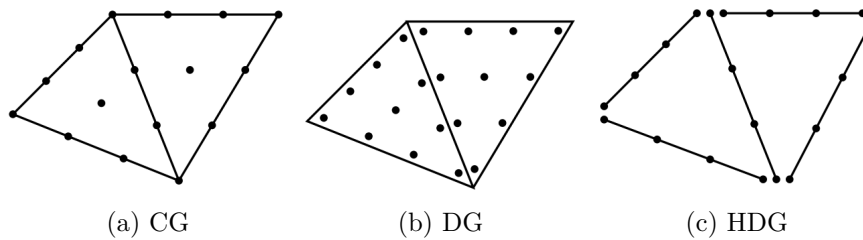


Figure 1: Polynomial interpolation of degree 3

In TABLE 1, we give a summary of the choice of local spaces for the different variations of the HDG method considered in this report.

Variable	Space	HDG $(p_h, \boldsymbol{\sigma})$	HDG $(p_h, \mathbf{q}_h)$	HDG+ $(p_h, \mathbf{q}_h)$
Pressure $p_h$	$W_h(K)$	$\mathcal{P}_k(K)$		$\mathcal{P}_{k+1}(K)$
Flux $\mathbf{q}_h$ or $\boldsymbol{\sigma}_h$	$\mathbf{V}_h(K)$	$\mathcal{P}_k(K)$		
Trace $\widehat{p}_h$	$M_h(e)$	$\mathcal{P}_k(e)$		

Table 1: Choice of local spaces and penalization parameter for the different methods

For the HDG+ method, we will also need the orthogonal projection on the space of piecewise polynomial functions on the edges

$$P_M : \prod_{K \in \mathcal{T}_h} L^2(\partial K) \longrightarrow \prod_{K \in \mathcal{T}_h} \prod_{e \in \mathcal{E}(K)} \mathcal{P}_k(e).$$

It is important to emphasize the difference between this space and  $M_h$ . Indeed as

$$M_h = \prod_{e \in \mathcal{E}_h} \mathcal{P}_k(e) \neq \prod_{K \in \mathcal{T}_h} \prod_{e \in \mathcal{E}(K)} \mathcal{P}_k(e),$$

the functions in  $M_h$  are single-valued on the skeleton of the mesh, whereas the functions in the other space are multi-valued on the interior edges. Functions in both spaces are discontinuous at the vertices.

**Remark 2.1:** It is also possible to choose a continuous space for  $\widehat{p}_h$ , this leads to the so-called *Locally Discontinuous but Globally Continuous method* (LDGC), see eg. [ALA13, FLd14]. However this choice does not seem to improve the convergence rate of the method.

## 2.2 Hermitian products and norms

For an element  $K \in \mathcal{T}_h$ , we denote the standard  $L^2$ -hermitian product<sup>1</sup> and its associated norm by

$$(u, v)_K := \int_K u \cdot v^* d\mathbf{x} \quad \text{and} \quad \|u\|_K^2 := (u, u)_K,$$

we then introduce the broken hermitian product and norm

$$(u, v)_{\mathcal{T}_h} := \sum_{K \in \mathcal{T}_h} (u, v)_K \quad \text{and} \quad \|u\|_{\mathcal{T}_h}^2 := \sum_{K \in \mathcal{T}_h} \|u\|_K^2.$$

<sup>1</sup>For vector fields, the  $\mathbb{R}^n$  dot-product is used inside the integral as the conjugate is already applied.

On the boundary of an element  $K$ , we also denote the local hermitian product by

$$\langle u, v \rangle_{\partial K} := \sum_{e \in \mathcal{E}(K)} \int_e u \cdot v^* d\sigma \quad \text{and} \quad \|u\|_{\partial K}^2 := \langle u, u \rangle_{\partial K},$$

and the broken hermitian product is denoted by

$$\langle u, v \rangle_{\partial \mathcal{T}_h} := \sum_{K \in \mathcal{T}_h} \langle u, v \rangle_{\partial K} \quad \text{and} \quad \|u\|_{\partial \mathcal{T}_h}^2 := \sum_{K \in \mathcal{T}_h} \|u\|_{\partial K}^2.$$

We have chosen to use angle brackets  $\langle \cdot, \cdot \rangle$  to denote the boundary integrals as they should formally be interpreted as the duality bracket between  $H^{-\frac{1}{2}}(\partial K)$  and  $H^{\frac{1}{2}}(\partial K)$ .

We also define the following weighted norms

$$\begin{aligned} \|u\|_{\rho_0, K}^2 &:= (\rho_0 u, u)_K && \text{which satisfies} \quad \|u\|_{\rho_0, K} \leq \|\rho_0\|_{L^\infty(K)}^{\frac{1}{2}} \|u\|_K \\ \|\mathbf{q}\|_{\mathbf{W}_0, K}^2 &:= (\mathbf{W}_0 \mathbf{q}, \mathbf{q})_K && \text{which satisfies} \quad \|\mathbf{q}\|_{\mathbf{W}_0, K} \leq C_{\mathbf{W}_0, K} \|\mathbf{q}\|_K \end{aligned}$$

where

$$C_{\mathbf{W}_0, K} = \left( \max_K \frac{1}{\rho_0 (c_0^2 - |\mathbf{v}_0|^2)} \right)^{\frac{1}{2}}$$

is the largest eigenvalue of  $\mathbf{W}_0$  in  $K$ , see [LEMMA 1.1](#).

### 2.3 Faces, jumps and averages

In this subsection, we will introduce notations for the faces quantities. As usual with methods belonging to the DG family, we will need to define jumps and averages which link the unknowns between two elements.

**Faces and normals:** For an interior face  $\mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-$ , we denote by  $\mathbf{n}^+$  (resp.  $\mathbf{n}^-$ ) a unitary outgoing normal vector of  $\partial K_+$  (resp.  $\partial K_-$ ). We will always assume that the flow  $\mathbf{v}_0$  goes from  $K_-$  to  $K_+$ , as depicted on [FIGURE 2](#).

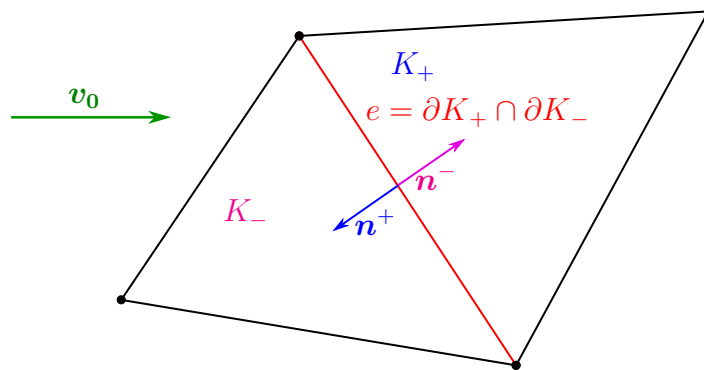


Figure 2: Normal vectors on an interior face

When the orientation of the face does not matter, we will denote by  $\mathbf{n}$  any unitary normal vector to  $e$ .

If  $e$  is a boundary edge, then  $\mathbf{n}$  denotes the outward-pointing unitary normal vector.

**Jumps and averages:** We will often use the *average operator* defined by

$$\begin{aligned} \text{On } \mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-, & \quad \{\!\!\{ \varphi \}\!\!\}_e := \frac{1}{2} (\varphi^+ + \varphi^-), \\ \text{On } \mathcal{E}_h^b \ni e = \partial K \cap \Gamma, & \quad \{\!\!\{ \varphi \}\!\!\}_e := \frac{1}{2} \varphi, \end{aligned}$$

where  $\varphi$  can either be a scalar or vectorial quantity.

We will also make frequent use of the *jump operator* defined by

$$\begin{aligned} \text{On } \mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-, & \quad [\![ \mathbf{q} ]\!]_e := \mathbf{q}^+ \cdot \mathbf{n}^+ + \mathbf{q}^- \cdot \mathbf{n}^-, \\ \text{On } \mathcal{E}_h^b \ni e = \partial K \cap \Gamma, & \quad [\![ \mathbf{q} ]\!]_e := \mathbf{q} \cdot \mathbf{n}, \end{aligned}$$

for a vectorial quantity. Notice that with this definition, the jump operator only controls the normal part of the vector. For a scalar quantity, the *jump operator* is defined by

$$\begin{aligned} \text{On } \mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-, & \quad [p]_e := p^+ - p^-, \\ \text{On } \mathcal{E}_h^b \ni e = \partial K \cap \Gamma, & \quad [p]_e := p, \end{aligned}$$

for a scalar quantity. A sketch of those quantities is given in [FIGURE 3](#).

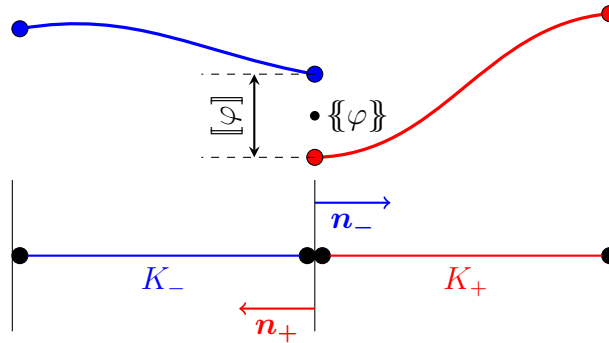


Figure 3: 1D-sketch of the jump and average on an interior node

### 3 HDG method for the total flux formulation

In this section, we will focus on the *total flux formulation*. We will first construct the HDG method and we will then discuss its most important properties.

#### 3.1 Constructing the formulation

On an element  $K \in \mathcal{T}_h$ , we recall that the weak formulation (9a)–(9b) reads : seek  $(\boldsymbol{\sigma}, p) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$  such that for all  $(\mathbf{r}, w) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$

$$\int_K \mathbf{W}_0 \boldsymbol{\sigma} \cdot \mathbf{r}^* \, d\mathbf{x} - \int_K p \operatorname{div}(\mathbf{r}^*) \, d\mathbf{x} + 2i\omega \int_K p \mathbf{W}_0 \mathbf{b}_0 \cdot \mathbf{r}^* \, d\mathbf{x} + \int_{\partial K} p \mathbf{r}^* \cdot \mathbf{n} \, d\sigma = 0, \quad (11a)$$

$$-\omega^2 \int_K \rho_0 p w^* \, d\mathbf{x} + \int_K \operatorname{div}(\boldsymbol{\sigma}) w^* \, d\mathbf{x} = \int_K s w^* \, d\mathbf{x} \quad (11b)$$

**Choice of approximation spaces:** We denote by  $\boldsymbol{\sigma}_h$  and  $p_h$  the approximations of  $\boldsymbol{\sigma}$  and  $p$  on  $K$ .

For this method, we choose to use the following local approximation spaces

$$\begin{aligned} \mathbf{V}_h(K) &= \mathcal{P}_k(K) && \text{for the flux } \boldsymbol{\sigma}_h, \\ W_h(K) &= \mathcal{P}_k(K) && \text{for the pressures } p_h, \end{aligned}$$

where  $k \geq 3$  is the degree of the method. We recall that if  $k \leq 2$ , then there are no interior degrees of freedom and the HDG method has no interest over the DG methods from a computational point of view. Notice that we use the same interpolation degree for both unknowns, which may lead to unstable continuous Galerkin methods, see *eg.* [EG04, "Checkerboard-like instability" p.188]. As we will discuss later, this is not a problem for HDG methods.

**Introduction of the hybrid unknown:** To reach a HDG formulation, we introduce a new unknown  $\widehat{p}_h$  which is an approximation of  $p$  on  $\mathcal{E}_h$ , the skeleton of the mesh  $\mathcal{T}_h$ . We will usually refer to  $\widehat{p}_h$  as the *numerical trace*. This unknown is the main unknown of the HDG method. Indeed, we will be able to use a *static condensation process* to eliminate the interior degrees of freedom and to obtain a so-called *global problem* for  $\widehat{p}_h$  only. To introduce this unknown in the formulation, the boundary integral in (11a) is discretized as follows

$$\int_{\partial K} \mathbf{p} \mathbf{r}^* \cdot \mathbf{n} d\sigma \quad \text{becomes} \quad \int_{\partial K} \widehat{p}_h \mathbf{r}_h^* \cdot \mathbf{n} d\sigma.$$

For this new unknown, we will use the following approximation space

$$M_h(e) = \mathcal{P}_k(e), \quad \forall e \in \mathcal{E}(K).$$

**Penalization parameter:** The unknown  $\widehat{p}_h$  is often called a *Lagrange multiplier*. Indeed, when going from a continuous Galerkin method to a HDG one, the continuity of the numerical solution is not strongly enforced anymore and it is added in the method as a constraint. The quantity  $\widehat{p}_h$  is therefore the Lagrange multiplier that enforces this weak continuity requirement. To enforce this constrain, we introduce a *penalization parameter* denoted by  $\tau$ . and the following boundary term

$$\langle \tau(p_h - \widehat{p}_h), w_h \rangle_{\partial K}$$

will be added to the local problem. This boundary term can be interpreted as a weak enforcement of the following Dirichlet boundary condition

$$p_h = \widehat{p}_h, \quad \text{on } \partial K.$$

Practical choice of  $\tau$  will be discussed later.

**Local problem:** We approximate the variational formulation (11a)–(11b) on an element  $K \in \mathcal{T}_h$  leading to the *local problem* : seek  $(\boldsymbol{\sigma}_h, p_h) \in \mathbf{V}_h(K) \times W_h(K)$  such that

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_K - (p_h, \operatorname{div}(\mathbf{r}_h))_K + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_K + \langle \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial K} = 0, \quad (12a)$$

$$-\omega^2 (\rho_0 p_h, w_h)_K + (\operatorname{div}(\boldsymbol{\sigma}_h), w_h)_K + i\omega \langle \tau(p_h - \widehat{p}_h), w_h \rangle_{\partial K} = (s, w_h)_K, \quad (12b)$$

for all  $(\mathbf{r}_h, w_h) \in \mathbf{V}_h(K) \times W_h(K)$ .

Notice that (12a)–(12b) is the variational formulation of the convected Helmholtz equation on  $K$  with weak Dirichlet boundary conditions on  $\partial K$ .

**Transmission condition:** Due to the discontinuous nature of the approximation spaces, we need to link all the local problems together. To this end, we introduce the *numerical flux* for  $\boldsymbol{\sigma}_h$

$$\widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n} := \boldsymbol{\sigma}_h \cdot \mathbf{n} + i\omega\tau(p_h - \widehat{p}_h), \quad (13)$$

which satisfies the following *transmission condition*

$$\langle \widehat{\boldsymbol{\sigma}}_h, \mu_h \rangle_{\partial\mathcal{T}_h \setminus \Gamma_D} + \langle \widehat{p}_h - g_D, \mu_h \rangle_{\Gamma_D} = \langle g_N, \mu_h \rangle_{\Gamma_N} \quad (14)$$

for all  $\mu_h \in M_h$ .

Notice that (14) enforces the normal continuity of  $\widehat{\boldsymbol{\sigma}}_h$  on the interior faces as well as the Neumann and Dirichlet boundary conditions on  $\Gamma_N$  and  $\Gamma_D$ .

Indeed on the interior faces we have

$$\langle \widehat{\boldsymbol{\sigma}}_h, \mu_h \rangle_{\partial\mathcal{T}_h \setminus \Gamma_D} = \sum_{e \in \mathcal{E}_h^i} \langle \llbracket \widehat{\boldsymbol{\sigma}}_h \rrbracket, \mu_h \rangle_e = 0,$$

as  $\boldsymbol{\sigma}_h \in \mathbf{V}_h$ ,  $p_h \in W_h$  and  $\widehat{p}_h \in M_h$ , we have  $\llbracket \widehat{\boldsymbol{\sigma}}_h \rrbracket \in M_h$  and we therefore conclude that  $\llbracket \widehat{\boldsymbol{\sigma}}_h \rrbracket = 0$ , as  $\llbracket \widehat{\boldsymbol{\sigma}}_h \rrbracket$  is a polynomial of degree up to  $k$  orthogonal to all polynomials of degree up to  $k$ . We recall that on an interior edge  $\mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-$ , the jump operator is defined as

$$\llbracket \boldsymbol{\sigma}_h \rrbracket := \boldsymbol{\sigma}_h^+ \cdot \mathbf{n}^+ + \boldsymbol{\sigma}_h^- \cdot \mathbf{n}^-.$$

**Remark 3.1:** The transmission condition (14) can be understood as a weak requirement of  $\mathbf{H}_{\text{div}}(\mathcal{O})$ -conformity. Indeed it is shown in [PE12, Lemma 1.2.4] that  $\boldsymbol{\sigma}_h \in \mathbf{H}_{\text{div}}(\mathcal{O})$  means

$$\forall K \in \mathcal{T}_h, \boldsymbol{\sigma}_h^K \in \mathbf{H}_{\text{div}}(K) \quad \text{and} \quad \forall e \in \mathcal{E}_h^i, \llbracket \boldsymbol{\sigma}_h \rrbracket_e \equiv 0.$$

The former is a consequence of the polynomial nature of the approximation spaces, and we will now focus on the latter. Owing to the transmission condition, we have

$$\forall e \in \mathcal{E}_h^i, 0 = \llbracket \widehat{\boldsymbol{\sigma}}_h \rrbracket = \llbracket \boldsymbol{\sigma}_h \rrbracket + i\omega \llbracket \tau(p_h - \widehat{p}_h) \rrbracket.$$

As  $p_h$  and  $\widehat{p}_h$  are two approximations of the same unknown  $p$ , the quantity  $p_h - \widehat{p}_h$  is expected to be small. We can therefore conclude that  $\llbracket \boldsymbol{\sigma}_h \rrbracket$  is small and that

$$\llbracket \boldsymbol{\sigma}_h \rrbracket \xrightarrow{h_K \rightarrow 0} 0.$$

For applications where a precise approximation of the flux is required, it is possible to post-process  $\boldsymbol{\sigma}_h$  to obtain a new approximate  $\widetilde{\boldsymbol{\sigma}}_h$  with strong  $\mathbf{H}_{\text{div}}$ -conformity, see [CGS10, Sec. 5.1].

### 3.1.1 Compact formulation:

HDG methods are usually stated in a compact form that can be obtained by summing the local problems (12a)–(12b) over the mesh elements and by adding the transmission condition (14). This formulation reads : seek  $(\boldsymbol{\sigma}_h, p_h, \widehat{p}_h) \in \mathbf{V}_h \times W_h \times M_h$  such that

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_{\mathcal{T}_h} - (p_h, \text{div}(\mathbf{r}_h))_{\mathcal{T}_h} + 2i\omega(p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} + \langle \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} = 0, \quad (15a)$$

$$-\omega^2(\rho_0 p_h, w_h)_{\mathcal{T}_h} + (\text{div}(\boldsymbol{\sigma}_h), w_h)_{\mathcal{T}_h} + i\omega \langle \tau(p_h - \widehat{p}_h), w_h \rangle_{\partial\mathcal{T}_h} = (s, w_h)_{\mathcal{T}_h}, \quad (15b)$$

$$\langle \boldsymbol{\sigma}_h \cdot \mathbf{n} + i\omega\tau(p_h - \widehat{p}_h), \mu_h \rangle_{\partial\mathcal{T}_h \setminus \Gamma_D} + \langle \widehat{p}_h - g_D, \mu_h \rangle_{\Gamma_D} = \langle g_N, \mu_h \rangle_{\Gamma_N}, \quad (15c)$$

for all  $(\mathbf{r}_h, w_h, \mu_h) \in \mathbf{V}_h \times W_h \times M_h$ . This formulation will be useful to perform the numerical analysis of the method.



**Remark 3.2:** At this point, to completely define the HDG method, it only remains to choose the penalization parameter  $\tau$ , this will be done in the next section.

### 3.1.2 Condensed variational formulation

The compact formulation (15a)–(15b)–(15c) cannot directly be used to efficiently implement the HDG method. Indeed it is not clear how a formulation involving only  $\widehat{p}_h$  can be reached. To emphasize how it can be done, we will now write a *condensed* variational formulation for  $\widehat{p}_h$  only.

We introduce the so-called *local solvers*

$$\begin{aligned} \mathbf{P}^K &: (\widehat{p}_h, s) \mapsto p_h^K, \\ \Sigma^K &: (\widehat{p}_h, s) \mapsto \boldsymbol{\sigma}_h^K, \\ \widehat{\Sigma}^K &: (\widehat{p}_h, s) \mapsto \widehat{\boldsymbol{\sigma}}_h^K, \end{aligned}$$

where  $(\boldsymbol{\sigma}_h^K, p_h^K)$  is the solution of (12a)–(12b) and  $\widehat{\boldsymbol{\sigma}}_h^K$  is defined by (13).

We can therefore rewrite the transmission condition (15c) as

$$a_h(\widehat{p}_h, \mu) = \ell_h(\mu), \quad (16)$$

where

$$\begin{aligned} a_h(\widehat{p}_h, \mu_h) &:= \langle \Sigma^K(\widehat{p}_h, s) \cdot \mathbf{n} + i\omega\tau(\mathbf{P}^K(\widehat{p}_h, s) - \widehat{p}_h), \mu_h \rangle_{\partial\mathcal{T}_h \setminus \Gamma_D} + \langle \widehat{p}_h, \mu_h \rangle_{\Gamma_D}, \\ \ell_h(\mu_h) &:= \langle g_N, \mu_h \rangle_{\Gamma_N} + \langle g_D, \mu_h \rangle_{\Gamma_D}. \end{aligned}$$

Equation (16) is the so-called *global problem* and is the main equation of the HDG method. From a computational point of view, we proceed as described in [ALGORITHM 1](#).

---

#### Algorithm 1: Solving HDG- $\boldsymbol{\sigma}_h$

---

- 1 **for**  $K \in \mathcal{T}_h$  **do**
  - 2     Construct the local solvers  $\mathbf{P}^K, \Sigma^K, \widehat{\Sigma}^K$
  - 3     Add local contribution to the global problem (16)
  - 4 Solve (16) for  $\widehat{p}_h$  // This is the main step
  - 5 **for**  $K \in \mathcal{T}_h$  **do**
  - 6     Reconstruct the local unknowns  $p_h^K = \mathbf{P}^k(\widehat{p}_h, s)$  and  $\boldsymbol{\sigma}_h^K = \Sigma(\widehat{p}_h, s)$
- 

This algorithm is the blueprint of the practical implementation of the HDG method which will be discussed in [SECTION 5](#).

## 3.2 Choice of penalization parameter

In this section, we will show how the penalization parameter can be chosen to obtain an upwinding mechanism with physical meaning. To do that, we will first need to rewrite the HDG method as a DG one, we will then solve a Riemann problem to obtain the value of  $\tau$ .

### 3.2.1 DG formulation

In this section, we will rewrite the HDG method (15a)–(15b)–(15c) as a standard discontinuous Galerkin method.

We introduce the following bilinear form

$$\begin{aligned} \mathcal{B}_h([\boldsymbol{\sigma}_h, p_h]; [\mathbf{r}_h, w_h]) &:= (\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_{\mathcal{T}_h} - (p_h, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} \\ &\quad - \omega^2 (\rho_0 p_h, w_h)_{\mathcal{T}_h} + (\operatorname{div}(\boldsymbol{\sigma}_h), w_h)_{\mathcal{T}_h} \\ &\quad + \sum_{e \in \mathcal{E}_h} \left( \langle \widehat{p}_h, \llbracket \mathbf{r}_h \rrbracket \rangle_e + \langle \widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n}, \llbracket w_h \rrbracket \rangle_e \right), \end{aligned} \quad (17)$$

from which all the mixed DG methods can be generated by choosing  $\widehat{p}_h$  and  $\widehat{\boldsymbol{\sigma}}_h$ . Notice that  $\widehat{p}_h$  was an unknown of the HDG method whereas it now should be chosen by the user of the method.

For example, the LDG method is obtained by choosing

$$\widehat{p}_h = \{\!\{ p_h \}\!\} + \alpha \llbracket p_h \rrbracket \quad \text{and} \quad \widehat{\boldsymbol{\sigma}}_h = \{\!\{ \boldsymbol{\sigma}_h \}\!\} + \beta \llbracket \boldsymbol{\sigma}_h \rrbracket \mathbf{n} + \gamma \llbracket p_h \rrbracket \mathbf{n},$$

and the DG method with central flux is obtained by choosing

$$\widehat{p}_h = \{\!\{ p_h \}\!\} \quad \text{and} \quad \widehat{\boldsymbol{\sigma}}_h = \{\!\{ \boldsymbol{\sigma} \}\!\} - \alpha \llbracket p_h \rrbracket.$$

### Proposition 3.1:

The HDG method (15a)–(15b)–(15c) and the DG method associated to the bilinear form (17) are equivalent if and only if

$$\widehat{p}_h = \{\!\{ p_h \}\!\} + \frac{\tau^+ - \tau^-}{2(\tau^+ + \tau^-)} \llbracket p_h \rrbracket + \frac{1}{i\omega(\tau^+ + \tau^-)} \llbracket \boldsymbol{\sigma}_h \rrbracket, \quad (18a)$$

$$\widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n} = \{\!\{ \boldsymbol{\sigma}_h \}\!\} \cdot \mathbf{n} + i\omega \frac{\tau^+ \tau^-}{\tau^+ + \tau^-} \llbracket p_h \rrbracket - \frac{\tau^+ - \tau^-}{2(\tau^+ + \tau^-)} \llbracket \boldsymbol{\sigma}_h \rrbracket, \quad (18b)$$

for all interior edges  $\mathcal{E}_h^b \ni e = \partial K_+ \cap \partial K_-$  and where  $\tau^\pm = \tau|_{\partial K_\pm}$ .

We recall the convention for the labelling  $\partial K_\pm : \mathbf{v}_0$  is directed from  $\partial K_-$  toward  $\partial K_+$  and we denote by  $\mathbf{n}$  any normal vector to the face when the orientation does not matter.

**Proof:** Writing down the transmission condition (15c) on an interior face  $e$ , we have

$$\llbracket \boldsymbol{\sigma}_h \rrbracket + 2i\omega \left( \{\!\{ \tau \}\!\} \{\!\{ p_h \}\!\} + \frac{1}{4} \llbracket \tau \rrbracket \llbracket p_h \rrbracket \right) - 2i\omega \{\!\{ \tau \}\!\} \widehat{p}_h = 0,$$

which leads to

$$\widehat{p}_h = \{\!\{ p_h \}\!\} + \frac{\llbracket \tau \rrbracket}{4\{\!\{ \tau \}\!\}} \llbracket p_h \rrbracket + \frac{1}{2i\omega \{\!\{ \tau \}\!\}} \llbracket \boldsymbol{\sigma}_h \rrbracket,$$

and we obtain (18a) by developing the jumps and average terms.

As the numerical flux  $\widehat{\boldsymbol{\sigma}}_h$  is continuous across the interface, we have

$$\begin{aligned} \widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n} &= \{\!\{ \widehat{\boldsymbol{\sigma}}_h \}\!\} \cdot \mathbf{n} \\ &= \{\!\{ \boldsymbol{\sigma}_h \}\!\} \cdot \mathbf{n} + \frac{i\omega}{2} \llbracket \tau \rrbracket (\{\!\{ p_h \}\!\} - \widehat{p}_h) + \frac{i\omega}{2} \{\!\{ \tau \}\!\} \llbracket p_h \rrbracket \\ &= \{\!\{ \boldsymbol{\sigma}_h \}\!\} \cdot \mathbf{n} - \frac{i\omega}{2} \llbracket \tau \rrbracket \left( \frac{\llbracket \tau \rrbracket}{4\{\!\{ \tau \}\!\}} \llbracket p_h \rrbracket + \frac{1}{2i\omega \{\!\{ \tau \}\!\}} \llbracket \boldsymbol{\sigma}_h \rrbracket \right) + \frac{i\omega}{2} \{\!\{ \tau \}\!\} \llbracket p_h \rrbracket \end{aligned}$$

and we obtain (18b) as

$$-\frac{\llbracket \tau \rrbracket^2}{4\{\!\{ \tau \}\!\}} + \{\!\{ \tau \}\!\} = -\frac{(\tau^+)^2 - 2\tau^+ \tau^- + (\tau^-)^2}{2(\tau^+ + \tau^-)} + \frac{(\tau^+)^2 + 2\tau^+ \tau^- + (\tau^-)^2}{2(\tau^+ + \tau^-)} = 2 \frac{\tau^+ \tau^-}{\tau^+ + \tau^-}.$$

We now have to show that this choice of numerical flux  $\widehat{\boldsymbol{\sigma}}_h$  is compatible with the HDG method. Starting from (18b) we have

$$\begin{aligned}\widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n}^+ &= \boldsymbol{\sigma}_h^+ \cdot \mathbf{n}^+ - \frac{1}{2} \left( 1 + \frac{\tau^+ - \tau^-}{\tau^+ + \tau^-} \right) \llbracket \boldsymbol{\sigma}_h \rrbracket + i\omega \frac{\tau^+ \tau^-}{\tau^+ + \tau^-} \llbracket p_h \rrbracket \\ &= \boldsymbol{\sigma}_h^+ \cdot \mathbf{n}^+ - \frac{\tau^+}{\tau^+ + \tau^-} \llbracket \boldsymbol{\sigma}_h \rrbracket + i\omega \frac{\tau^+ \tau^-}{\tau^+ + \tau^-} \llbracket p_h \rrbracket,\end{aligned}$$

on the other hand, rewriting (18a) gives

$$\begin{aligned}-\frac{\tau^+}{\tau^+ + \tau^-} \llbracket \boldsymbol{\sigma}_h \rrbracket &= i\omega \tau^+ \left[ \{\{p_h\}\} - \widehat{p}_h + \frac{\tau^+ - \tau^-}{2(\tau^+ + \tau^-)} \llbracket p_h \rrbracket \right] \\ &= i\omega \tau^+ (p_h^+ - \widehat{p}_h) + i\omega \frac{\tau^+}{2} \left( \frac{\tau^+ - \tau^-}{\tau^+ + \tau^-} - 1 \right) \llbracket p_h \rrbracket \\ &= i\omega \tau^+ (p_h^+ - \widehat{p}_h) - i\omega \frac{\tau^+ \tau^-}{\tau^+ + \tau^-} \llbracket p_h \rrbracket,\end{aligned}$$

so we finally have

$$\widehat{\boldsymbol{\sigma}}_h \cdot \mathbf{n}^+ = \boldsymbol{\sigma}_h^+ \cdot \mathbf{n}^+ + i\omega \tau^+ (p_h^+ - \widehat{p}_h).$$

Similar computations can be carried out on  $\partial K_-$ . ■

**Particular form of the HDG fluxes:** We would like to point out that  $\widehat{p}_h$  depends on  $\llbracket \boldsymbol{\sigma}_h \rrbracket$ , this is a distinctive feature of HDG methods among the family of DG methods. To understand this, let us consider DG method with the following fluxes

$$\widehat{p}_h = \{\{p_h\}\} + \alpha \llbracket p_h \rrbracket \quad \text{and} \quad \widehat{\boldsymbol{\sigma}}_h = \{\{\boldsymbol{\sigma}_h\}\} + \beta \llbracket \boldsymbol{\sigma}_h \rrbracket \mathbf{n} + \gamma \llbracket p_h \rrbracket \mathbf{n}, \quad (19)$$

where  $\alpha, \beta, \gamma$  are arbitrary constants. This construction is adapted from [HW08, Sec 7.2.2]. Testing (17) with  $[\mathbf{r}_h, 0]$  leads to

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_{\mathcal{T}_h} - (p_h, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} + \sum_{e \in \mathcal{E}_h} \langle \widehat{p}_h, [\mathbf{r}_h] \rangle_e = 0.$$

Integrating by parts leads to

$$\begin{aligned}(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_{\mathcal{T}_h} &= -(\nabla p_h, \mathbf{r}_h)_{\mathcal{T}_h} - 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} + \sum_{e \in \mathcal{E}_h} [\langle \llbracket p_h \rrbracket \mathbf{n}, \{\{\mathbf{r}_h\}\} \rangle_e - \langle \widehat{p}_h, [\mathbf{r}_h] \rangle_e] \\ &\quad + \sum_{e \in \mathcal{E}_h^i} \langle \{\{p_h\}\}, [\mathbf{r}_h] \rangle_e,\end{aligned} \quad (20)$$

where we used the identity

$$\langle p_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = \sum_{e \in \mathcal{E}_h} \langle \llbracket p_h \rrbracket \mathbf{n}, \{\{\mathbf{r}_h\}\} \rangle_e + \sum_{e \in \mathcal{E}_h^i} \langle \{\{p_h\}\}, [\mathbf{r}_h] \rangle_e,$$

coming from [HW08, Lemma 7.9]. Using the definition of  $\widehat{p}_h$  given in (19), the surfacic terms in (20) become

$$-\sum_{e \in \mathcal{E}_h^i} \langle \llbracket p_h \rrbracket, \{\{\mathbf{r}_h\}\} \cdot \mathbf{n} - \alpha [\mathbf{r}_h] \rangle_e - \sum_{e \in \mathcal{E}_h^b} \langle p_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_e.$$

We now introduce the *lifting operator*  $\mathcal{L}$  defined by

$$(\mathcal{L}(p_h), \mathbf{r}_h)_{\mathcal{T}_h} = \sum_{e \in \mathcal{E}_h^i} \langle \llbracket p_h \rrbracket, \{\{\mathbf{r}_h\}\} \cdot \mathbf{n} - \alpha [\mathbf{r}_h] \rangle_e + \sum_{e \in \mathcal{E}_h^b} \langle p_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_e, \quad \forall \mathbf{r}_h \in \mathbf{V}_h,$$

and (20) becomes

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_{\mathcal{T}_h} = (-\nabla p_h - 2i\omega p_h \mathbf{W}_0 \mathbf{b}_0 - \mathcal{L}(p_h), \mathbf{r}_h)_{\mathcal{T}_h}.$$

We can see that  $\boldsymbol{\sigma}_h$  is completely defined in terms of  $p_h$  and it is therefore not possible to add a transmission condition, which is required to allow the static condensation process.

On the other hand, if the HDG flux

$$\widehat{p}_h = \{\{p_h\}\} + \alpha \llbracket p_h \rrbracket + \delta \llbracket \boldsymbol{\sigma}_h \rrbracket$$

is used, we will obtain an expression of  $\boldsymbol{\sigma}_h$  in terms of  $p_h$  and  $\llbracket \boldsymbol{\sigma}_h \rrbracket$ . We will therefore need to add the transmission condition to close the discrete system and it will be possible to perform the static condensation.

### 3.2.2 Computing the penalization parameter

#### Proposition 3.2:

On an interior face  $\mathcal{E}_h^i \ni e = \partial K_+ \cap \partial K_-$  the following choice of penalization parameter

$$\tau^\pm = \rho_0(c_0 + \mathbf{v}_0 \cdot \mathbf{n}^\pm), \quad (21)$$

where  $\tau^\pm = \tau|_{\partial K_\pm}$ , leads to an upwinding mechanism.

To prove this proposition, we will need to solve a Riemann problem and compare its solution with PROPOSITION 3.1 to obtain a value for  $\tau^\pm$  with physical meaning. The first step to be able to solve the Riemann problem is to rewrite the original equation as a time-domain hyperbolic system.

**Hyperbolic system:** We start from the convected acoustic wave equation

$$\rho_0 \left( \frac{\partial}{\partial t} + \mathbf{v}_0 \cdot \nabla \right)^2 p - \operatorname{div} (\rho_0 c_0^2 \nabla p) = 0$$

and we write it as a hyperbolic system. First we have

$$\rho_0 \frac{\partial^2 p}{\partial t^2} - \operatorname{div} \left( \mathbf{K}_0 \nabla p - 2\rho_0 \frac{\partial p}{\partial t} \mathbf{v}_0 \right) = 0,$$

we therefore introduce the *total flux*

$$\frac{\partial \tilde{\boldsymbol{\sigma}}}{\partial t} = -\mathbf{K}_0 \nabla p + 2\rho_0 \frac{\partial p}{\partial t} \mathbf{v}_0,$$

leading to the following first-order formulation

$$\frac{\partial p}{\partial t} = -\frac{1}{\rho_0} \operatorname{div} (\tilde{\boldsymbol{\sigma}}), \quad (22a)$$

$$\frac{\partial \tilde{\boldsymbol{\sigma}}}{\partial t} = -\mathbf{K}_0 \nabla p + 2\rho_0 \frac{\partial p}{\partial t} \mathbf{v}_0. \quad (22b)$$

However this formulation does not have the form of a hyperbolic system.

Using (22a) in (22b), we have

$$\frac{\partial p}{\partial t} = -\frac{1}{\rho_0} \operatorname{div}(\tilde{\boldsymbol{\sigma}}), \quad (23a)$$

$$\frac{\partial \tilde{\boldsymbol{\sigma}}}{\partial t} = -\mathbf{K}_0 \nabla p - 2 \operatorname{div}(\tilde{\boldsymbol{\sigma}}) \mathbf{v}_0. \quad (23b)$$

Notice that we need to work with a first-order in time formulation whereas our methods are written for second-order in time (or equivalently in frequency) formulations. However we have the following relationship between  $\boldsymbol{\sigma}$  and  $\tilde{\boldsymbol{\sigma}}$

$$\boldsymbol{\sigma} = i\omega \tilde{\boldsymbol{\sigma}},$$

making it possible to go back to a second-order formulation.

The system (23a)–(23b) can be written as

$$\frac{\partial \mathbf{U}}{\partial t} = \mathbb{A}_x \frac{\partial \mathbf{U}}{\partial x} + \mathbb{A}_y \frac{\partial \mathbf{U}}{\partial y}, \quad (24)$$

where

$$\mathbf{U} := \begin{bmatrix} p \\ \tilde{\boldsymbol{\sigma}} \end{bmatrix} ; \quad \mathbb{A}_x := \begin{bmatrix} 0 & -\frac{1}{\rho_0} & 0 \\ -M_{0,xx} & -2v_{0,x} & 0 \\ -M_{0,yx} & -2v_{0,y} & 0 \end{bmatrix} ; \quad \mathbb{A}_y := \begin{bmatrix} 0 & 0 & -\frac{1}{\rho_0} \\ -M_{0,xy} & 0 & -2v_{0,x} \\ -M_{0,yy} & 0 & -2v_{0,y} \end{bmatrix}.$$

To check that (24) is a hyperbolic system, one needs to show that for all  $\alpha, \beta \in \mathbb{R}$  the matrix

$$\mathbb{A}_{\alpha,\beta} := \alpha \mathbb{A}_x + \beta \mathbb{A}_y = - \begin{bmatrix} 0 & \frac{\alpha}{\rho_0} & \frac{\beta}{\rho_0} \\ \alpha M_{0,xx} + \beta M_{0,xy} & 2\alpha v_{0,x} & 2\beta v_{0,x} \\ \alpha M_{0,yx} + \beta M_{0,yy} & 2\alpha v_{0,y} & 2\beta v_{0,y} \end{bmatrix}$$

is diagonalizable with real eigenvalues, see [FIGURE 4](#).

Input:

eigenvalues	$\begin{pmatrix} 0 & \frac{a}{r} & \frac{b}{r} \\ a e + b f & 2 a u & 2 b u \\ a f + b g & 2 a v & 2 b v \end{pmatrix}$
-------------	-------------------------------------------------------------------------------------------------------------------------

Results:

Exact forms

Step-by-step solution

$$\lambda_1 = 0$$

$$\lambda_2 \approx \frac{-\sqrt{r(a^2 r u^2 + 2.71828 a^2 + 2 a b f + 2 a b r u v + b^2 g + b^2 r v^2)} + a r u + b r v}{r}$$

$$\lambda_3 \approx \frac{\sqrt{r(a^2 r u^2 + 2.71828 a^2 + 2 a b f + 2 a b r u v + b^2 g + b^2 r v^2)} + a r u + b r v}{r}$$

Figure 4: Computation of the eigenvalues of  $\mathbb{A}_{\alpha,\beta}$  with WolframAlpha

**Riemann solver:** To compute the upwind penalization parameters, we consider a vertical interface located at  $x = 0$  and we assume that the background flow is uniform.

We will solve the problem (24) with the following initial condition

$$\begin{aligned} \mathbf{U}(x, y, 0) &= \mathbf{U}^+, & \text{if } x > 0, \\ \mathbf{U}(x, y, 0) &= \mathbf{U}^-, & \text{if } x < 0. \end{aligned}$$

With this choice of initial condition, we obtain a well-posed problem which is invariant with respect to  $y$ .

Our goal is to compute  $\mathbf{U}$  at  $x = 0$ .

Due to the invariance with respect to  $y$ , we can rewrite (24) as

$$\frac{\partial \mathbf{U}}{\partial t} = \mathbb{A}_x \frac{\partial \mathbf{U}}{\partial x}.$$

Furthermore, we can obtain the following system for  $[p, \tilde{\sigma}_x]^T$  only

$$\frac{\partial}{\partial t} \begin{bmatrix} p \\ \tilde{\sigma}_x \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & -\frac{1}{\rho_0} \\ -M_{0,xx} & -2v_{0,x} \end{bmatrix}}_{=: \mathbb{A}} \frac{\partial}{\partial x} \begin{bmatrix} p \\ \tilde{\sigma}_x \end{bmatrix},$$

as the DG method is only written in terms of  $\tilde{\boldsymbol{\sigma}} \cdot \mathbf{n} = \tilde{\sigma}_x$ .

To compute the eigenvalues of  $\mathbb{A}$ , we need to solve

$$\begin{vmatrix} -\lambda & -\frac{1}{\rho_0} \\ -M_{0,xx} & -2v_{0,x} - \lambda \end{vmatrix} = 0 \iff \lambda^2 + 2v_{0,x}\lambda - \frac{M_{0,xx}}{\rho_0} = 0.$$

Recalling that

$$M_{0,xx} = \rho_0 c_0^2 - \rho_0 v_{0,x}^2,$$

we obtain the two following eigenvalues

$$\begin{aligned} \lambda_1 &= -(c_0 + v_{0,x}), \\ \lambda_2 &= c_0 - v_{0,x}, \end{aligned}$$

and the associated eigenvectors are

$$\mathbf{w}_1 := \begin{bmatrix} 1 \\ \rho_0(c_0 + v_{0,x}) \end{bmatrix} \quad \text{and} \quad \mathbf{w}_2 := \begin{bmatrix} 1 \\ \rho_0(v_{0,x} - c_0) \end{bmatrix}.$$

We can now define

$$\mathbf{W} := \begin{bmatrix} 1 & 1 \\ \rho_0(c_0 + v_{0,x}) & \rho_0(v_{0,x} - c_0) \end{bmatrix},$$

and therefore

$$\mathbf{W}^{-1} = \frac{1}{2\rho_0 c_0} \begin{bmatrix} \rho_0(c_0 - v_{0,x}) & 1 \\ \rho_0(c_0 + v_{0,x}) & -1 \end{bmatrix} = \begin{bmatrix} \ell_1 \\ \ell_2 \end{bmatrix}.$$

We have

$$\begin{aligned} \begin{bmatrix} p \\ \tilde{\sigma}_x \end{bmatrix} (0, t) &= \ell_1 \begin{bmatrix} p^+ \\ \tilde{\sigma}_x^+ \end{bmatrix} \mathbf{w}_1 + \ell_2 \begin{bmatrix} p^- \\ \tilde{\sigma}_x^- \end{bmatrix} \mathbf{w}_2 \\ &= \frac{\rho_0(c_0 - v_{0,x})p^+ + \tilde{\sigma}_x^+}{2\rho_0 c_0} \begin{bmatrix} 1 \\ \rho_0(c_0 + v_{0,x}) \end{bmatrix} + \frac{\rho_0(c_0 + v_{0,x})p^- - \tilde{\sigma}_x^-}{2\rho_0 c_0} \begin{bmatrix} 1 \\ \rho_0(v_{0,x} - c_0) \end{bmatrix}, \end{aligned}$$

therefore

$$\begin{aligned}\widehat{p} &= \frac{1}{2} (p^+ + p^-) - \frac{v_{0,x}}{2c_0} (p^+ - p^-) + \frac{1}{2\rho_0 c_0} (\widetilde{\sigma}_x^+ - \widetilde{\sigma}_x^-), \\ \widehat{\sigma}_x &= \frac{1}{2} (\widetilde{\sigma}_x^+ + \widetilde{\sigma}_x^-) + \frac{v_{0,x}}{2c_0} (\widetilde{\sigma}_x^+ - \widetilde{\sigma}_x^-) + \rho_0 \frac{c_0^2 - v_{0,x}^2}{2c_0} (p^+ - p^-).\end{aligned}$$

Finally, we can infer the form of the DG flux for a generic interface

$$\widehat{p} = \{\!\{p\}\!\} - \frac{\mathbf{v}_0 \cdot \mathbf{n}^-}{2c_0} \llbracket p \rrbracket + \frac{1}{2\rho_0 c_0} \llbracket \widetilde{\sigma} \rrbracket, \quad (26a)$$

$$\widehat{\sigma} \cdot \mathbf{n}^- = \{\!\{\widetilde{\sigma}\}\!\} \cdot \mathbf{n}^- + \frac{\mathbf{v}_0 \cdot \mathbf{n}^-}{2c_0} \llbracket \widetilde{\sigma} \rrbracket + \rho_0 \frac{c_0^2 - (\mathbf{v}_0 \cdot \mathbf{n}^-)^2}{2c_0} \llbracket p \rrbracket. \quad (26b)$$

Notice that we had to chose an orientation of the normal vector. Following our convention, we have chosen to use  $\mathbf{n}^-$  as it has the same orientation as  $\mathbf{v}_0$ .

Rewriting (26a)–(26b) in terms of  $\sigma$  instead of  $\widetilde{\sigma}$  leads to

$$\widehat{p}_h = \{\!\{p\}\!\} - \frac{\mathbf{v}_0 \cdot \mathbf{n}^-}{2c_0} \llbracket p \rrbracket + \frac{1}{2i\omega\rho_0 c_0} \llbracket \sigma \rrbracket \quad (27a)$$

$$\widehat{\sigma} \cdot \mathbf{n}^- = \{\!\{\widetilde{\sigma}\}\!\} \cdot \mathbf{n}^- + \frac{\mathbf{v}_0 \cdot \mathbf{n}^-}{2c_0} \llbracket \widetilde{\sigma} \rrbracket + i\omega\rho_0 \frac{c_0^2 - (\mathbf{v}_0 \cdot \mathbf{n}^-)^2}{2c_0} \llbracket p \rrbracket. \quad (27b)$$

Comparing (27a)–(27b) with (18a)–(18b), we see that

$$\tau^+ + \tau^- = 2\rho_0 c_0, \quad (28a)$$

$$\frac{\tau^+ \tau^-}{\tau^+ + \tau^-} = \rho_0 \frac{c_0^2 - (\mathbf{v}_0 \cdot \mathbf{n}^-)^2}{2c_0}. \quad (28b)$$

The system (28a)–(28b) leads to the following second-order equation

$$(\tau^+)^2 - 2\rho_0 c_0 \tau^+ + \rho_0^2 (c_0^2 - (\mathbf{v}_0 \cdot \mathbf{n}^-)^2) = 0,$$

and to the two following families for  $\tau^\pm$

$$\begin{aligned}\tau_1^+ &= \rho_0 (c_0 + \mathbf{v}_0 \cdot \mathbf{n}^-), & \tau_1^- &= \rho_0 (c_0 - \mathbf{v}_0 \cdot \mathbf{n}^-), \\ \tau_2^+ &= \rho_0 (c_0 - \mathbf{v}_0 \cdot \mathbf{n}^-), & \tau_2^- &= \rho_0 (c_0 + \mathbf{v}_0 \cdot \mathbf{n}^-).\end{aligned}$$

To discriminate between  $\tau_1^\pm$  and  $\tau_2^\pm$  we once again go back to (18a)–(18b) and we see that the solution must satisfy

$$\frac{\tau^+ - \tau^-}{2(\tau^+ + \tau^-)} = -\frac{\mathbf{v}_0 \cdot \mathbf{n}^-}{2c_0}.$$

We can therefore conclude that the upwind fluxes are obtained by using the  $\tau_2^\pm$  solution. We can make this choice independent of the orientation convention by noticing that  $\mathbf{n}^+ = -\mathbf{n}^-$ , leading to

$$\tau_2^\pm = \rho_0 (c_0 + \mathbf{v}_0 \cdot \mathbf{n}^\pm).$$

**Remark 3.3:** To keep polynomial fluxes on the interfaces, the background quantities will be approximated by their value at the center of the interface.

**Remark 3.4:** In the context of DG and HDG methods,  $\tau$  is usually chosen to be of the «order of unity» to ensure optimal convergence rate. In the error analysis of the method, we allow the dependency to the background coefficient to be hidden in the constants, so the choice (21) is actually possible.

### 3.3 Local solvability

We will now show the local solvability for the *total flux* formulation. Proving the well-posedness of the local problems is always very important when working with HDG methods. For the strongly coercive problems, for which HDG methods were initially designed, this property usually comes directly from the continuous problem. However for harmonic wave equations, which are only weakly coercive, things are more complicated : indeed solving the local problem amounts to solving a wave problem with Dirichlet boundary conditions. We therefore need to ensure that the local problem does not introduce resonance into the method, which is the case when the elements are small enough. In this section, we will prove that the static condensation process is well-defined when the mesh is fine enough.

Notice that in this case the proof relies on an *absorption technique* and is therefore very technical. Readers who are not familiar with HDG theory should probably begin with the proof for the *diffusive flux* formulation which is easier. It will be detailed in [SUBSECTION 4.2](#).

First, we need to show the

**Lemma 3.1:**

For  $p_h \in \mathcal{P}_k(K)$  with  $k > 0$ , the following inverse inequality holds

$$\|\nabla p_h \cdot \mathbf{n}\|_{\partial K} \lesssim \|\nabla p_h\|_{\partial K} \lesssim h_K^{-\frac{1}{2}} \|\nabla p_h\|_K.$$

**Proof:**

First, we notice that if  $p_h$  is constant the desired inequality reduces to  $0 \lesssim 0$ . We therefore only consider non-constant  $p_h$ .

Let  $\tilde{K}$  be the reference unit element. We consider the map  $F : \tilde{K} \rightarrow K$ . We use  $\tilde{\cdot}$  to denote quantities on the reference element instead of the more standard notation  $\hat{\cdot}$  to avoid confusion, as we already used  $\hat{\cdot}$  to denote the numerical fluxes.

Let  $\tilde{\gamma}^1 : H^2(\tilde{K}) \rightarrow L^2(\partial\tilde{K})$  be the normal derivative operator in the reference element. As  $\tilde{\gamma}^1$  is continuous, we have

$$\|\tilde{\gamma}^1(\tilde{p}_h)\|_{\partial\tilde{K}} \lesssim \|\tilde{p}_h\|_{2,\tilde{K}} \lesssim |\tilde{p}_h|_{1,\tilde{K}}$$

The second inequality holds as  $\tilde{p}_h \in \mathcal{P}_k(\tilde{K})$  which is a finite-dimensional vector-space on which all the norms are equivalent and  $p_h$  is not constant.

We now recall the following scaling inequalities, see [\[DS19, Eq \(1.6\), \(1.7\) & \(1.8\)\]](#)

$$\begin{aligned} |\tilde{p}_h|_{1,\tilde{K}} &\lesssim |\det \text{Jac}(F)|^{-\frac{1}{2}} \|\text{Jac}(F)\| |p_h|_{1,K} \lesssim |p_h|_{1,K}; \\ h_K^{\frac{1}{2}} \|\mu_h\|_{\partial K} &\lesssim \|\tilde{\mu}_h\|_{\partial\tilde{K}} \end{aligned}$$

Due to the regularity of the mesh, we have

$$h_K^{\frac{1}{2}} \|\nabla p_h \cdot \mathbf{n}\|_{\partial K} \lesssim \|\nabla p_h\|_K.$$

■



**Theorem 1** : *Local solvability for the total flux HDG method*

If  $\tau$  is chosen such that

$$\exists \tau_0 > 0, \quad \forall e \in \mathcal{E}(K), \quad 0 < \tau_0 \leq \tau + \mathbf{b}_0 \cdot \mathbf{n},$$

then there exists a constant  $\alpha_+ > 0$  such that the local problem is well-posed if  $\omega h_K < \alpha_+$ .

**Proof:** As (12a)–(12b) is a finite-dimensional problem, we only need to prove uniqueness of the solution. We therefore assume that  $\widehat{p}_h = 0$  and  $s = 0$ , and we need to show that the system

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, \mathbf{r}_h)_K - (p_h, \operatorname{div}(\mathbf{r}_h))_K + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_K = 0, \quad \forall \mathbf{r}_h \in \mathbf{V}_h(K) \quad (30a)$$

$$-\omega^2 (\rho_0 p_h, w_h)_K + (\operatorname{div}(\boldsymbol{\sigma}_h), w_h)_K + i\omega \langle \tau p_h, w_h \rangle_{\partial K} = 0, \quad \forall w_h \in W_h(K). \quad (30b)$$

has only one solution  $(\boldsymbol{\sigma}_h, p_h) = (\mathbf{0}, 0)$ .

We will prove the theorem by contradiction. We therefore assume that there is a non-zero solution  $(\boldsymbol{\sigma}_h, p_h)$  to (30a)–(30b).

*Step 1:* Energy-like identity.

We test (30a) with  $\mathbf{r}_h = \boldsymbol{\sigma}_h$  and conjugate the resulting equation, we then test (30b) with  $w_h = p_h$  and add the two resulting equations leading to

$$\|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K}^2 - \omega^2 \|p_h\|_{\rho_0, K}^2 - 2i\omega (\boldsymbol{\sigma}_h, p_h \mathbf{W}_0 \mathbf{b}_0)_K + i\omega \langle \tau p_h, p_h \rangle_{\partial K} = 0.$$

We then focus on the third term, as  $\mathbf{W}_0$  is real and symmetric we have

$$(\boldsymbol{\sigma}_h, p_h \mathbf{W}_0 \mathbf{b}_0)_K = (\mathbf{W}_0 \boldsymbol{\sigma}_h, p_h \mathbf{b}_0)_K.$$

Taking  $\mathbf{r}_h = p_h \{\mathbf{b}_0\}$ , where  $\{\mathbf{b}_0\}$  is the average of  $\mathbf{b}_0$  on  $K$ , in (30a), we have

$$(\mathbf{W}_0 \boldsymbol{\sigma}_h, p_h \{\mathbf{b}_0\})_K - (p_h, \{\mathbf{b}_0\} \cdot \nabla p_h)_K + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, p_h \{\mathbf{b}_0\})_K = 0,$$

leading to

$$-(\mathbf{W}_0 \boldsymbol{\sigma}_h, p_h \mathbf{b}_0)_K + (p_h, \mathbf{b}_0 \cdot \nabla p_h)_K - 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, p_h \mathbf{b}_0)_K = -\varepsilon,$$

where

$$\varepsilon := (\mathbf{W}_0 \boldsymbol{\sigma}_h, p_h (\mathbf{b}_0 - \{\mathbf{b}_0\}))_K - (p_h, (\mathbf{b}_0 - \{\mathbf{b}_0\}) \cdot \nabla p_h)_K + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, p_h (\mathbf{b}_0 - \{\mathbf{b}_0\}))_K.$$

We chose the notation  $\varepsilon$  to emphasize that this quantity is small, as it will be discussed in *Step 3*.

So we have

$$\begin{aligned} \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K}^2 - \omega^2 (\|p_h\|_{\rho_0, K}^2 + 4 \|p_h \mathbf{b}_0\|_{\mathbf{W}_0, K}^2) \\ - 2i\omega (p_h, \mathbf{b}_0 \cdot \nabla p_h)_K - 2i\omega \varepsilon + i\omega \langle \tau p_h, p_h \rangle_{\partial K} = 0 \end{aligned} \quad (31)$$

*Step 2:* Boundary condition

Taking the imaginary part (31) leads to

$$\begin{aligned} 2\omega \Re (p_h \mathbf{b}_0, \nabla p_h)_K + \omega \langle \tau p_h, p_h \rangle_{\partial K} + 2\omega \Re \varepsilon = 0, \\ \text{(by LEMMA 4.1)} \quad \langle (\tau + \mathbf{b}_0 \cdot \mathbf{n}) p_h, p_h \rangle_{\partial K} = -2\Re \varepsilon. \end{aligned}$$

As we have chosen  $\tau$  such that there is  $0 < \tau_0 \leq \tau + \mathbf{b}_0 \cdot \mathbf{n}$  where  $\tau_0$  does not depend on  $h_K$ , we have

$$\|p_h\|_{\partial K}^2 \lesssim |\varepsilon|. \quad (32)$$

As we do not have shown that  $p|_{\partial K} = 0$ , we cannot use Poincaré's inequality, however according to [EG04, Lemma B.63 & Example B.64] we have

$$\|p_h\|_K \leq \|p_h\|_{1,K} \leq C_K \|\nabla p_h\|_K + \frac{C_K}{\text{meas}(\partial K)} \|p_h\|_{\partial K},$$

where  $C_K$  is the Poincaré constant of  $K^2$ . Using standard scaling inequalities, we have

$$C_K \lesssim h_K \quad \text{and} \quad \frac{C_K}{\text{meas}(\partial K)} \lesssim h_K^{\frac{1}{2}}.$$

Using (32) to estimate the boundary term, this leads to

$$\|p_h\|_K \lesssim h_K \|\nabla p_h\|_K + h_K^{1/2} |\varepsilon|^{1/2}. \quad (33)$$

*Step 3:* Estimating  $|\varepsilon|$  and  $\|p_h\|_K$ . We have

$$\begin{aligned} |\varepsilon| &\lesssim h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K} \|p_h\|_K + h_K \|p_h\|_K \|\nabla p_h\|_K + \omega h_K \|p_h\|_K^2 \\ \text{(by Young)} &\lesssim h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + (1 + \omega h_K) \|p_h\|_K^2 + h_K^2 \|\nabla p_h\|_K^2 \\ &\lesssim h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + (1 + \omega h_K) \left( h_K^2 \|\nabla p_h\|_K^2 + h_K |\varepsilon| \right) + h_K^2 \|\nabla p_h\|_K^2 \end{aligned}$$

If  $h_K$  is small enough, the term  $h_K |\varepsilon|$  in the right-hand side can be absorbed by the left-hand side leading to

$$|\varepsilon| \lesssim h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + (1 + \omega h_K) h_K^2 \|\nabla p_h\|_K^2.$$

Assuming that  $h_K$  is small enough, we can overestimate  $\omega h_K \lesssim 1$ , leading to

$$|\varepsilon| \lesssim h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + h_K^2 \|\nabla p_h\|_K^2. \quad (34)$$

Together with the generalized Poincaré's inequality (33) we have

$$\|p_h\|_K^2 \lesssim h_K^2 \|\nabla p_h\|_K^2 + h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2, \quad (35)$$

and using that  $\sqrt{a^2 + b^2} \leq |a| + |b|$  we also have

$$\|p_h\|_K \lesssim h_K \|\nabla p_h\|_K + h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}.$$

*Step 4:* Estimating  $\|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}$

Taking the real part of the Garding's identity (31), we have

$$\begin{aligned} \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 &\lesssim \omega^2 \|p_h\|_K^2 + \omega \|p_h\|_K \|\nabla p_h\|_K + \omega |\varepsilon| \\ \text{(by (35))} &\lesssim \omega^2 \left( h_K^2 \|\nabla p_h\|_K^2 + h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 \right) + \omega \left( h_K \|\nabla p_h\|_K + h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K} \right) \|\nabla p_h\|_K + \omega |\varepsilon| \\ \text{(by Young)} &\lesssim \left( \omega^2 h_K^2 + \omega h_K \right) \|\nabla p_h\|_K^2 + \omega^2 h_K^2 \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + \omega h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 + \omega |\varepsilon| \\ \text{(by (34))} &\lesssim \left( \omega^2 h_K^2 + \omega h_K \right) \|\nabla p_h\|_K^2 + \left( \omega^2 h_K^2 + \omega h_K \right) \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 \\ &\lesssim \omega h_K \|\nabla p_h\|_K^2 + \omega h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{w}_0,K}^2 \end{aligned}$$

<sup>2</sup>The constant used by the authors of [EG04] is the inverse of the usual Poincaré constant.

We obtained the last line by assuming that  $\omega^2 h_K^2 \lesssim \omega h_K$  which is true if  $h_K$  is small enough. If  $h_K$  is small enough the term  $\omega h_K \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K}^2$  in the right-hand side can be absorbed by the left-hand side, leading to

$$\|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K}^2 \lesssim \omega h_K \|\nabla p_h\|_K^2. \quad (36)$$

*Step 5: Estimating  $\|\nabla p_h\|_K$*

Taking  $\mathbf{r}_h = \nabla p_h$  in (30a) and reverting the integration by parts, we have

$$\begin{aligned} \|\nabla p_h\|_K^2 &= |(\mathbf{W}_0 \boldsymbol{\sigma}_h, \nabla p_h)_K + 2i\omega (p_h \mathbf{W}_0 \mathbf{b}_0, \nabla p_h)_K + \langle p_h, \nabla p_h \cdot \mathbf{n} \rangle_{\partial K}| \\ &\lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} \|\nabla p_h\|_K + \omega \|p_h\|_K \|\nabla p_h\|_K + \|p_h\|_{\partial K} \|\nabla p_h\|_{\partial K} \\ (\text{by LEMMA 3.1}) \quad &\lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} \|\nabla p_h\|_K + \omega \|p_h\|_K \|\nabla p_h\|_K + h_K^{-1/2} \|p_h\|_{\partial K} \|\nabla p_h\|_K \end{aligned}$$

So we have

$$\begin{aligned} \|\nabla p_h\|_K &\lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} + \omega \|p_h\|_K + h_K^{-1/2} \|p_h\|_{\partial K} \\ &\lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} + \omega \|p_h\|_K + (1 + \omega h_K) h_K \|\nabla p_h\|_K \\ &\lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} + \omega h_K \|\nabla p_h\|_K + \omega h_K^{3/2} \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} + (1 + \omega h_K) h_K \|\nabla p_h\|_K \end{aligned}$$

Finally we have

$$\|\nabla p_h\|_K \lesssim \|\boldsymbol{\sigma}_h\|_{\mathbf{W}_0, K} \lesssim \|\boldsymbol{\sigma}_h\|_K. \quad (37)$$

*Step 6: Contradiction*

Combining (36) and (37), we have

$$\|\boldsymbol{\sigma}_h\|_K^2 \lesssim \omega h_K \|\boldsymbol{\sigma}_h\|_K^2$$

as we assumed that  $\boldsymbol{\sigma}_h \neq 0$ , we can divide by  $\|\boldsymbol{\sigma}_h\|_K$ , leading to

$$1 \lesssim \omega h_K,$$

which does not hold if  $\omega h_K$  is small enough.

This is the desired contradiction, and we can therefore conclude that  $(\boldsymbol{\sigma}_h, p_h) = (\mathbf{0}, 0)$  is the only solution of the system (30a)–(30b).  $\blacksquare$

### 3.4 Error analysis

The error analysis can be carried out by following the projection analysis for the Helmholtz equation given in [DS19, Sec. 3.5.1 & 3.5.2] with some minor changes.

This error analysis relies on the tailored HDG projection that fits the structure of the numerical trace. This projection  $(\boldsymbol{\Pi}, \Pi)$ , with

$$(\boldsymbol{\Pi}, \Pi) : \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O}) \longrightarrow \mathbf{V}_h \times W_h := \mathcal{P}_k(\mathcal{T}_h) \times \mathcal{P}_k(\mathcal{T}_h)$$

is defined by the following equations

$$\begin{aligned} (\boldsymbol{\Pi}\boldsymbol{\sigma}, \mathbf{r}_h)_K &= (\boldsymbol{\sigma}, \mathbf{r}_h)_K, & \forall \mathbf{r}_h \in \mathcal{P}_{k-1}(K), \\ (\Pi p, w_h)_K &= (p, w_h)_K, & \forall w_h \in \mathcal{P}_{k-1}(K), \\ \langle \boldsymbol{\Pi}\boldsymbol{\sigma} \cdot \mathbf{n} + i\omega\tau\Pi p, \mu_h \rangle_{\partial K} &= \langle \boldsymbol{\sigma} \cdot \mathbf{n} + i\omega\tau p, \mu_h \rangle_{\partial K}, & \forall \mu_h \in \mathcal{R}_k(\partial K). \end{aligned}$$

Notice that denoting the image of  $(\boldsymbol{\sigma}, p)$  under  $(\boldsymbol{\Pi}, \Pi)$  by  $(\boldsymbol{\Pi}\boldsymbol{\sigma}, \Pi p)$  is a slight abuse of notation as both components depend on  $\boldsymbol{\sigma}$  and  $p$ . However it is very convenient and often found in the literature.

We define the following error quantities

$$\boldsymbol{\delta}_h^\sigma := \mathbf{\Pi}\boldsymbol{\sigma} - \boldsymbol{\sigma} \ ; \ \delta_h^p := \mathbf{\Pi}p - p \ ; \ \widehat{\delta}_h^p := p - P_{MP}$$

and

$$\boldsymbol{\varepsilon}_h^\sigma := \mathbf{\Pi}\boldsymbol{\sigma} - \boldsymbol{\sigma}_h \in \mathbf{V}_h \ ; \ \varepsilon_h^p := \mathbf{\Pi}p - p_h \in W_h \ ; \ \widehat{\varepsilon}_h^p := P_{MP} - \widehat{p}_h \in M_h.$$

We will split the errors as

$$\begin{aligned} \|\boldsymbol{\sigma}_h - \boldsymbol{\sigma}\|_{\mathbf{W}_0, \mathcal{T}_h} &\leq \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h} + \|\boldsymbol{\delta}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}, \\ \|p_h - p\|_{\rho_0, \mathcal{T}_h} &\leq \|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h} + \|\delta_h^p\|_{\rho_0, \mathcal{T}_h}, \end{aligned}$$

Notice that the following estimates hold

$$\begin{aligned} \|\delta_h^p\|_K &\lesssim h_K^{k+1} \left( |p|_{k+1, K} + \tau_{\max}^{-1} |\operatorname{div} \boldsymbol{\sigma}|_{k, K} \right), \\ \|\boldsymbol{\delta}_h^\sigma\|_K &\lesssim h_K^{k+1} \left( |\boldsymbol{\sigma}|_{k+1, K} + \tau^* |p|_{k+1, K} \right), \end{aligned}$$

for some constants  $\tau_{\max}$  and  $\tau^*$ . So we only need to prove estimates for  $\|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}$ ,  $\|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h}$ .

The error analysis can be summarized as follows

1. we derive an estimate for  $\|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}$  using the energy-like inequality,
2. we use a *dual problem* to estimate  $\|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h}$  in terms of  $\|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}$ ,
3. those estimates are combined through a *bootstrapping process*.

This analysis is therefore strongly related to the *Aubin-Nitsche method* and only works for regular solutions.

We will now give the main changes needed to adapt the error analysis from [DS19].

The error equations (3.30) become

$$\begin{aligned} (\mathbf{W}_0 \boldsymbol{\varepsilon}_h^\sigma, \mathbf{r}_h)_{\mathcal{T}_h} - (\varepsilon_h^p, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + 2i\omega (\varepsilon_h^p \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} + \langle \widehat{\varepsilon}_h^p, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} &= (\mathbf{W}_0 \boldsymbol{\delta}_h^\sigma, \mathbf{r}_h)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{W}_0 \mathbf{b}_0, \mathbf{r}_h)_{\mathcal{T}_h} \\ -\omega^2 (\rho_0 \varepsilon_h^p, w_h)_{\mathcal{T}_h} + (\operatorname{div}(\boldsymbol{\varepsilon}_h^\sigma), w_h)_{\mathcal{T}_h} + i\omega \langle \tau(\varepsilon_h^p - \widehat{\varepsilon}_h^p), w_h \rangle_{\partial \mathcal{T}_h} &= -\omega^2 (\rho_0 \delta_h^p, w_h)_{\mathcal{T}_h} \\ -\langle \boldsymbol{\varepsilon}_h^\sigma \cdot \mathbf{n} + i\omega \tau(\varepsilon_h^p - \widehat{\varepsilon}_h^p), \boldsymbol{\mu}_h \rangle_{\partial \mathcal{T}_h} &= 0. \end{aligned}$$

The energy-like identity of Prop. 3.7 becomes

$$\begin{aligned} \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}^2 - \omega^2 \|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h}^2 - 2i\omega (\boldsymbol{\varepsilon}_h^\sigma, \varepsilon_h^p \mathbf{W}_0 \mathbf{b}_0)_{\mathcal{T}_h} + i\omega \langle \tau(\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} &= (\boldsymbol{\varepsilon}_h^\sigma, \mathbf{W}_0 \boldsymbol{\delta}_h^\sigma)_{\mathcal{T}_h} \\ &\quad - 2i\omega (\boldsymbol{\varepsilon}_h^\sigma, \delta_h^p \mathbf{W}_0 \mathbf{b}_0)_{\mathcal{T}_h} \\ &\quad - \omega^2 (\rho_0 \delta_h^p, \varepsilon_h^p)_{\mathcal{T}_h}, \end{aligned}$$

leading to the following estimate

$$\begin{aligned} \left| \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h}^2 + i\omega \langle \tau(\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} \right| &\lesssim \omega^2 \|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h}^2 + \omega \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h} \|\delta_h^p\|_{\rho_0, \mathcal{T}_h} + \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h} \|\boldsymbol{\delta}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h} \\ &\quad + \omega \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0, \mathcal{T}_h} \|\delta_h^p\|_{\rho_0, \mathcal{T}_h} + \omega^2 \|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h} \|\delta_h^p\|_{\rho_0, \mathcal{T}_h}. \end{aligned}$$

The adjoint problem (3.31) becomes

$$\begin{aligned} \mathbf{W}_0 \boldsymbol{\xi} - \nabla \theta - 2i\omega \theta \mathbf{W}_0 \mathbf{b}_0 &= 0 \\ -\rho_0 \omega^2 \theta - \operatorname{div}(\boldsymbol{\xi}) &= \varepsilon_h^p. \end{aligned}$$

We state the *elliptic regularity* assumption which is a key ingredient in the error analysis

$$\|\theta\|_{2,\mathcal{O}} + \|\boldsymbol{\xi}\|_{1,\mathcal{O}} \leq C_{\text{reg}} \|\varepsilon_h^p\|_{\mathcal{O}}$$

The identity of Prop. 3.8 becomes

$$\begin{aligned} \|\varepsilon_h^p\|_{\rho_0,\mathcal{T}_h}^2 &= \omega^2 (\rho_0(\Pi\theta - \theta), \varepsilon_h^p - \delta_h^p)_{\mathcal{T}_h} - \omega^2 (\rho_0\theta - \{\rho_0\theta\}, \delta_h^p)_{\mathcal{T}_h} \\ &\quad - (\mathbf{W}_0(\Pi\boldsymbol{\xi} - \boldsymbol{\xi}), \boldsymbol{\varepsilon}_h^\sigma - \boldsymbol{\delta}_h^\sigma)_{\mathcal{T}_h} + (\mathbf{W}_0\boldsymbol{\xi} - \{\mathbf{W}_0\boldsymbol{\xi}\}, \boldsymbol{\delta}_h^\sigma)_{\mathcal{T}_h} \\ &\quad - 2i\omega (\theta\mathbf{b}_0, \mathbf{W}_0\varepsilon_h^\sigma)_{\mathcal{T}_h} + 2i\omega (\mathbf{W}_0\boldsymbol{\xi}, (\varepsilon_h^p - \delta_h^p)\mathbf{b}_0)_{\mathcal{T}_h} \\ &\quad + 2i\omega (\mathbf{W}_0(\Pi\boldsymbol{\xi} - \boldsymbol{\xi}), (\varepsilon_h^p - \delta_h^p)\mathbf{b}_0)_{\mathcal{T}_h}, \end{aligned}$$

leading to the following estimate

$$\begin{aligned} \|\varepsilon_h^p\|_{\rho_0,\mathcal{T}_h} &\lesssim (\omega^3 + \omega^2 + \omega) h \left( \|\varepsilon_h^p\|_{\rho_0,\mathcal{T}_h} + \|\delta_h^p\|_{\rho_0,\mathcal{T}_h} \right) + (1 + \omega)h \left( \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0,\mathcal{T}_h} + \|\boldsymbol{\delta}_h^\sigma\|_{\mathbf{W}_0,\mathcal{T}_h} \right) \\ &\quad + \|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathbf{W}_0,K} + \omega \|\varepsilon_h^p\|_{\rho_0,K} + \omega \|\delta_h^p\|_{\rho_0,K}. \end{aligned}$$

Notice that in contrast to the HDG method for the standard Helmholtz equation, both  $p_h$  and  $\boldsymbol{\sigma}_h$  have the *same* convergence rate.

It is now straightforward to follow the bootstrapping argument of Sec. 3.5.2 to obtain the

**Theorem 2** : *Convergence of the HDG method with total flux*

Assuming that  $\omega h$  is small enough and under the elliptic regularity assumption, we have

$$\|p_h - \Pi p\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+1}) \quad ; \quad \|\boldsymbol{\sigma}_h - \Pi\boldsymbol{\sigma}\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+1}).$$

**Remark 3.5:** Some HDG methods are known to achieve *super-convergence*, ie taking  $p_h \in \mathcal{P}_k$  leads to the following error estimate

$$\|\Pi p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+2}).$$

Superconvergence is an attractive property for a numerical scheme, indeed using a postprocessing scheme it is possible to use the solution  $(\boldsymbol{\sigma}_h, p_h)$  to construct a new approximation  $\widehat{p}_h$  which converges with order  $\mathcal{O}(h^{k+2})$ , see [Ste91], [CGS10, Sec. 5] for more details.

Here we were only able to prove *optimal convergence* and the use of post-processing schemes to improve the convergence rate is therefore not possible.

### 3.5 Global solvability

The analysis that we have carried out in the previous subsection works for any solution  $(\boldsymbol{\sigma}_h, p_h, \widehat{p}_h)$  of the discrete system (15a)–(15b)–(15c) provided that such solution exists. We already discussed the well-posedness of the local problems in THEOREM 1, but we have not yet proved that the global problem (16) for  $\widehat{p}_h$  was well-posed.

To do that we can either directly show the well-posedness of the global problem (16). Or we can choose take advantage of the error estimates of THEOREM 2 as we will describe below<sup>3</sup>.

**Resonant frequencies:** We recall that the convected Helmholtz equation is a problem of Fredholm type. It is therefore uniquely solvable except on a set of *resonant frequencies*. For those frequencies, there exists non-zero solutions to the homogenous equation and unique solvability cannot be guaranteed.

<sup>3</sup>In [DS19] this idea is attributed to B. Cockburn.

**Main result:** We can now state and prove the main result of this section.

**Theorem 3 :** *Global solvability*

Under the assumptions of [THEOREM 1](#) and [THEOREM 2](#) and if  $\omega$  is not a resonant frequency of the convected Helmholtz equation (1) then the global problem is well-posed, ie  $\widehat{p}_h$  is uniquely defined by (16).

**Proof:** First we recall that (15a)–(15b)–(15c), or equivalently (16), is a square system of linear equations, we therefore only need to show the uniqueness of the solution of the homogenous system (when  $g_N = g_D = s = 0$ ).

Assuming that  $\omega$  is not a resonant frequency of (1), the exact solution is  $p = 0$  and  $\boldsymbol{\sigma} = \mathbf{0}$ , and therefore

$$\|p\|_{s,\mathcal{O}} = 0 \quad \text{and} \quad \|\boldsymbol{\sigma}\|_{t,\mathcal{O}} = 0$$

and

$$\varepsilon_h^p = -p_h \quad ; \quad \boldsymbol{\varepsilon}_h^\sigma = -\boldsymbol{\sigma}_h \quad ; \quad \widehat{\varepsilon}_h^p = -\widehat{p}_h.$$

The aim of the error analysis was to prove the following inequalities when  $h$  is small enough :

$$\|\varepsilon_h^p\|_{\mathcal{T}_h} \lesssim \|p\|_{s,\mathcal{O}} + \|\boldsymbol{\sigma}_h\|_{t,\mathcal{O}} = 0 \tag{38a}$$

$$\|\boldsymbol{\varepsilon}_h^\sigma\|_{\mathcal{T}_h} \lesssim \|p\|_{s,\mathcal{O}} + \|\boldsymbol{\sigma}_h\|_{t,\mathcal{O}} = 0 \tag{38b}$$

Notice that we have hidden the powers of  $h$  in  $\lesssim$  as they do not play an important part here. Therefore using (38a) and (38b) we have shown that

$$p_h \equiv 0 \quad \text{and} \quad \boldsymbol{\sigma}_h \equiv \mathbf{0}$$

when  $h$  is small enough.

For all  $K \in \mathcal{T}_h$ , we can now rewrite (12a) as

$$\langle \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial K} = 0, \quad \forall \mathbf{r}_h \in \mathbf{V}_h(K),$$

which leads to

$$\widehat{p}_h \equiv 0. \quad \blacksquare$$

## 4 HDG(+) methods for the diffusive flux formulation

In this section, we will construct HDG methods based on the *diffusive flux* formulation, where  $\mathbf{q}$  is used instead of  $\boldsymbol{\sigma}$ . We will mostly describe the HDG+ method where different polynomial degrees are used for the different unknowns, as it the most important novelty of this paper. Adaptation of the formulation construction and theoretical results to a more standard HDG method with the same polynomial interpolation for all the unknowns is straightforward. The main differences between the HDG and HDG+ methods are stated but the details are left out. The *global solvability* will not be included in this section as the adaptation of the result from the previous section is immediate.

## 4.1 Construction of the method

We recall that on element  $K \in \mathcal{T}_h$ , the weak formulation (7a)–(7b) reads : Seek  $(\mathbf{q}, p) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$  such that

$$\int_K \mathbf{W}_0 \mathbf{q} \cdot \mathbf{r}^* d\mathbf{x} - \int_K p \operatorname{div}(\mathbf{r}^*) d\mathbf{x} + \int_{\partial K} p \mathbf{r}^* \cdot \mathbf{n} d\sigma = 0, \quad (39a)$$

$$-\omega^2 \int_K \rho_0 p w^* d\mathbf{x} - 2i\omega \int_K \mathbf{b}_0 \cdot \nabla p w^* d\mathbf{x} + \int_K \operatorname{div}(\mathbf{q}) w^* d\mathbf{x} = \int_K s w^* d\mathbf{x}, \quad (39b)$$

for all  $(\mathbf{r}, w) \in \mathbf{H}_{\text{div}}(\mathcal{O}) \times H^1(\mathcal{O})$ .

**Choice of approximation spaces:** For the HDG+ method, the choice of approximation spaces is different from the choice made for the previous HDG method. We consider the following *local approximation spaces*

$$\begin{aligned} \mathbf{V}_h(K) &= \mathcal{P}_k(K), & \text{for the flux } \mathbf{q}_h, \\ W_h(K) &= \mathcal{P}_{k+1}(K), & \text{for the pressure } p_h, \end{aligned}$$

where  $k \geq 2$  is the degree of the method. The use of a higher polynomial degree for  $p_h$  is the distinctive feature of the HDG+ method.

**Introduction of the hybrid unknown:** As we did before, we introduce the *numerical trace*  $\widehat{p}_h$  which approximates  $p$  on the skeleton  $\mathcal{E}_h$  of the mesh. As before the boundary integral in (39a) will be discretized as

$$\int_{\partial K} p \mathbf{r}^* \cdot \mathbf{n} d\sigma \quad \text{becomes} \quad \int_{\partial K} \widehat{p}_h \mathbf{r}_h^* \cdot \mathbf{n} d\sigma.$$

For the HDG+ method, we use the following approximation space for  $\widehat{p}_h$

$$M_h(e) = \mathcal{P}_k(e), \quad \forall e \in \mathcal{E}(K).$$

With this choice,  $p_h$  and  $\widehat{p}_h$  do not have the same polynomial degree and we therefore have two approximations of  $p$  with different polynomial degrees on the skeleton of the mesh. We therefore need to change the *penalization term* to

$$\tau(P_M p_h - \widehat{p}_h), \quad (40)$$

where  $P_M$  is the  $L^2$ -orthogonal projection onto  $M_h$ . This is called the *reduced stabilization* and it was introduced in [Leh10]. It allows to get convergence rate of  $k + 2$  for  $p_h$  for the cost of a method of degree  $k$ . A large penalization parameter  $\tau \sim h_K^{-1}$  is needed to obtain optimal convergence as it will be detailed in [SUBSECTION 4.3](#).

**Local problem:** We approximate the weak formulation (39a)–(39b) on an element  $K \in \mathcal{T}_h$  leading to the so-called *local problem* : seek  $(\mathbf{q}_h, p_h) \in \mathbf{V}_h(K) \times W_h(K)$  such that

$$(\mathbf{W}_0 \mathbf{q}_h, \mathbf{r}_h)_K - (p_h, \operatorname{div}(\mathbf{r}_h))_K + \langle \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial K} = 0, \quad (41a)$$

$$\begin{aligned} -\omega^2 (\rho_0 p_h, w_h)_K - 2i\omega (\mathbf{b}_0 \cdot \nabla p_h, w_h)_K + (\operatorname{div}(\mathbf{q}_h), w_h)_K \\ + 2i\omega \langle \tau(P_M p_h - \widehat{p}_h) - \tau_{\text{upw}}(p_h - \widehat{p}_h), w_h \rangle_{\partial K} = (s, w_h)_K, \end{aligned} \quad (41b)$$

for all  $(\mathbf{r}_h, w_h) \in \mathbf{V}_h(K) \times W_h(K)$ .

Following [QS16a], we have introduced a second penalization parameter  $\tau_{\text{upw}}$  defined by

$$\tau_{\text{upw}} := \max(\mathbf{b}_0 \cdot \mathbf{n}, 0).$$

To understand why this second parameter is required, we recall that in HDG methods the penalization serves two purposes:

1. it enforces the Dirichlet boundary condition for the local problems,
2. it controls the stability of the method.

Here as  $\mathbf{q}_h$  does not take the convection into account, the penalization term (40) with  $\tau$  only stabilizes the diffusion. We therefore need to add a second penalization to stabilize the convection. We denoted it  $\tau_{\text{upw}}$  as it leads to an upwinding behavior that will be detailed in the next paragraph.

**Transmission condition:** Following the previous example, we introduce the *following numerical flux*

$$\widehat{\mathbf{q}}_h \cdot \mathbf{n} := \mathbf{q}_h \cdot \mathbf{n} + 2i\omega\tau(P_M p_h - \widehat{p}_h), \quad (42)$$

where  $\tau = \mathcal{O}(h_K^{-1})$ . As discussed before, we need to require the normal continuity of the *total flux* on the interface between two elements, and the quantity  $\widehat{\mathbf{q}}_h \cdot \mathbf{n}$  only takes the diffusion into account. To deal with convection we add a second numerical flux

$$2i\omega p_h \widehat{\mathbf{b}}_0 \cdot \mathbf{n} := 2i\omega(\mathbf{b}_0 \cdot \mathbf{n})\widehat{p}_h + 2i\omega\tau_{\text{upw}}(p_h - \widehat{p}_h). \quad (43)$$

It is important to notice that this flux has an upwind behavior. Let  $e = \partial K_+ \cap \partial K_-$  be an interior edge with  $\mathbf{b}_0 \cdot \mathbf{n}_- > 0$  on  $\partial K_-$ . We have

$$\text{On } \partial K_-: \quad \tau_{\text{upw}} := \max(\mathbf{b}_0 \cdot \mathbf{n}, 0) = \mathbf{b}_0 \cdot \mathbf{n}, \quad \text{so} \quad 2i\omega p_h \widehat{\mathbf{b}}_0 \cdot \mathbf{n} = 2i\omega(\mathbf{b}_0 \cdot \mathbf{n})p_h \quad (44a)$$

$$\text{On } \partial K_+: \quad \tau_{\text{upw}} := \max(\mathbf{b}_0 \cdot \mathbf{n}, 0) = 0, \quad \text{so} \quad 2i\omega p_h \widehat{\mathbf{b}}_0 \cdot \mathbf{n} = 2i\omega(\mathbf{b}_0 \cdot \mathbf{n})\widehat{p}_h \quad (44b)$$

So on the outflow boundary we use the interior value  $p_h$ , whereas on the inflow boundary we use the trace value  $\widehat{p}_h$ .

Finally we write the *transmission condition* as

$$\left\langle (\widehat{\mathbf{q}}_h - 2i\omega p_h \widehat{\mathbf{b}}_0) \cdot \mathbf{n}, \mu_h \right\rangle_{\partial \mathcal{T}_h \setminus \Gamma_D} + \langle \widehat{p}_h - g_D, \mu_h \rangle_{\Gamma_D} = \langle g_N, \mu_h \rangle_{\Gamma_N}. \quad (45)$$

This formulation enforces normal continuity of the total flux between the elements and the boundary conditions on  $\Gamma_D$  and  $\Gamma_N$ .

**Remark 4.1:** To ensure the well-posedness of the *local problems*, the second penalization must be

$$\tau_{\text{upw}}(p_h - \widehat{p}_h),$$

and not

$$\tau_{\text{upw}}(P_M p_h - \widehat{p}_h),$$

see [THEOREM 4](#).

**Remark 4.2:** We would like to point out the main theoretical difficulty of this method : when the background flow is not constant, the second flux (43) leads to non-polynomial terms on the skeleton. This is usually avoided as much as possible in HDG methods.

**Adaptation to a standard HDG method:** With this formulation, it is also possible to consider a standard HDG method by using the same polynomial degree for  $p_h$  and  $\mathbf{q}_h$ , *ie.* by using the following local approximation spaces

$$\begin{aligned} \mathbf{V}_h(K) &= \mathcal{P}_k(K), & \text{for the flux } \mathbf{q}_h, \\ W_h(K) &= \mathcal{P}_k(K), & \text{for the pressure } p_h. \end{aligned}$$



In this case, as  $W_h$  and  $M_h$  have the same polynomial degree, the projection term becomes simpler, indeed

$$P_M p_h = p_h.$$

For this formulation, we do not require a large penalization parameter anymore, and we only need  $\tau = \mathcal{O}(1)$ .

#### 4.1.1 Compact formulation of the methods:

HDG methods are usually stated in a compact way that can be obtained by summing the local problems (41a)–(41b) over the mesh elements and by adding the transmission condition (45). This formulation reads : seek  $(\mathbf{q}_h, p_h, \widehat{p}_h) \in \mathbf{V}_h \times W_h \times M_h$ , such that

$$(\mathbf{W}_0 \mathbf{q}_h, \mathbf{r}_h)_{\mathcal{T}_h} - (p_h, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + \langle \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = 0, \quad (46a)$$

$$-\omega^2 (\rho_0 p_h, w_h)_{\mathcal{T}_h} - 2i\omega (\mathbf{b}_0 \cdot \nabla p_h, w_h)_{\mathcal{T}_h} + (\operatorname{div}(\mathbf{q}_h), w_h)_{\mathcal{T}_h} \quad (46b)$$

$$+ 2i\omega \langle \tau(P_M p_h - \widehat{p}_h) - \tau_{\text{upw}}(p_h - \widehat{p}_h), w_h \rangle_{\partial \mathcal{T}_h} = (s, w_h)_{\mathcal{T}_h}$$

$$\langle (\widehat{\mathbf{q}}_h - 2i\omega p_h \widehat{\mathbf{b}}_0) \cdot \mathbf{n}, \mu \rangle_{\partial \mathcal{T}_h \setminus \Gamma_D} + \langle \widehat{p}_h - g_D, \mu_h \rangle_{\Gamma_D} = \langle g_N, \mu_h \rangle_{\Gamma_N}, \quad (46c)$$

for all  $(\mathbf{r}_h, w_h, \mu_h) \in \mathbf{V}_h \times W_h \times M_h$ .

#### 4.1.2 Condensed variational formulation

The compact formulation (46a)–(46b)–(46c) cannot be used to efficiently implement the HDG method, indeed with this formulation it is not clear how the *global problem* for  $\widehat{p}_h$  only can be obtained. To describe this process, we will now write a *condensed* variational formulation for  $\widehat{p}_h$  only.

We introduce the so-called *local solvers*

$$\begin{aligned} \mathbf{P}^K &: (\widehat{p}_h, s) \mapsto p_h^K, \\ \mathbf{Q}^K &: (\widehat{p}_h, s) \mapsto \mathbf{q}_h^K, \\ \widehat{\mathbf{Q}}^K &: (\widehat{p}_h, s) \mapsto \widehat{\mathbf{q}}_h^K, \end{aligned}$$

where  $(\mathbf{q}_h^K, p_h^K)$  is the solution of (41a)–(41b) and  $\widehat{\mathbf{q}}_h^K$  is defined by (42).

We can therefore rewrite the transmission condition (46c) as

$$a_h(\widehat{p}_h, \mu) = \ell_h(\mu), \quad (47)$$

where

$$\begin{aligned} a_h(\widehat{p}_h, \mu_h) &:= \langle \mathbf{Q}^K(\widehat{p}_h, s) \cdot \mathbf{n} + 2i\omega \tau (P_M \mathbf{P}^K(\widehat{p}_h, s) - \widehat{p}_h) + 2i\omega \tau_{\text{upw}}(\mathbf{P}^K(\widehat{p}_h, s) - \widehat{p}_h), \mu_h \rangle_{\partial \mathcal{T}_h \setminus \Gamma_D} \\ &\quad + \langle \widehat{p}_h, \mu_h \rangle_{\Gamma_D}, \\ \ell_h(\mu_h) &:= \langle g_N, \mu_h \rangle_{\Gamma_N} + \langle g_D, \mu_h \rangle_{\Gamma_D}. \end{aligned}$$

Equation (47) is the so-called *global problem* and is the main equation of the HDG method. From a computational point of view, we proceed as described in [ALGORITHM 2](#).

---

**Algorithm 2:** Solving HDG+
 

---

```

1 for  $K \in \mathcal{T}_h$  do
2   | Construct the local solvers  $\mathbf{P}^K, \mathbf{Q}^K, \widehat{\mathbf{Q}}^K$ 
3   | Add local contribution to the global problem (16)
4 Solve (16) for  $\widehat{p}_h$  // This is the main step
5 for  $K \in \mathcal{T}_h$  do
6   | Reconstruct the local unknowns  $p_h^K = \mathbf{P}^k(\widehat{p}_h, s)$  and  $\mathbf{q}_h^K = \mathbf{Q}^K(\widehat{p}_h, s)$ 
    
```

---

## 4.2 Local solvability

It is worth remembering that HDG methods were originally developed for elliptic problems and that harmonic wave equations are only coercive. It is well-known that solving those equations with Dirichlet boundary conditions<sup>4</sup> leads to numerical pollution due to the resonance phenomenon. In this section we will show that the static condensation process is well-defined when the mesh is fine enough, *ie.* the local problem does not produce resonance.

Before actually showing the local solvability, we need to prove the following lemma.

**Lemma 4.1:**

If  $p \in H^1(K)$  and  $\mathbf{b}_0 \in \mathbf{L}^\infty(K) \cap \mathcal{C}(\mathcal{O})$ , where  $\mathcal{C}(\mathcal{O})$  is the space of vector functions continuous in the domain  $\mathcal{O}$ , then the following identity holds

$$\Re (p\mathbf{b}_0, \nabla p)_K = \frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n})p, p \rangle_{\partial K}.$$

**Proof:** We use an integration by parts to obtain a relationship between  $(p\mathbf{b}_0, \nabla p)_K$  and its complex conjugate :

$$\begin{aligned}
 2\Re (p\mathbf{b}_0, \nabla p)_K &= (p\mathbf{b}_0, \nabla p)_K + (p\mathbf{b}_0, \nabla p)_K^* \\
 &= (p\mathbf{b}_0, \nabla p)_K + (\nabla p, p\mathbf{b}_0)_K \\
 &= -(\operatorname{div}(p\mathbf{b}_0), p)_K + \langle (\mathbf{b}_0 \cdot \mathbf{n})p, p \rangle_{\partial K} + (\nabla p, p\mathbf{b}_0)_K \\
 (\operatorname{div}(\mathbf{b}_0) = 0) &= -(\nabla p, p\mathbf{b}_0)_K + \langle (\mathbf{b}_0 \cdot \mathbf{n})p, p \rangle_{\partial K} + (\nabla p, p\mathbf{b}_0)_K \\
 &= \langle (\mathbf{b}_0 \cdot \mathbf{n})p, p \rangle_{\partial K}.
 \end{aligned}$$

■

We can now state and prove the main result of this section.

**Theorem 4 :** *Local solvability for the HDG+ methods*


---

If

$$\forall e \in \mathcal{E}(K), \quad \tau|_e < 0 \tag{48}$$

and if

$$\omega_{h,K} < \frac{-C_{\mathbf{w}_0,K} \|\mathbf{b}_0\|_{L^\infty(K)} + \left( C_{\mathbf{w}_0,K}^2 \|\mathbf{b}_0\|_{L^\infty(K)}^2 + \|\rho_0\|_{L^\infty(K)} \right)^{\frac{1}{2}}}{C_{\mathbf{w}_0,K} C \|\rho_0\|_{L^\infty(K)}} \tag{49}$$

where  $C > 0$  is a constant that depends only on the shape regularity of  $K$ , then the local solver

$$(\widehat{p}_h, s) \longmapsto (p_h, \mathbf{q}_h)$$

is well-posed.

---

<sup>4</sup>Which is what the local solver does.

**Proof:**

As the local problems have a finite dimension, we only need to prove uniqueness of the solution. We therefore assume that  $\widehat{p}_h = s = 0$  and we need to prove that the system

$$(\mathbf{W}_0 \mathbf{q}_h, \mathbf{r}_h)_K - (p_h, \operatorname{div}(\mathbf{r}_h))_K = 0, \quad (50a)$$

$$\begin{aligned} & -\omega^2 (\rho_0 p_h, w_h)_K - 2i\omega (\mathbf{b}_0 \cdot \nabla p_h, w_h)_K \\ & + (\operatorname{div}(\mathbf{q}_h), w_h)_K + 2i\omega \langle \tau P_M p_h - \tau_{\text{upw}} p_h, w_h \rangle_{\partial K} = 0, \end{aligned} \quad (50b)$$

has only one solution :  $(p_h, \mathbf{q}_h) = (0, \mathbf{0})$ .

We will prove the theorem by contradiction. We therefore assume that the system (50a)–(50b) has a non-zero solution  $(p_h, \mathbf{q}_h)$ .

Step 1: An energy-like system

We begin by testing (50b) with  $w_h = p_h$

$$-\omega^2 \|p_h\|_{\rho_0, K}^2 + 2i\omega (p_h \mathbf{b}_0, \nabla p_h)_K + (\operatorname{div}(\mathbf{q}_h), p_h)_K + \langle 2i\omega \tau P_M p_h - 2i\omega \tau_{\text{upw}} p_h, p_h \rangle_{\partial K} = 0 \quad (51)$$

Then, (50a) is tested with  $\mathbf{r}_h = \mathbf{q}_h$  and conjugated :

$$\|\mathbf{q}_h\|_{\mathbf{W}_0, K}^2 - (\operatorname{div}(\mathbf{q}_h), p_h)_K = 0 \quad (52)$$

We now add (51) and (52) leading to

$$\|\mathbf{q}_h\|_{\mathbf{W}_0, K}^2 - \omega^2 \|p_h\|_{\rho_0, K}^2 + 2i\omega (p_h \mathbf{b}_0, \nabla p_h)_K + \langle 2i\omega \tau P_M p_h - 2i\omega \tau_{\text{upw}} p_h, p_h \rangle_{\partial K} = 0 \quad (53)$$

We now obtain the following system by taking the real and imaginary parts of (53)

$$\Re : \quad \|\mathbf{q}_h\|_{\mathbf{W}_0, K}^2 - \omega^2 \|p_h\|_{\rho_0, K}^2 - 2\omega \Im (p_h \mathbf{b}_0, \nabla p_h)_K = 0 \quad (54)$$

$$\Im : \quad \Re (p_h \mathbf{b}_0, \nabla p_h)_K + \langle \tau P_M p_h, p_h \rangle_{\partial K} = \langle \tau_{\text{upw}} p_h, p_h \rangle_{\partial K} \quad (55)$$

Indeed, as  $P_M p_h \in M_h$  and  $\tau$  is constant on each edge, one has

$$\langle \tau P_M p_h, p_h \rangle_{\partial K} = \langle \tau P_M p_h, P_M p_h \rangle_{\partial K} \in \mathbb{R}$$

Step 2: We focus on (55) to express  $p_h|_{\partial K}$ .

By the [LEMMA 4.1](#) we have

$$\Re (p_h \mathbf{b}_0, \nabla p_h)_K = \frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n}) p_h, p_h \rangle_{\partial K}$$

and (55) becomes

$$\frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n}) p_h, p_h \rangle_{\partial K} + \langle \tau P_M p_h, p_h \rangle_{\partial K} = \langle \tau_{\text{upw}} p_h, p_h \rangle_{\partial K}$$

For the sake of simplicity, we assume that the sign of  $\mathbf{b}_0 \cdot \mathbf{n}$  is constant on each edge. It amounts to assuming that  $h_K$  is small enough. For a given edge  $e \in \mathcal{E}(K)$ , the three following cases are exhaustive:

- *Case 1:*  $\mathbf{b}_0 \cdot \mathbf{n} < 0$ : therefore  $\tau_{\text{upw}} := \max(\mathbf{b}_0 \cdot \mathbf{n}, 0) = 0$  and

$$\underbrace{\langle \tau P_M p_h, P_M p_h \rangle_{\partial K}}_{\leq 0 \text{ by (48) as } \tau < 0} - \frac{1}{2} \overbrace{\langle \mathbf{b}_0 \cdot \mathbf{n} | p_h, p_h \rangle_{\partial K}}^{\leq 0} = 0$$

- *Case 2:*  $\mathbf{b}_0 \cdot \mathbf{n} > 0$ : therefore  $\tau_{\text{upw}} := \max(\mathbf{b}_0 \cdot \mathbf{n}, 0) = \mathbf{b}_0 \cdot \mathbf{n}$  and

$$\begin{aligned} & \frac{1}{2} \langle \mathbf{b}_0 \cdot \mathbf{n} | p_h, p_h \rangle_{\partial K} + \langle \tau P_M p_h, P_M p_h \rangle_{\partial K} = \langle \mathbf{b}_0 \cdot \mathbf{n} | p_h, p_h \rangle_{\partial K} \\ \iff & \underbrace{\langle \tau P_M p_h, P_M p_h \rangle_{\partial K}}_{\leq 0 \text{ by (48) as } \tau < 0} - \overbrace{\frac{1}{2} \langle \mathbf{b}_0 \cdot \mathbf{n} | p_h, p_h \rangle_{\partial K}}^{\leq 0} = 0 \end{aligned}$$

- *Case 3:*  $\mathbf{b}_0 \cdot \mathbf{n} = 0$ : in this case we only have

$$\langle \tau P_M p_h, P_M p_h \rangle_{\partial K} = 0.$$

In the first two cases we have  $p_h|_{\partial K} = P_M p_h = 0$  and in the third one we only have  $P_M p_h = 0$ . In particular, the following identity holds for all the three previous cases

$$\int_{\partial K} p_h d\sigma = 0,$$

indeed as  $P_M p_h$  is the  $L^2$ -orthogonal projection of  $p_h$  onto  $M_h$ , we have

$$\int_{\partial K} P_M p_h \mu_h^* d\sigma = \int_{\partial K} p_h \mu_h^* d\sigma, \quad \forall \mu_h \in M_h := \prod_{e \in \mathcal{E}(K)} \mathcal{P}_k(e),$$

and the previous identity is obtained by taking  $\mu_h = 1$ .

Step 3: Contradiction

As  $p_h \in \mathcal{P}_{k+1}(K)$ , we have  $p_h \in H^1(K)$  and the following Poincaré-Friedrichs inequality holds<sup>5</sup>

$$\|p_h\|_K \leq Ch_K \|\nabla p_h\|_K,$$

see [EG04, Lemma B.66] with  $f(v) = \int_{\partial K} v d\sigma$ . The constant  $C$  is the same one as in (49).

Going back to (50a), integrating by parts and testing it with  $\mathbf{r}_h = \nabla p_h$  we have

$$\begin{aligned} \|\nabla p_h\|_K^2 &= |(\mathbf{W}_0 \mathbf{q}_h, \nabla p_h)_K| \\ &\leq C_{\mathbf{W}_0, K} \|\mathbf{q}_h\|_{\mathbf{W}_0, K} \|\nabla p_h\|_K \\ \|\nabla p_h\|_K &\leq C_{\mathbf{W}_0, K} \|\mathbf{q}_h\|_{\mathbf{W}_0, K} \end{aligned} \tag{56}$$

On the other hand, from (54) we see that

$$\begin{aligned} \|\mathbf{q}_h\|_{\mathbf{W}_0, K}^2 &= \omega^2 \|p_h\|_{\rho_0, K}^2 + 2\omega \Im (p_h \mathbf{b}_0, \nabla p_h)_K \\ &\leq \omega^2 \|\rho_0\|_{L^\infty(K)} \|p_h\|_K^2 + 2\omega \|\mathbf{b}_0\|_{L^\infty(K)} \|p_h\|_K \|\nabla p_h\|_K \\ \|\mathbf{q}_h\|_{\mathbf{W}_0, K}^2 &\leq C^2 \|\rho_0\|_{L^\infty(K)} \omega^2 h_K^2 \|\nabla p_h\|_K^2 + 2C \|\mathbf{b}_0\|_{L^\infty(K)} \omega h_K \|\nabla p_h\|_K^2 \end{aligned} \tag{57}$$

Combining (56) and (57) we have

$$\|\nabla p_h\|_K^2 \leq C_{\mathbf{W}_0, K}^2 \left[ C^2 \|\rho_0\|_{L^\infty(K)} \omega^2 h_K^2 \|\nabla p_h\|_K^2 + 2C \|\mathbf{b}_0\|_{L^\infty(K)} \omega h_K \|\nabla p_h\|_K^2 \right],$$

as we assumed  $(\mathbf{q}_h, p_h) \neq (\mathbf{0}, 0)$  we can divide by  $\|\nabla p_h\|_K$  to obtain

$$1 \leq C_{\mathbf{W}_0, K}^2 \left[ C^2 \|\rho_0\|_{L^\infty(K)} \omega^2 h_K^2 + 2C \|\mathbf{b}_0\|_{L^\infty(K)} \omega h_K \right]. \tag{58}$$

<sup>5</sup>When  $\mathbf{b}_0 \cdot \mathbf{n} \neq 0$  we can use the standard Poincaré inequality instead.-

We now define the function

$$f : \alpha \mapsto C_{\mathbf{w}_0, K}^2 C^2 \|\rho_0\|_{L^\infty(K)} \alpha^2 + 2C_{\mathbf{w}_0, K} C \|\mathbf{b}_0\|_{L^\infty(K)} \alpha - 1$$

Rewriting (58) in terms of  $f$  gives

$$f(\omega h_K) \geq 0.$$

We notice that  $f$  is a second-order polynomial whose roots are

$$\alpha_{\pm} = \frac{-C_{\mathbf{w}_0, K} \|\mathbf{b}_0\|_{L^\infty(K)} \pm \left( C_{\mathbf{w}_0, K}^2 \|\mathbf{b}_0\|_{L^\infty(K)}^2 + \|\rho_0\|_{L^\infty(K)} \right)^{\frac{1}{2}}}{C_{\mathbf{w}_0, K} C \|\rho_0\|_{L^\infty(K)}}$$

As the leading coefficient of  $f$  is positive, we know that

$$\forall \alpha \in (\alpha_-, \alpha_+), \quad f(\alpha) < 0$$

and it is obvious that  $\alpha_- < 0$  and  $\alpha_+ > 0$ .

Finally, we can see that the assumption on  $\omega h_K$  (49) is exactly

$$0 < \omega h_K < \alpha_+,$$

which means

$$f(\omega h_K) < 0.$$

This is the desired contradiction and concludes the proof, as we necessarily have  $p_h \equiv 0$  and  $\mathbf{q}_h \equiv \mathbf{0}$ . ■

**Remark 4.3:** For triangular elements, the constant  $C$  satisfies

$$C < \frac{1}{\pi}.$$

**Remark 4.4:** when  $\mathbf{b}_0 = \mathbf{0}$ , the solvability assumption (48) becomes

$$\omega h_K < \frac{1}{C C_{\mathbf{w}_0, K} \|\rho_0\|_{L^\infty(K)}^{\frac{1}{2}}}$$

which is similar to the ones given in [DS19, Prop. 3.9] and [Hun19, Prop. 3.4.2].

**Remark 4.5:** this proof is written for the HDG+ method, for the more standard HDG method only minor changes are needed : in *Step 2*,  $P_{Mp}$  should be replaced by  $p$ . Assumption (48) can therefore be replaced with

$$\forall e \in \mathcal{E}(K), \quad \tau|_e < 0 \quad \text{or} \quad \tau|_e > \max_e \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right).$$

### 4.3 Error analysis of the HDG+ method

In this section we will carry out a detailed error analysis of the HDG+ method. The adaptation of this process to the HDG method is straightforward, see [SUBSECTION 4.4](#).

We chose to use the orthogonal  $L^2$  projections instead of the tailored HDG(+) projections that fit the numerical trace. As we study problems involving convection, the design of a new projection would be required as using the standard HDG(+) projection would not lead to cleaner error system. The design of such a projection seems very difficult when  $\mathbf{b}_0$  is not constant.

We denote by  $\pi_{\mathbf{V}}$ ,  $\pi_W$  and  $P_M$  the  $L^2$ -orthogonal projections onto  $\mathbf{V}_h$ ,  $W_h$  and  $M_h$  respectively. We recall the following estimates due to standard approximation theory for polynomials and trace inequalities which will be useful for our analysis, see *eg.* [EG04, Prop. 1.135] :

$$\|p - \pi_W p\|_{\mathcal{O}} \lesssim h^s \|p\|_{s,\mathcal{O}}, \quad 0 \leq s \leq k+2, \quad (59a)$$

$$\|\mathbf{q} - \pi_{\mathbf{V}} \mathbf{q}\|_{\mathcal{O}} \lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}}, \quad 0 \leq t \leq k+1, \quad (59b)$$

$$\|p - P_M p\|_{\partial\mathcal{T}_h} \lesssim h^{s-\frac{1}{2}} \|p\|_{s,\mathcal{O}}, \quad 1 \leq s \leq k+1, \quad (59c)$$

$$\|p - \pi_W p\|_{\partial K} \lesssim h^{s-\frac{1}{2}} \|p\|_{s,K}, \quad 1 \leq s \leq k+2, \quad (59d)$$

$$\|\mathbf{q} \cdot \mathbf{n} - \pi_{\mathbf{V}} \mathbf{q} \cdot \mathbf{n}\|_{\partial K} \lesssim h^{t-\frac{1}{2}} \|\mathbf{q}\|_{t,K}, \quad 1 \leq t \leq k+1, \quad (59e)$$

where  $a \lesssim b$  means that there exists a constant  $C > 0$  independent of the mesh size and frequency such that  $a \leq Cb$ .

We will also frequently use the following inverse inequality

$$\|w\|_{\partial K} \lesssim h^{-\frac{1}{2}} \|w\|_K, \quad \forall w \in W_h. \quad (60)$$

### 4.3.1 Error equations

Let  $(p, \mathbf{q})$  be the solution of the original problem (7a)–(7b). We define the projection errors

$$\delta_h^{\mathbf{q}} := \pi_{\mathbf{V}} \mathbf{q} - \mathbf{q} \ ; \ \delta_h^p := \pi_W p - p \ ; \ \widehat{\delta}_h^p := p - P_M p$$

and

$$\varepsilon_h^{\mathbf{q}} := \pi_{\mathbf{V}} \mathbf{q} - \mathbf{q}_h \in \mathbf{V}_h \ ; \ \varepsilon_h^p := \pi_W p - p_h \in W_h \ ; \ \widehat{\varepsilon}_h^p := P_M p - \widehat{p}_h \in M_h$$

**Lemma 4.2:**

The error quantities  $(\boldsymbol{\varepsilon}_h^q, \varepsilon_h^p, \widehat{\varepsilon}_h^p)$  satisfy the following error equations:

$$(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \mathbf{r}_h)_{\mathcal{T}_h} - (\varepsilon_h^p, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + \langle \widehat{\varepsilon}_h^p, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \mathbf{r}_h)_{\mathcal{T}_h} \quad (61a)$$

$$-\omega^2 (\rho_0 \varepsilon_h^p, w_h)_{\mathcal{T}_h} + 2i\omega (\varepsilon_h^p \mathbf{b}_0, \nabla w_h)_{\mathcal{T}_h} - (\boldsymbol{\varepsilon}_h^q, \nabla w_h)_{\mathcal{T}_h} + \langle \widehat{\mathbf{Q}} \cdot \mathbf{n} - \widehat{\mathbf{Q}}_h \cdot \mathbf{n}, w_h \rangle_{\partial \mathcal{T}_h} = -\omega^2 (\rho_0 \delta_h^p, w_h)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{b}_0, w_h)_{\mathcal{T}_h} \quad (61b)$$

$$\langle \widehat{\mathbf{Q}} \cdot \mathbf{n} - \widehat{\mathbf{Q}}_h \cdot \mathbf{n}, \mu_h \rangle_{\partial \mathcal{T}_h} = 0 \quad (61c)$$

where

$$\widehat{\mathbf{Q}} \cdot \mathbf{n} = \mathbf{q} \cdot \mathbf{n} - 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) p \quad \text{on } \partial \mathcal{T}_h$$

and

$$\begin{aligned} \widehat{\mathbf{Q}} \cdot \mathbf{n} - \widehat{\mathbf{Q}}_h \cdot \mathbf{n} &= \boldsymbol{\varepsilon}_h^q \cdot \mathbf{n} - 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p + 2i\omega \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) - 2i\omega \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \\ &\quad - \boldsymbol{\delta}_h^q \cdot \mathbf{n} - 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \delta_h^p - 2i\omega \tau P_M \delta_h^p + 2i\omega \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p) \end{aligned} \quad (62)$$

**Proof:** Notice that  $(p, \mathbf{q})$  satisfy the equations (46a)–(46b)–(46c), introduce the projections wherever possible and subtract the actual discrete equations. ■

**A useful estimate:** We will need to use the following estimate for  $\|\nabla \varepsilon_h^p\|_{\partial \mathcal{T}_h}$  to carry out our analysis

**Lemma 4.3:**

The following estimate holds

$$\|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} \leq C_{\mathbf{W}_0, \mathcal{T}_h} \left( \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h} + \|\boldsymbol{\delta}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h} \right) + C \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}$$

**Proof:** Going back to (61a), testing with  $\mathbf{r}_h = \nabla \varepsilon_h^p$  and integrating by parts leads to

$$(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \nabla \varepsilon_h^p)_{\mathcal{T}_h} + \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h}^2 + \langle \widehat{\varepsilon}_h^p - \varepsilon_h^p, \nabla \varepsilon_h^p \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \nabla \varepsilon_h^p)_{\mathcal{T}_h}$$

As  $\varepsilon_h^p \in W_h$ ,  $\nabla \varepsilon_h^p \cdot \mathbf{n} \in \mathcal{P}_k$  and we can use the following property of the projection  $P_M$  :

$$\langle P_M \varepsilon_h^p, \nabla \varepsilon_h^p \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = \langle \varepsilon_h^p, \nabla \varepsilon_h^p \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}$$

Using the Cauchy-Schwartz inequality we get

$$\begin{aligned} \left| \langle \widehat{\varepsilon}_h^p - \varepsilon_h^p, \nabla \varepsilon_h^p \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} \right| &= \left| \langle \widehat{\varepsilon}_h^p - P_M \varepsilon_h^p, \nabla \varepsilon_h^p \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} \right| \\ &\leq C \|P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} \|\nabla \varepsilon_h^p\|_{\partial \mathcal{T}_h} \end{aligned}$$

for some constant  $C > 0$ .

Using the following trace inequality (60)

$$\forall w \in W_h, \quad \|w\|_{\partial K} \leq Ch_K^{-\frac{1}{2}} \|w\|_K$$

we have

$$\|\nabla \varepsilon_h^p\|_{\mathcal{T}_h}^2 \leq C_{\mathbf{W}_0, K} \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, K} \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} + \|P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} Ch_K^{-\frac{1}{2}} \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} + C_{\mathbf{W}_0, K} \|\boldsymbol{\delta}_h^q\|_{\mathbf{W}_0, K} \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h}$$

which is the desired estimate as  $\tau|_K = \mathcal{O}(h_K^{-1})$ . ■

**Using the Poincaré-Wirtinger inequality:** We denote by  $\{u\}$  the  $L^2$ -projection of  $u$  on  $\mathcal{P}_0$ , ie

$$\forall K \in \mathcal{T}_h, \quad \{u\}|_K = \frac{1}{|K|} \int_K u d\mathbf{x}$$

We will need to subtract  $\{u\}$  from the equations to apply the Poincaré-Wirtinger inequality :

$$\|u - \{u\}\|_{\mathcal{T}_h} \leq Ch \|\nabla u\|_{\mathcal{T}_h}. \quad (63)$$

We can do that thanks to the following property of the projections  $\pi_W$  and  $\pi_V$ , indeed we have for  $\boldsymbol{\xi}$  and  $\mathbf{q}$

$$(\pi_V \mathbf{q}, \{\mathbf{W}_0 \boldsymbol{\xi}\})_{\mathcal{T}_h} = (\mathbf{q}, \{\mathbf{W}_0 \boldsymbol{\xi}\})_{\mathcal{T}_h} \quad \text{as} \quad \{\mathbf{W}_0 \boldsymbol{\xi}\} \in \mathcal{P}_0 \subset \mathcal{P}_k$$

therefore

$$(\boldsymbol{\delta}_h^q, \mathbf{W}_0 \boldsymbol{\xi})_{\mathcal{T}_h} = (\mathbf{q} - \pi_V \mathbf{q}, \mathbf{W}_0 \boldsymbol{\xi})_{\mathcal{T}_h} = (\mathbf{q} - \pi_V \mathbf{q}, \mathbf{W}_0 \boldsymbol{\xi} - \{\mathbf{W}_0 \boldsymbol{\xi}\})_{\mathcal{T}_h}. \quad (64)$$

Similar results can be obtained in the same way for the other quantities.

**Best approximation property of  $P_M$ :** During the analysis, we will often need to compare quantities like  $\|u - P_M u\|_{\partial K}$  and  $\|u - \{u\}\|_{\partial K}$ .

**Lemma 4.4:**

For  $u \in \mathcal{P}_{k+1}(\mathcal{T}_h)$ , the following inequality holds

$$\|u - P_M u\|_{\partial K} \lesssim \|u - \{u\}\|_{\partial K}$$

**Proof:**

We recall that  $M_h := \prod_{e \in \mathcal{E}_h} \mathcal{P}_k(e)$  is a finite-dimensional vector subspace of  $L^2(\partial \mathcal{T}_h)$ . We recalled that functions in  $M_h$  are bi-valued piecewise polynomials of degree up to  $k$  on the skeleton of the mesh.

On an internal edge  $e = \partial K_- \cap \partial K_+$ , we define

$$\{u\}_e := \begin{cases} \{u^-\} & \text{on } K_- \\ \{u^+\} & \text{on } K_+ \end{cases}, \quad \text{where } u^\pm = u|_{K^\pm}.$$

With this definition  $\{u\}_e$  is a bi-valued piecewise constant on the skeleton of the mesh, and therefore  $\{u\}_e \in M_h$ .

As  $P_M$  is the orthogonal projection onto  $M_h$ , we can use the Hilbert projection theorem to obtain

$$\|u - P_M u\|_{\partial \mathcal{T}_h} \leq \inf_{v \in M_h} \|u - v\|_{\partial \mathcal{T}_h}.$$

We can therefore conclude that

$$\|u - P_M u\|_{\partial \mathcal{T}_h} \leq \|u - \{u\}_e\|_{\partial \mathcal{T}_h}.$$

When no confusions are possible, we will denote  $\{u\}_e$  by  $\{u\}$ . ■

This property will often be referred to as the *best approximation property* of  $P_M$ .



**Discrete energy-like equality:** We will now establish a discrete energy-like equality which will be one of the key ingredients to study the convergence of our method.

**Lemma 4.5:**

The following discrete energy-like equality holds

$$\begin{aligned} & \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{W}_0, \mathcal{T}_h}^2 - \omega^2 \|\varepsilon_h^p\|_{\rho_0, \mathcal{T}_h}^2 - 2i\omega \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 - 2i\omega \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 \\ & \quad - 2\omega \Im (\varepsilon_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_K \\ & = -\omega^2 (\rho_0 \delta_h^p, \varepsilon_h^p)_K + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_K + (\boldsymbol{\varepsilon}_h^{\mathbf{q}}, \mathbf{W}_0 \boldsymbol{\delta}_h^{\mathbf{q}})_K \\ & \quad + \left\langle \boldsymbol{\delta}_h^{\mathbf{q}} \cdot \mathbf{n} + 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p + 2i\omega \tau P_M \delta_h^p - 2i\omega \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \right\rangle_{\partial \mathcal{T}_h} \end{aligned} \quad (65)$$

Furthermore if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s \in [1, k+2]$  and  $t \in [1, k+1]$  then the following estimate holds

$$\begin{aligned} & \left| \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{W}_0, \mathcal{T}_h}^2 - 2i\omega \left( \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 + \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 \right) \right| \\ & \lesssim \omega^2 \|\varepsilon_h^p\|_{\mathcal{T}_h}^2 + \omega \|\varepsilon_h^p\|_{\mathcal{T}_h} \left( \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{W}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} + \omega h^{s-1} \|p\|_{s, \mathcal{O}} \right) \\ & \quad + \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{W}_0, \mathcal{T}_h} \left( h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \omega h^s \|p\|_{s, \mathcal{O}} \right) + h^{2t} \|\mathbf{q}\|_{t, \mathcal{O}}^2 + \omega h^{s-1} \|p\|_{s, \mathcal{O}} h^t \|\mathbf{q}\|_{t, \mathcal{O}} \\ & \quad + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \left( \omega h^{s-1} \|p\|_{s, \mathcal{O}} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} \right) \end{aligned} \quad (66)$$

where

$$h := \max_{K \in \mathcal{T}_h} h_K.$$

**Proof:** Test (61a)–(61b)–(61c) with  $(\boldsymbol{\varepsilon}_h^{\mathbf{q}}, \varepsilon_h^p, \widehat{\varepsilon}_h^p)$  and sum the resulting equations to obtain

$$\begin{aligned} & \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{W}_0, K}^2 - \omega^2 \|\varepsilon_h^p\|_{\rho_0, K}^2 + 2i\omega (\varepsilon_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_K + \left\langle \widehat{\mathbf{Q}} \cdot \mathbf{n} - \widehat{\mathbf{Q}}_h \cdot \mathbf{n} - \boldsymbol{\varepsilon}_h^{\mathbf{q}} \cdot \mathbf{n}, \varepsilon_h^p - \widehat{\varepsilon}_h^p \right\rangle_{\partial K} = \\ & \quad -\omega^2 (\rho_0 \delta_h^p, \varepsilon_h^p)_K + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_K + (\boldsymbol{\varepsilon}_h^{\mathbf{q}}, \mathbf{W}_0 \boldsymbol{\delta}_h^{\mathbf{q}})_K \end{aligned}$$

We will now compute the boundary terms using (62).

Boundary terms involving  $P_M$ :

Notice that, as  $P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p \in M_h(K)$  and  $\tau \in \mathcal{R}_0$ , we have

$$\langle \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p), P_M \varepsilon_h^p \rangle_{\partial K} = \langle \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p \rangle_{\partial K}$$

and therefore

$$\langle \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial K} = \langle \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p), P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial K} = \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial K}^2$$

We also have the following estimate

$$\begin{aligned} 2\omega \left| \langle \tau P_M \delta_h^p, \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} \right| & \leq 2\omega \left| \langle |\tau| P_M \delta_h^p, P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} \right| \\ & \leq 2\omega \left| \left\langle |\tau|^{\frac{1}{2}} \delta_h^p, |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\rangle_{\partial \mathcal{T}_h} \right| \\ & \lesssim \omega \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \left\| |\tau|^{\frac{1}{2}} \delta_h^p \right\|_{\partial \mathcal{T}_h} \\ (\tau = \mathcal{O}(h^{-1}) \text{ and by (59d)}) & \lesssim \omega h^{s-1} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \|p\|_{s, \mathcal{O}} \end{aligned}$$

Boundary terms involving convection:

As in *Step 2* of the proof of [THEOREM 4](#), we will separate the volumetric term involving  $\mathbf{b}_0$  into its real and imaginary parts. By the [LEMMA 4.1](#) we have

$$\Re(\varepsilon_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_K = \frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \varepsilon_h^p \rangle_{\partial K}$$

and we can now obtain the second boundary norm

$$\begin{aligned} & \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p, \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} - \frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \varepsilon_h^p \rangle_{\partial \mathcal{T}_h} + \langle \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} \\ &= \left\langle -\frac{1}{2} (\mathbf{b}_0 \cdot \mathbf{n}) (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \right\rangle_{\partial \mathcal{T}_h} - \frac{1}{2} \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p, \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} + \langle \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} \\ &= \left\langle \left( \tau_{\text{upw}} - \frac{1}{2} (\mathbf{b}_0 \cdot \mathbf{n}) \right) (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \right\rangle_{\partial \mathcal{T}_h} \\ &= \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 \end{aligned}$$

indeed as  $\widehat{\varepsilon}_h^p$  is single-valued across the skeleton of the mesh we have

$$\langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p, \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} = \langle \llbracket \widehat{\varepsilon}_h^p \mathbf{b}_0 \rrbracket, \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} = 0$$

and we also use that

$$\tau_{\text{upw}} - \frac{1}{2} (\mathbf{b}_0 \cdot \mathbf{n}) = \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \geq 0 \quad (67)$$

to use the square root.

We will now eliminate the terms involving  $\tau_{\text{upw}}$  from the right-hand side.

$$\begin{aligned} & \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \lesssim \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} \\ & \lesssim \|\varepsilon_h^p - P_M \varepsilon_h^p\|_{\partial \mathcal{T}_h} + \|P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} \\ \text{(by [LEMMA 4.4](#) and } \tau = \mathcal{O}(h^{-1})\text{)} & \lesssim \|\varepsilon_h^p - \{\varepsilon_h^p\}\|_{\partial \mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \\ \text{(by [\(60\)](#))} & \lesssim h^{-\frac{1}{2}} \|\varepsilon_h^p - \{\varepsilon_h^p\}\|_{\mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \\ \text{(by [\(63\)](#))} & \lesssim h^{\frac{1}{2}} \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \\ \text{(by [LEMMA 4.3](#))} & \lesssim h^{\frac{1}{2}} \left( \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \|\boldsymbol{\delta}_h^q\|_{\mathcal{T}_h} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \\ \text{(by [\(59e\)](#))} & \lesssim h^{\frac{1}{2}} \left( \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \quad (68) \end{aligned}$$

Using the following inequalities that can be derived from [\(67\)](#)

$$\begin{aligned} \tau_{\text{upw}} & \lesssim \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \\ \mathbf{b}_0 \cdot \mathbf{n} & \lesssim \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \end{aligned}$$

and using (59d) with  $s - 1$  instead of  $s$  to keep  $s \leq k + 2$ , we deduce that

$$\begin{aligned}
2\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p, \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h} &= 2\omega \langle (\mathbf{b}_0 \cdot \mathbf{n})^{\frac{1}{2}} \delta_h^p, (\mathbf{b}_0 \cdot \mathbf{n})^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \rangle_{\partial\mathcal{T}_h} \\
&\lesssim 2\omega \left\langle (\mathbf{b}_0 \cdot \mathbf{n})^{\frac{1}{2}} \delta_h^p, \left(\frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}|\right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\rangle_{\partial\mathcal{T}_h} \\
&\lesssim \omega \|\widehat{\delta}_h^p\|_{\partial\mathcal{T}_h} \left\| \left(\frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}|\right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \\
&\lesssim \omega h^{s-\frac{3}{2}} \left\| \left(\frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}|\right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \|p\|_{s,\mathcal{O}} \\
&\lesssim \omega h^{s-1} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t,\mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \right) \|p\|_{s,\mathcal{O}}
\end{aligned}$$

and

$$\begin{aligned}
2\omega \langle \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h} &\lesssim \omega \|\delta_h^p - \widehat{\delta}_h^p\|_{\partial\mathcal{T}_h} \left\| \left( \tau_{\text{upw}} - \frac{1}{2} \mathbf{b}_0 \cdot \mathbf{n} \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \\
&\lesssim \omega h^{s-\frac{3}{2}} \left\| \left(\frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}|\right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \|p\|_{s,\mathcal{O}} \\
&\lesssim \omega h^{s-1} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t,\mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \right) \|p\|_{s,\mathcal{O}}
\end{aligned}$$

Boundary term involving  $\delta_h^q$ :

$$\begin{aligned}
\langle \delta_h^q \cdot \mathbf{n}, \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h} &= \langle \delta_h^q \cdot \mathbf{n}, \varepsilon_h^p - P_M \varepsilon_h^p + P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h} \\
&= \underbrace{\langle \delta_h^q \cdot \mathbf{n}, \varepsilon_h^p - P_M \varepsilon_h^p \rangle_{\partial\mathcal{T}_h}}_{=:T_1} + \underbrace{\langle \delta_h^q \cdot \mathbf{n}, P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h}}_{=:T_2}
\end{aligned}$$

Using a weighted Cauchy-Schwartz inequality and recalling that  $\tau = \mathcal{O}(h^{-1})$  we have

$$\begin{aligned}
T_2 &\lesssim \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} h^{\frac{1}{2}} \|\delta_h^q\|_{\partial\mathcal{T}_h} \\
(\text{by (59e)}) &\lesssim h^t \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \|\mathbf{q}\|_{t,\mathcal{O}}
\end{aligned}$$

$$\begin{aligned}
T_1 &= \langle \delta_h^q \cdot \mathbf{n}, \varepsilon_h^p - P_M \varepsilon_h^p \rangle_{\partial\mathcal{T}_h} \\
&\lesssim \|\delta_h^q\|_{\partial\mathcal{T}_h} \|\varepsilon_h^p - P_M \varepsilon_h^p\|_{\partial\mathcal{T}_h} \\
(\text{by LEMMA 4.4}) &\lesssim \|\delta_h^q\|_{\partial\mathcal{T}_h} \|\varepsilon_h^p - \{\varepsilon_h^p\}\|_{\partial\mathcal{T}_h} \\
(\text{by (59e) and (59d)}) &\lesssim h^{t-\frac{1}{2}} \|\mathbf{q}\|_{t,\mathcal{O}} h^{-\frac{1}{2}} \|\varepsilon_h^p - \{\varepsilon_h^p\}\|_{\mathcal{T}_h} \\
(\text{by (63)}) &\lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}} \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} \\
(\text{by LEMMA 4.3}) &\lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \|\delta_h^q\|_{\mathcal{T}_h} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \right) \\
(\text{by (59b)}) &\lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \right) + h^{2t} \|\mathbf{q}\|_{t,\mathcal{O}}^2
\end{aligned}$$

Therefore

$$\langle \delta_h^q \cdot \mathbf{n}, \varepsilon_h^p - \widehat{\varepsilon}_h^p \rangle_{\partial\mathcal{T}_h} \lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \right) + h^{2t} \|\mathbf{q}\|_{t,\mathcal{O}}^2$$

**Volumetric terms:** By similar computations using the Cauchy-Schwartz inequality, the projection estimates in (59) and LEMMA 4.3 we can show that

$$\begin{aligned}
 \omega^2 (\rho_0 \delta_h^p, \varepsilon_h^p)_{\mathcal{T}_h} &\lesssim \omega^2 h^s \|p\|_{s,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h} \\
 2\omega (\delta_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} &\lesssim \omega h^s \|p\|_{s,\mathcal{O}} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t,\mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \\
 (\mathbf{W}_0 \varepsilon_h^q, \delta_h^q)_{\mathcal{T}_h} &\lesssim h^t \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} \|\mathbf{q}\|_{t,\mathcal{O}} \\
 2\omega (\varepsilon_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} &\lesssim \|\varepsilon_h^p\|_{\mathcal{T}_h} \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t,\mathcal{O}} + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right)
 \end{aligned}$$

■

We can rewrite estimate (66) in a more readable form :

#### Corollary 4.1:

The following estimate holds

$$\begin{aligned}
 &\left| \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h}^2 - 2i\omega \left( \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 + \left\| \left( \frac{1}{2} \mathbf{b}_0 \cdot \mathbf{n} \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 \right) \right| \\
 &\lesssim \varepsilon \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h}^2 + \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}^2 \right) + \frac{1}{\varepsilon} \left( \omega^2 \|\varepsilon_h^p\|_{\mathcal{T}_h}^2 + h^{2t} \|\mathbf{q}\|_{t,\mathcal{O}}^2 + \omega^2 h^{2s-2} \|p\|_{s,\mathcal{O}}^2 \right)
 \end{aligned}$$

**Proof:** apply the weighted Young's inequality to the right-hand side of (66). The value of  $\varepsilon$  will be discussed later. ■

### 4.3.2 Adjoint problem

As the identity (65) does not allow us to directly obtain any error estimate, we need to use a duality argument. For an introduction to the *Aubin-Nitsche method* we refer to [EG04, Sec. 2.3.4], similar processes for HDG(+) methods in the context of wave equations have been carried out in [QSS16], [Hum19, Sec. 3.5] and [DS19, Sec. 3.5.2] and for coercive problems with convection in [QS16a].

The adjoint problem is

$$\begin{aligned}
 \mathbf{W}_0 \boldsymbol{\xi} - \nabla \theta &= 0 && \text{in } \mathcal{O} \\
 -\rho_0 \omega^2 \theta - 2i\omega \mathbf{b}_0 \cdot \nabla \theta - \operatorname{div}(\boldsymbol{\xi}) &= \varepsilon_h^p && \text{in } \mathcal{O} \\
 \theta &= 0 && \text{on } \Gamma_D \\
 \boldsymbol{\xi} \cdot \mathbf{n} - 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \theta &= 0 && \text{on } \Gamma_N
 \end{aligned}$$

and  $(\boldsymbol{\xi}, \theta) \in \mathbf{H}^1(\mathcal{O}) \times H^2(\mathcal{O})$  satisfy the following discrete problem for all  $(\mathbf{r}_h, w_h, \mu_h) \in \mathbf{V}_h \times W_h \times M_h$

$$(\mathbf{W}_0 \boldsymbol{\xi}, \mathbf{r}_h)_{\mathcal{T}_h} + (\theta, \operatorname{div}(\mathbf{r}_h))_{\mathcal{T}_h} + \langle \theta, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = 0 \quad (70a)$$

$$-\omega^2 (\rho_0 \theta, w_h)_{\mathcal{T}_h} + (2i\omega \theta \mathbf{b}_0 + \boldsymbol{\xi}, \nabla w_h)_{\mathcal{T}_h} + \langle \boldsymbol{\xi} \cdot \mathbf{n} - 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \theta, w_h \rangle_{\partial \mathcal{T}_h} = (\varepsilon_h^p, w_h)_{\mathcal{T}_h} \quad (70b)$$

$$\langle \boldsymbol{\xi} \cdot \mathbf{n}, \mu_h \rangle_{\partial \mathcal{T}_h \setminus \Gamma_D} = 0 \quad (70c)$$

The last equation (70c) translates the continuity of  $\boldsymbol{\xi} \cdot \mathbf{n}$  between the elements and should be interpreted as a jump term. Indeed by the same argument as when we discussed weak continuity of  $\mathbf{q}_h \cdot \mathbf{n}$  in SUBSECTION 4.1, we can show that

$$\sum_{K \in \mathcal{T}_h} \int_{\partial K} \boldsymbol{\xi} \cdot \mathbf{n} \mu_h^* d\sigma = \sum_{e \in \mathcal{E}_h^i} \int_e [[\boldsymbol{\xi}]] \mu_h^* d\sigma - \int_{\Gamma_D} \boldsymbol{\xi} \cdot \mathbf{n} \mu_h^* d\sigma. \quad (71)$$

**Remark 4.6:** The functional framework for (70c) is a bit complicated. Indeed the interior integrals should formally be interpreted as duality brackets between  $H^{-\frac{1}{2}}$  and  $H^{\frac{1}{2}}$  and the restriction of those distributions to a segment is not defined. Notice however that the right-hand side of (71) is well-defined, therefore giving meaning to the left-hand side. Moreover, as we assume additional regularity  $(\boldsymbol{\xi}, \theta) \in \mathbf{H}^1(\mathcal{O}) \times H^2(\mathcal{O})$  for the solution of the adjoint problem and as we will work with polynomial quantities at the discrete level, this is not problematic.

In our analysis we will need to use the following elliptic regularity estimate for the dual problem

$$\|\theta\|_{2,\mathcal{O}} + \|\boldsymbol{\xi}\|_{1,\mathcal{O}} \leq C_{\text{reg}} \|\varepsilon_h^p\|_{\mathcal{O}}. \quad (72)$$

This estimate holds when  $\boldsymbol{\xi}$  and  $\theta$  are regular enough, which amounts to requiring enough regularity on the background quantities  $\rho_0$ ,  $c_0$  and  $\mathbf{b}_0$ , and the convexity of the domain  $\mathcal{O}$ .

**Lemma 4.6:**

We have the following dual identity :

$$\begin{aligned} \|\varepsilon_h^p\|_{\mathcal{T}_h}^2 &= -(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi} - \boldsymbol{\xi})_{\mathcal{T}_h} + \omega^2 (\rho_0 \varepsilon_h^p, \pi_W \theta - \theta)_{\mathcal{T}_h} + 2i\omega (\nabla \varepsilon_h^p, (\pi_W \theta - \theta) \mathbf{b}_0)_{\mathcal{T}_h} \\ &\quad - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \pi_W \theta - \theta \rangle_{\partial \mathcal{T}_h} \\ &\quad + (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} - \omega^2 (\rho_0 \delta_h^p, \pi_W \theta)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla (\pi_W \theta))_{\mathcal{T}_h} \\ &\quad + 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p - \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) + \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} \\ &\quad - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p - \tau P_M \delta_h^p + \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} + \langle \boldsymbol{\delta}_h^q \cdot \mathbf{n}, \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} \end{aligned}$$

**Proof:** Introducing the projections in (70a)–(70b)–(70c) and testing with  $(\boldsymbol{\varepsilon}_h^q, \varepsilon_h^p, \widehat{\varepsilon}_h^p)$

$$(\mathbf{W}_0 \boldsymbol{\pi}_V \boldsymbol{\xi}, \boldsymbol{\varepsilon}_h^q)_{\mathcal{T}_h} + (\pi_W \theta, \text{div}(\boldsymbol{\varepsilon}_h^q))_{\mathcal{T}_h} - \langle P_M \theta, \boldsymbol{\varepsilon}_h^q \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = (\mathbf{W}_0 (\boldsymbol{\pi}_V \boldsymbol{\xi} - \boldsymbol{\xi}), \boldsymbol{\varepsilon}_h^q)_{\mathcal{T}_h} \quad (73a)$$

$$\begin{aligned} -\omega^2 (\rho_0 \pi_W \theta, \varepsilon_h^p)_{\mathcal{T}_h} + 2i\omega ((\pi_W \theta) \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} - (\text{div}(\boldsymbol{\pi}_V \boldsymbol{\xi}), \varepsilon_h^p)_{\mathcal{T}_h} - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \pi_W \theta, \varepsilon_h^p \rangle_{\partial \mathcal{T}_h} = \\ (\varepsilon_h^p, \varepsilon_h^p)_{\mathcal{T}_h} - \omega^2 (\rho_0 (\pi_W \theta - \theta), \varepsilon_h^p)_{\mathcal{T}_h} + 2i\omega ((\pi_W \theta - \theta) \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) (\pi_W \theta - \theta), \varepsilon_h^p \rangle_{\partial \mathcal{T}_h} \end{aligned} \quad (73b)$$

$$\langle \boldsymbol{\pi}_V \boldsymbol{\xi} \cdot \mathbf{n}, \widehat{\varepsilon}_h^p \rangle_{\partial \mathcal{T}_h} = 0 \quad (73c)$$

Now conjugate and sum those equations and compare with the sum of (61a)–(61b)–(61c) tested with  $(\boldsymbol{\pi}_V \boldsymbol{\xi}, \pi_W \theta, P_M \theta)$ .

$$(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} - (\varepsilon_h^p, \text{div}(\boldsymbol{\pi}_V \boldsymbol{\xi}))_{\mathcal{T}_h} + \langle \widehat{\varepsilon}_h^p, \boldsymbol{\pi}_V \boldsymbol{\xi} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = \ell_1(\boldsymbol{\pi}_V \boldsymbol{\xi}) \quad (74a)$$

$$-\omega^2 (\rho_0 \varepsilon_h^p, \pi_W \theta)_{\mathcal{T}_h} - 2i\omega (\nabla \varepsilon_h^p, (\pi_W \theta) \mathbf{b}_0)_{\mathcal{T}_h} \quad (74b)$$

$$\begin{aligned} + (\text{div}(\boldsymbol{\varepsilon}_h^q), \pi_W \theta)_{\mathcal{T}_h} + 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \pi_W \theta \rangle_{\partial \mathcal{T}_h} = \ell_2(\pi_W \theta) \\ - \langle \boldsymbol{\varepsilon}_h^q \cdot \mathbf{n}, P_M \theta \rangle_{\partial \mathcal{T}_h} = \ell_3(P_M \theta) \end{aligned} \quad (74c)$$

where

$$\begin{aligned}
 \ell_1(\boldsymbol{\pi}_V \boldsymbol{\xi}) &:= (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} \\
 \ell_2(\pi_W \theta) &:= -\omega^2 (\rho_0 \delta_h^p, \pi_W \theta)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla \pi_W \theta)_{\mathcal{T}_h} \\
 &\quad + 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p - \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) + \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta \rangle_{\partial \mathcal{T}_h} \\
 &\quad + \langle \boldsymbol{\delta}_h^q \cdot \mathbf{n} + 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p + 2i\omega \tau P_M \delta_h^p - 2i\omega \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \pi_W \theta \rangle_{\partial \mathcal{T}_h} \\
 \ell_3(P_M \theta) &:= -2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p - \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) + \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), P_M \theta \rangle_{\partial \mathcal{T}_h} \\
 &\quad - \langle \boldsymbol{\delta}_h^q \cdot \mathbf{n} + 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p + 2i\omega \tau P_M \delta_h^p - 2i\omega \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), P_M \theta \rangle_{\partial \mathcal{T}_h}.
 \end{aligned}$$

Notice that an integration by parts has been carried out in  $\ell_2$ .

As  $\mathbf{W}_0$  is real and symmetric, we have

$$(\mathbf{W}_0 \boldsymbol{\pi}_V \boldsymbol{\xi}, \boldsymbol{\varepsilon}_h^q)_{\mathcal{T}_h} = (\boldsymbol{\pi}_V \boldsymbol{\xi}, \mathbf{W}_0 \boldsymbol{\varepsilon}_h^q)_{\mathcal{T}_h}$$

and we can therefore notice that the left-hand sides of (73a)–(73b)–(73c) and (74a)–(74b)–(74c) (after being conjugated) are the same, leading to

$$\begin{aligned}
 &(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi} - \boldsymbol{\xi})_{\mathcal{T}_h} + (\varepsilon_h^p, \varepsilon_h^p)_{\mathcal{T}_h} - \omega^2 (\rho_0 \varepsilon_h^p, \pi_W \theta - \theta)_{\mathcal{T}_h} \\
 &- 2i\omega (\nabla \varepsilon_h^p, (\pi_W \theta - \theta) \mathbf{b}_0)_{\mathcal{T}_h} + 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \pi_W \theta - \theta \rangle_{\partial \mathcal{T}_h} \\
 &= \ell_1(\boldsymbol{\pi}_V \boldsymbol{\xi}) + \ell_2(\pi_W \theta) + \ell_3(P_M \theta)
 \end{aligned}$$

And the identity is obtained by a reorganisation of the different terms.  $\blacksquare$

**Remark 4.7:** conjugating gives the good sign for the convection term, indeed :  $[iz]^* = -iz^*$ , therefore

$$\left[ 2i\omega \langle (\pi_W \theta) \mathbf{b}_0, \nabla \varepsilon_h^p \rangle_{\mathcal{T}_h} \right]^* = -2i\omega (\nabla \varepsilon_h^p, (\pi_W \theta) \mathbf{b}_0)_{\mathcal{T}_h}$$

#### Lemma 4.7:

Assuming that the regularity assumption (72) holds and that  $\omega^2 h^2 \|\rho_0\|_\infty C_{\text{reg}} C$  (where  $C$  is the constant of [THEOREM 4](#)) is small enough, if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s \in [1, k+2]$  and  $t \in [1, k+1]$  then

$$\|\varepsilon_h^p\|_{\mathcal{T}_h} \lesssim h^{t+1} (1+\omega) \|\mathbf{q}\|_{t,\mathcal{O}} + h^s (1+\omega+\omega^2) \|p\|_{s,\mathcal{O}} + \omega h \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} + h (1+\omega) \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h}$$

where

$$h := \max_{K \in \mathcal{T}_h} h_K$$

**Proof:** we are going to estimate the terms in the right hand side of the [LEMMA 4.6](#).

Volumetric terms involving  $\boldsymbol{\varepsilon}_h^q$  :

$$(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\xi} - \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} \lesssim h \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, K} \|\boldsymbol{\xi}\|_{1,\mathcal{O}} \lesssim h \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h} \|\varepsilon_h^p\|_{\mathcal{T}_h} \quad (75a)$$

and

$$\begin{aligned}
 &\left| (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} \right| = \left| (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\xi})_{\mathcal{T}_h} - (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\xi} - \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} \right| \\
 &\quad \text{(by (64))} = \left| (\boldsymbol{\delta}_h^q, \mathbf{W}_0 \boldsymbol{\xi} - \{\mathbf{W}_0 \boldsymbol{\xi}\})_{\mathcal{T}_h} - (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\xi} - \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} \right| \\
 &\quad \lesssim \|\mathbf{q} - \boldsymbol{\pi}_V \mathbf{q}\|_{\mathcal{T}_h} \|\mathbf{W}_0 \boldsymbol{\xi} - \{\mathbf{W}_0 \boldsymbol{\xi}\}\|_{\mathcal{T}_h} + \|\mathbf{q} - \boldsymbol{\pi}_V \mathbf{q}\|_{\mathcal{T}_h} \|\boldsymbol{\xi} - \boldsymbol{\pi}_V \boldsymbol{\xi}\|_{\mathcal{T}_h} \\
 &\quad \text{(by (63) and (59b))} \lesssim h^t \|\mathbf{q}\|_{t,\mathcal{O}} h \|\boldsymbol{\xi}\|_{1,\mathcal{O}} + h^t \|\mathbf{q}\|_{t,\mathcal{O}} h \|\boldsymbol{\xi}\|_{1,\mathcal{O}} \\
 &\quad \text{(by regularity (72))} \lesssim h^{t+1} \|\mathbf{q}\|_{t,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h} \quad (75b)
 \end{aligned}$$

Volumetric term involving  $\varepsilon_h^p$  :

$$\omega^2 (\rho_0 \varepsilon_h^p, \pi_W \theta - \theta)_{\mathcal{T}_h} \lesssim \omega^2 h^2 \|\varepsilon_h^p\|_{\rho_0, K} \|\theta\|_{2, \mathcal{O}} \lesssim \omega^2 h^2 \|\varepsilon_h^p\|_{\mathcal{O}} \quad (75c)$$

and

$$\begin{aligned} \left| \omega^2 (\rho_0 \delta_h^p, \pi_W \theta)_{\mathcal{T}_h} \right| &= \left| (\rho_0 \delta_h^p, \theta)_{\mathcal{T}_h} - (\rho_0 \delta_h^p, \theta - \pi_W \theta)_{\mathcal{T}_h} \right| \\ &\stackrel{\text{(by (64))}}{=} \omega^2 \left| (\delta_h^p, \rho_0 \theta - \{\rho_0 \theta\})_{\mathcal{T}_h} - (\rho_0 \delta_h^p, \theta - \pi_W \theta)_{\mathcal{T}_h} \right| \\ &\lesssim \omega^2 \|p - \pi_W p\|_{\mathcal{T}_h} \|\rho_0 \theta - \{\rho_0 \theta\}\|_{\mathcal{T}_h} + \omega^2 \|p - \pi_W p\|_{\mathcal{T}_h} \|\theta - \pi_W \theta\|_{\mathcal{T}_h} \\ &\stackrel{\text{(by (63) and (59a))}}{\lesssim} \omega^2 h^s \|p\|_{s, \mathcal{O}} h \|\theta\|_{1, \mathcal{O}} + \omega^2 h^s \|p\|_{s, \mathcal{O}} h \|\theta\|_{1, \mathcal{O}} \\ &\stackrel{\text{(by regularity (72))}}{\lesssim} \omega^2 h^{s+1} \|p\|_{s, \mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h} \end{aligned} \quad (75d)$$

Volumetric convection term :

$$\begin{aligned} 2\omega (\nabla \varepsilon_h^p, (\pi_W \theta - \theta) \mathbf{b}_0)_{\mathcal{T}_h} &\lesssim \omega \|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} \|\pi_W \theta - \theta\|_{\mathcal{T}_h} \\ &\stackrel{\text{(by LEMMA 4.3)}}{\lesssim} \omega \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \|\delta_h^q\|_{\mathcal{T}_h} + \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \|\pi_W \theta - \theta\|_{\mathcal{T}_h} \\ &\stackrel{\text{(by (59b), (59a), (63) and (72))}}{\lesssim} \omega \left( \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) h \|\varepsilon_h^p\|_{\mathcal{T}_h} \\ &\lesssim \omega \left( h \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^{t+1} \|\mathbf{q}\|_{t, \mathcal{O}} + h \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \|\varepsilon_h^p\|_{\mathcal{T}_h} \end{aligned} \quad (75e)$$

and

$$\begin{aligned} 2\omega \left| (\delta_h^p \mathbf{b}_0, \nabla (\pi_W \theta))_{\mathcal{T}_h} \right| &= 2\omega \left| (\delta_h^p \mathbf{b}_0, \nabla \theta)_{\mathcal{T}_h} - (\delta_h^p \mathbf{b}_0, \nabla (\theta - \pi_W \theta))_{\mathcal{T}_h} \right| \\ &\stackrel{\text{(by (59a))}}{\lesssim} \omega \|p - \pi_W p\|_{\mathcal{T}_h} \|\theta\|_{1, \mathcal{O}} \\ &\stackrel{\text{(by (59a) and (72))}}{\lesssim} \omega h^s \|p\|_{s, \mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h} \end{aligned} \quad (75f)$$

Boundary terms involving  $\mathbf{b}_0 \cdot \mathbf{n}$  : As we want to keep  $s \leq k + 2$ , we will use (59c) with  $s - 1$  instead of  $s$ .

As

$$\langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p, \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} = \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p, \pi_W \theta - \theta \rangle_{\partial \mathcal{T}_h}$$

we focus on

$$\begin{aligned} 2\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta - \theta \rangle_{\partial \mathcal{T}_h} &\lesssim \omega \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} \|\pi_W \theta - \theta\|_{\partial \mathcal{T}_h} \\ &\stackrel{\text{(by (68))}}{\lesssim} \omega h^2 \|\varepsilon_h^p\|_{\mathcal{T}_h} \left[ \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right] \end{aligned} \quad (75g)$$

and

$$\begin{aligned} 2\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p, \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &\lesssim \omega \left\| \widehat{\delta}_h^p \right\|_{\partial \mathcal{T}_h} \|\pi_W \theta - P_M \theta\|_{\partial \mathcal{T}_h} \\ &\stackrel{\text{(by (59d) and (59c))}}{\lesssim} \omega h^{s-\frac{1}{2}} \|p\|_{s, \mathcal{O}} h^{\frac{3}{2}} \|\theta\|_{2, \mathcal{O}} \\ &\stackrel{\text{(by (72))}}{\lesssim} \omega h^s \|p\|_{s, \mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h} \end{aligned} \quad (75h)$$

Boundary terms involving  $\tau$  :

$$\begin{aligned}
 2\omega \langle \tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &\lesssim \omega \|\tau (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p)\|_{\partial \mathcal{T}_h} \|\pi_W \theta - P_M \theta\|_{\partial \mathcal{T}_h} \\
 (\tau = \mathcal{O}(h^{-1})) &\lesssim \omega h^{-\frac{1}{2}} \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} h^{\frac{3}{2}} \|\theta\|_{2,\mathcal{O}} \\
 &\lesssim \omega h \left\| \tau^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \|\varepsilon_h^p\|_{\mathcal{T}_h}
 \end{aligned} \tag{75i}$$

and

$$\begin{aligned}
 2\omega \langle \tau P_M \delta_h^p, \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &= 2\omega \langle \tau P_M \delta_h^p, P_M (\pi_W \theta - \theta) \rangle_{\partial \mathcal{T}_h} \\
 &\lesssim \omega \|\tau P_M \delta_h^p\|_{\partial \mathcal{T}_h} \|P_M (\theta - \pi_W \theta)\|_{\partial \mathcal{T}_h} \\
 (\tau = \mathcal{O}(h^{-1})) &\lesssim \omega h^{-1} \|p - \pi_W p\|_{\partial \mathcal{T}_h} \|\theta - \pi_W \theta\|_{\partial \mathcal{T}_h} \\
 (\text{by (59d)}) &\lesssim \omega h^{-1} h^{s-\frac{1}{2}} \|p\|_{s,\mathcal{O}} h^{2-\frac{1}{2}} \|\theta\|_{2,\mathcal{O}} \\
 &\lesssim \omega h^s \|p\|_{s,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h}
 \end{aligned} \tag{75j}$$

Boundary terms involving  $\tau_{\text{upw}}$ : As we want to keep  $s \leq k + 2$ , we will use (59c) with  $s - 1$  instead of  $s$ .

$$\begin{aligned}
 2\omega \langle \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &\lesssim \omega \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h} \|\pi_W \theta - P_M \theta\|_{\partial \mathcal{T}_h} \\
 &\lesssim \omega h^{s-\frac{3}{2}} \|p\|_{s,\mathcal{O}} h^{\frac{3}{2}} \|\theta\|_{2,\mathcal{O}} \\
 &\lesssim \omega h^s \|p\|_{s,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h}
 \end{aligned} \tag{75k}$$

and

$$\begin{aligned}
 2\omega \langle \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &\lesssim \omega h^s \|p\|_{s,\mathcal{O}} \|\theta\|_{2,\mathcal{O}} \\
 &\lesssim \omega h^s \|p\|_{s,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h}
 \end{aligned} \tag{75l}$$

Boundary term involving  $\delta_h^q$ :

$$\begin{aligned}
 \langle \delta_h^q \cdot \mathbf{n}, \pi_W \theta - P_M \theta \rangle_{\partial \mathcal{T}_h} &\lesssim h^{t-\frac{1}{2}} \|\mathbf{q}\|_{t,\mathcal{O}} h^{\frac{3}{2}} \|\theta\|_{2,\mathcal{O}} \quad (\text{by (59e) and (59d)}) \\
 &\lesssim h^{t+1} \|\mathbf{q}\|_{t,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h}
 \end{aligned} \tag{75m}$$

The desired estimate can now be obtained by collecting the estimates (75a), (75b), (75c), (75d), (75e), (75f), (75g), (75h), (75i), (75j), (75k), (75l) and (75m), and using that  $h \rightarrow 0$ .  $\blacksquare$

### 4.3.3 Bootstrapping process

We will now combine the results of COROLLARY 4.1 of LEMMA 4.5 and LEMMA 4.7 through a bootstrapping process to obtain a convergence result.



**Theorem 5** : *Convergence of the HDG+ method*

Assuming that the regularity assumption (72) holds and that  $\omega^2 h^2 \|\rho_0\|_\infty C_{\text{reg}} C$  (where  $C$  is the constant of [THEOREM 4](#)) is small enough, if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s \in [1, k+2]$  and  $t \in [1, k+1]$  then

$$\begin{aligned} \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{w}_0, \mathcal{T}_h} + \sqrt{2\omega} \left( \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} + \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \right) \\ \lesssim h^t \|\mathbf{q}\|_{t, \mathcal{O}} + h^{s-1} \|p\|_{s, \mathcal{O}} \end{aligned}$$

and

$$\|\varepsilon_h^p\|_{\mathcal{T}_h} \lesssim (1 + \omega) \left( h^{t+1} \|\mathbf{q}\|_{t, \mathcal{O}} + h^s \|p\|_{s, \mathcal{O}} \right)$$

Optimal error estimates are

$$\|\pi_W p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+2}) \quad \text{and} \quad \|\boldsymbol{\pi}_V \mathbf{q} - \mathbf{q}_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+1}),$$

and

$$\|p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+2}) \quad \text{and} \quad \|\mathbf{q} - \mathbf{q}_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+1}).$$

**Proof:** to make the computations easier we introduce the following notations

$$P := \omega \|\varepsilon_h^p\|_{\mathcal{T}_h} \quad ; \quad Q := \|\boldsymbol{\varepsilon}_h^{\mathbf{q}}\|_{\mathbf{w}_0, \mathcal{T}_h} \quad ; \quad T := \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h}$$

and

$$\begin{aligned} P_h &:= \omega h^{s-1} \|p\|_{s, \mathcal{O}} \quad ; \quad Q_h := h^t \|\mathbf{q}\|_{t, \mathcal{O}} \\ B &:= \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} + \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial \mathcal{T}_h} \\ \alpha &:= (\omega + \omega^2) h \quad ; \quad \beta := (1 + \omega + \omega^2) h \quad ; \quad \gamma := \omega^2 h \end{aligned}$$

With those shorthands the estimate of [LEMMA 4.7](#) can be rewritten as

$$\begin{aligned} P &\lesssim (\omega + \omega^2) h Q_h + (1 + \omega + \omega^2) h P_h + \omega^2 h T + (\omega + \omega^2) h Q \\ &\lesssim \alpha Q_h + \beta P_h + \gamma B + \alpha Q \end{aligned} \tag{76}$$

and the estimate of [COROLLARY 4.1](#) can be rewritten as

$$\begin{aligned} Q^2 + 2\omega B^2 &\lesssim |Q^2 - 2i\omega B^2| \lesssim \varepsilon (Q^2 + T^2) + \frac{1}{\varepsilon} (P^2 + Q_h^2 + P_h^2) \\ &\lesssim \varepsilon (Q^2 + B^2) + \frac{1}{\varepsilon} (P^2 + Q_h^2 + P_h^2) \end{aligned} \tag{77}$$

Taking the square of (76) and using Young's inequality leads to

$$P^2 \lesssim \alpha^2 Q_h^2 + \beta^2 P_h^2 + \gamma^2 B^2 + \alpha^2 Q^2. \tag{78}$$

Using (78) in (77) gives

$$Q^2 + 2\omega B^2 \lesssim \left(1 + \frac{1}{\varepsilon}\right) \left[ (1 + \alpha^2) Q_h^2 + (1 + \beta^2) P_h^2 \right] + \left( \varepsilon + \gamma^2 \left(1 + \frac{1}{\varepsilon}\right) \right) B^2 + \left( \varepsilon + \alpha^2 \left(1 + \frac{1}{\varepsilon}\right) \right) Q^2.$$

Let  $C_1$  denote the constant hidden in  $\lesssim$ . Choosing  $\varepsilon$  so that  $C_1\varepsilon < 1$  and assuming that  $h$  is small enough the last two terms of the right hand side can be absorbed by the left hand side, leading to

$$Q^2 + 2\omega B^2 \lesssim \left(1 + \frac{1}{\varepsilon}\right) \left[ (1 + \alpha^2) Q_h^2 + (1 + \beta^2) P_h^2 \right].$$

As  $\varepsilon$  does not depend on  $\omega$  and  $h$ , we can hide the first factor of the right-hand side into  $\lesssim$  leading to

$$Q^2 + 2\omega B^2 \lesssim (1 + \alpha^2) Q_h^2 + (1 + \beta^2) P_h^2.$$

As  $\alpha, \beta = \mathcal{O}(h)$ , we can overestimate  $\alpha, \beta \lesssim 1$  leading to

$$Q^2 + 2\omega B^2 \lesssim Q_h^2 + P_h^2.$$

And finally we have

$$Q + \sqrt{2\omega}B \lesssim P_h + Q_h. \quad (79)$$

Now by taking  $s = k + 2$  and  $t = k + 1$  in (59b) and (59d), we can see that

$$Q = \mathcal{O}(h^{k+1}) \quad \text{and} \quad B = \mathcal{O}(h^{k+1}) \quad (80)$$

and finally by using (79) in (76), we have

$$P = \mathcal{O}(h^{k+2}). \quad (81)$$

It is also possible to obtain a convergence result for the trace  $\widehat{p}_h$  which is the main unknown of the method :

#### Corollary 4.2:

Under the assumptions of [THEOREM 5](#), the following error estimates for  $\widehat{p}_h$  hold

$$\|\widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{3}{2}}) \quad \text{and} \quad \|p - \widehat{p}_h\|_{\partial\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{1}{2}}).$$

**Proof:** We have

$$\begin{aligned} \|\varepsilon_h^p\|_{\partial\mathcal{T}_h} &\leq \|P_M \varepsilon_h^p\|_{\partial\mathcal{T}_h} + \|P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h} \\ (\tau = \mathcal{O}(h^{-1})) &\lesssim \|P_M \varepsilon_h^p\|_{\partial\mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \\ (P_M \text{ continuous}) &\lesssim \|\varepsilon_h^p\|_{\partial\mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \\ (\text{by (60)}) &\lesssim h^{-\frac{1}{2}} \|\varepsilon_h^p\|_{\mathcal{T}_h} + h^{\frac{1}{2}} \left\| |\tau|^{\frac{1}{2}} (P_M \varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h} \\ (\text{by (80) and (81)}) &= \mathcal{O}(h^{-\frac{1}{2}} h^{k+2} + h^{\frac{1}{2}} h^{k+1}) \\ &= \mathcal{O}(h^{k+\frac{3}{2}}), \end{aligned}$$

which is the first estimate. The second one comes from (59c) with  $s = k + 1$ . ■

## 4.4 Error analysis of the HDG method with diffusive flux

As the error analysis for the HDG method is very similar to the one for the HDG+ method we only state the main theorem. The intermediate results (error equations, Garding's equality, dual estimate, ...) are stated without proof in [APPENDIX A](#).

**Theorem 6** : *Convergence of the HDG method with diffusive flux*

Assuming that the regularity assumption (72) holds and that  $\omega^2 h^2 \|\rho_0\|_\infty C_{\text{reg}} C$  (where  $C$  is the constant of THEOREM 4) is small enough, if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s, t \in [1, k+1]$  then

$$\|\pi_W p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{3}{2}}) ; \quad \|P_M p - \widehat{p}_h\|_{\partial\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{1}{2}}) ; \quad \text{and} \quad \|\pi_V \mathbf{q} - \mathbf{q}_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{1}{2}}),$$

and

$$\|p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+1}) ; \quad \|p - \widehat{p}_h\|_{\partial\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{1}{2}}) ; \quad \text{and} \quad \|\mathbf{q} - \mathbf{q}_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+\frac{1}{2}}).$$

**Remark 4.8:** Some HDG methods are known to achieve superconvergence, *ie* taking  $p_h \in \mathcal{P}_k$  leads to the following error estimate

$$\|\pi_W p - p_h\|_{\mathcal{T}_h} = \mathcal{O}(h^{k+2}),$$

this is not possible for this method because of the convection term

$$2\omega \left| (\delta_h^p \mathbf{b}_0, \nabla(\pi_W \theta))_{\mathcal{T}_h} \right| \lesssim \omega h^s \|p\|_{s,\mathcal{O}} \|\varepsilon_h^p\|_{\mathcal{T}_h}$$

in the dual estimate which locks the convergence rate to  $\mathcal{O}(h^{k+1})$  for the scalar variable  $p_h$ . Superconvergence is an attractive property for a numerical scheme, indeed using a postprocessing scheme it is possible to use the solution  $(p_h, \mathbf{q}_h)$  to construct a new approximation  $\widetilde{p}_h$  which converges with order  $\mathcal{O}(h^{k+2})$ , see [Ste91], [CGS10, Sec. 5] for more details.

## 5 Implementation

In this section, we will give details on how the HDG(+) methods can be implemented. For the sake of simplicity, we will assume that the physical coefficients are constant on each element. In this case, the integrals may be computed using

- an analytic integration procedure that relies on the decomposition of the basis functions in the monomial basis,
- high-order quadrature rules.

Taking into account the variations of the physical parameters inside the elements only requires a straightforward generalization of the material presented here. However this implementation is only possible if quadrature rules are used to evaluate the integrals.

To make the notations lighter, we will drop the subscript  $h$  in this section. The quantities  $(\mathbf{q}, p, \widehat{p})$  will denote the solutions of (46a)–(46b)–(46c) and  $(\mathbf{q}^K, p^K, \widehat{p}^e)$  their restrictions to an element  $K \in \mathcal{T}_h$  and an edge  $e \in \mathcal{E}_h$  respectively.

In TABLE 2, we recall the main differences between the three variants of the HDG methods considered in this work.

Variable	Space	HDG $(p_h, \boldsymbol{\sigma}_h)$	HDG $(p_h, \mathbf{q}_h)$	HDG+ $(p_h, \mathbf{q}_h)$
Pressure $p_h$	$W_h(K)$	$\mathcal{P}_k(K)$	$\mathcal{P}_k(K)$	$\mathcal{P}_{k+1}(K)$
Flux $\mathbf{q}_h$ or $\boldsymbol{\sigma}_h$	$\mathbf{V}_h(K)$	$\mathcal{P}_k(K)$	$\mathcal{P}_k(K)$	$\mathcal{P}_k(K)$
Trace $\widehat{p}_h$	$M_h(e)$	$\mathcal{P}_k(e)$	$\mathcal{P}_k(e)$	$\mathcal{P}_k(e)$
Penalization $\tau _K$		$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(h_K^{-1})$

Table 2: Choice of local spaces and penalization parameter for the different methods

## 5.1 Framework and notations

In this section we present the notations used in [BCDL15, FS20] which are very close to the ones used in hawen<sup>6</sup>. We will then give details on how the matrices are assembled.

**Local problem:** The discretization of the local problem (46a)–(46b) leads to the following system

$$\mathbb{A}^K \underline{W}^K + \mathbb{C}^K \underline{\Lambda}^K = \mathcal{S}^K, \quad (82)$$

where

$$\underline{W}^K := \begin{bmatrix} p^K & q_x^K & q_y^K & q_z^K \end{bmatrix}^T \quad \text{and} \quad \underline{\Lambda}^K := \begin{bmatrix} \widehat{p}^{g(K,1)} & \widehat{p}^{g(K,2)} & \widehat{p}^{g(K,3)} & \widehat{p}^{g(K,4)} \end{bmatrix}^T,$$

and  $\underline{p}^K$  denotes the vector of the coefficients of  $p_h^K$  in the basis of  $W_h(K)$ .

**Global problem:** The discretization of the transmission condition (46c) leads to

$$\sum_{K \in \mathcal{T}_h} \mathbb{B}^K \underline{W}^K + \mathbb{L}^K \underline{\Lambda}^K = 0 \quad (83)$$

We denote by  $\underline{\Lambda}$  the global trace approximation and for each element  $K \in \mathcal{T}_h$  we introduce the *connectivity map* which is the operator  $\mathcal{A}_K$  such that<sup>7</sup>

$$\mathcal{A}_K \underline{\Lambda} = \underline{\Lambda}^K.$$

Using (82), we can express  $\underline{W}^K$  in terms of  $\underline{\Lambda}$

$$\underline{W}^K = \left(\mathbb{A}^K\right)^{-1} \mathcal{S}^K - \left(\mathbb{A}^K\right)^{-1} \mathbb{C}^K \mathcal{A}_K \underline{\Lambda}$$

and we can finally construct the global problem using (83)

$$\sum_{K \in \mathcal{T}_h} \mathcal{A}_K^T \left[ -\mathbb{B}^K \left(\mathbb{A}^K\right)^{-1} \mathbb{C}^K + \mathbb{L}^K \right] \mathcal{A}_K \underline{\Lambda} = - \sum_{K \in \mathcal{T}_h} \mathcal{A}_K^T \mathbb{B}^K \left(\mathbb{A}^K\right)^{-1} \mathcal{S}^K. \quad (84)$$

For conciseness, we will denote

$$\mathbb{K}^K := -\mathbb{B}^K \left(\mathbb{A}^K\right)^{-1} \mathbb{C}^K + \mathbb{L}^K, \quad \mathbb{K} := \sum_{K \in \mathcal{T}_h} \mathcal{A}_K^T \mathbb{K}^K \mathcal{A}_K \quad \text{and} \quad \mathcal{S} := - \sum_{K \in \mathcal{T}_h} \mathcal{A}_K^T \mathbb{B}^K \left(\mathbb{A}^K\right)^{-1} \mathcal{S}^K.$$

**Remark 5.1:** the invertibility of matrix  $\mathbb{A}^K$  was the object of [THEOREM 1](#) (HDG- $\sigma_h$ ) and [THEOREM 4](#) (HDG+ and HDG- $\mathbf{q}_h$ ) and the invertibility of the global matrix  $\mathbb{K}$  was the object of [THEOREM 3](#).

**Remark 5.2:** to efficiently implement the HDG(+) methods, we will need to compute the *inverse* of  $\mathbb{A}^K$ . As this matrix is not too large, it is possible to perform this operation using `lapack`. However the discretization of the global problem will lead to a large sparse system which will be solved using `mumps`.

<sup>6</sup>See <https://ffaucher.gitlab.io/hawen-website/>.

<sup>7</sup>The aim of operator  $\mathcal{A}_K$  is to copy the global informations to the local solver, depending on the element it may be a simple copy or also involve reordering.

**Solving HDG:** In [ALGORITHM 3](#), we describe how the HDG method can be efficiently implemented. Details on the construction of the local matrices  $\mathbb{A}^K$ ,  $\mathbb{C}^K$ ,  $\mathbb{B}^K$  and  $\mathbb{L}^K$  and on the implementation of boundary conditions will be given in the next sections.

---

**Algorithm 3:** Solving HDG

---

```

/* Step 1: Construction of the local and global problems */
1 for  $K \in \mathcal{T}_h$  do
2   Construct the local matrices  $\mathbb{A}^K$ ,  $\mathbb{C}^K$ ,  $\mathbb{B}^K$  and  $\mathbb{L}^K$ 
3   Compute  $(\mathbb{A}^K)^{-1}$  using lapack
4   Compute  $\mathbb{K}^K = \mathbb{B}^K (\mathbb{A}^K)^{-1} \mathbb{C}^K + \mathbb{L}^K$ 
5   Modify  $\mathbb{K}^K$  to enforce the boundary conditions
6   Use the connectivity operator  $\mathcal{A}_K$  to add the local contribution to the global matrix  $\mathbb{K}$ 
/* Step 2: Construction of the source term */
7 Localize the source
8 for  $K \in \mathcal{T}_h$  where  $s \neq 0$  do
9   Construct the local source term  $\mathcal{S}^K$ 
10  Compute  $-\mathbb{B}^K (\mathbb{A}^K)^{-1} \mathcal{S}^K$ 
11  Use the connectivity operator  $\mathcal{A}_K$  to add the local contribution to the global source  $\mathcal{S}$ 
/* Step 3: Resolution of the global system */
12 Solve  $\mathbb{K}\underline{\Lambda} = \mathcal{S}$  with mumps
/* Step 4: Reconstruction of the solution */
13 for  $K \in \mathcal{T}_h$  do
14  Compute the local unknowns  $\underline{W}^K = (\mathbb{A}^K)^{-1} \mathcal{S}^K - (\mathbb{A}^K)^{-1} \mathbb{C}^K \mathcal{A}_K \underline{\Lambda}$ 

```

---

**Convention for the indices:** we will always use the following convention for the indices :

- index of basis functions :  $i$  (test),  $j$  (trial),  $r$  (other)
- component of a vector :  $u, v \in \{x, y, z\}$
- global number of an edge :  $m$
- local number of an edge :  $\ell$
- polynomial degree :  $k$

The global number of the  $\ell$ -th face of element  $K$  will be denoted by  $g(K, \ell)$ .

**Dimension of polynomial spaces:** We recall that the dimension of the space of polynomials of  $n$  variables and of degree up to  $k$  is given by

$$\dim(\mathcal{P}_k(K)) = \binom{k+n}{n} =: d_n(k), \quad \text{for } K \subset \mathbb{R}^n.$$

## 5.2 Implementation of the diffusive flux HDG method

In this section, we focus on the implementation of the HDG- $\mathbf{q}_h$  method.

### 5.2.1 Definition of the elementary matrices

We introduce the basis  $(\Phi_j^K)_{j=1}^{d_n(k)}$  of  $\mathcal{P}_k(K)$  and  $(\Psi_j^e)_{j=1}^{d_{n-1}(k)}$  of  $\mathcal{P}_k(e)$  and we decompose the unknowns as

$$p^K = \sum_{j=1}^{d_n(k)} p_j^K \Phi_j^K; \quad q_u^K = \sum_{j=1}^{d_n(k)} q_{u,j}^K \Phi_j^K; \quad \widehat{p}^e(x) = \sum_{j=1}^{d_{n-1}(k)} \widehat{p}_j^e \Psi_j^e$$

for  $u \in \{x, y, z\}$ ,  $K \in \mathcal{T}_h$  and  $e \in \mathcal{E}_h$ .

We introduce the following matrices

$$\begin{aligned} \mathbb{M}_{ij}^K &= \int_K \Phi_i^K \Phi_j^K d\mathbf{x} & \mathbb{D}_{u,ij}^K &= \int_K \Phi_j^K \partial_u \Phi_i^K d\mathbf{x} \\ \mathbb{E}_{\ell,ij}^K &= \int_{\partial K^\ell} \Phi_i^K \Phi_j^K d\sigma & \mathbb{F}_{\ell,ij}^K &= \int_{\partial K^\ell} \Psi_j^{g(K,\ell)} \Phi_i^K d\sigma \\ \mathbb{G}_{ij}^m &= \int_{e_m} \Psi_i^m \Psi_j^m d\sigma \end{aligned} \quad (85)$$

where  $g(K, \ell)$  is the global number of the  $\ell$ -th face of element  $K$  and  $m$  is also the global number of the edge  $e_m$ .

Notice that all those matrices except  $\mathbb{F}_\ell^K$  (which has dimension  $d_n(k) \times d_{n-1}(k)$ ) are square matrices and that all the elementary matrices are real.

## 5.2.2 Local problem

We introduce the following shorthand to make the algebra easier

$$\alpha_\ell = \min(\mathbf{b}_0 \cdot \mathbf{n}, 0) \quad \text{and} \quad t_{\text{upw},\ell} = \max_{\partial K^\ell}(\mathbf{b}_0 \cdot \mathbf{n}, 0)$$

it corresponds to using the expression (44a)–(44b) of the numerical flux. From an implementation point of view, this can be evaluated using array slicing or a ternary operator.

Using the elementary matrices introduced in (85) we can write (46a) as

$$\sum_u W_{0,vu} \mathbb{M}_u^K \underline{q}_u^K - \mathbb{D}_v^K \underline{p}^K + \sum_\ell n_v^{K,\ell} \mathbb{F}_\ell^K \widehat{\underline{p}}^{g(K,\ell)} = 0, \quad \forall v \in \{x, y, z\} \quad (86)$$

and (46b) as

$$\begin{aligned} & -\omega^2 \rho_0 \mathbb{M}^K \underline{p}^K + 2i\omega \sum_u b_{0,u} \mathbb{D}_u^K \underline{p}^K + \sum_u (\mathbb{D}_u^K)^T \underline{q}_u^K \\ & + 2i\omega \sum_\ell \left[ -\alpha_\ell \mathbb{F}_\ell^K \widehat{\underline{p}}^{g(K,\ell)} - t_{\text{upw},\ell} \mathbb{E}_\ell^K \underline{p}^K + \tau_\ell (\mathbb{E}_\ell^K \underline{p}^K - \mathbb{F}_\ell^K \widehat{\underline{p}}^{g(K,\ell)}) \right] = \mathbb{S}^K \end{aligned} \quad (87)$$

**Matrix form:** To construct the global problem (84) we need to construct the matrices  $\mathbb{A}^K$  and  $\mathbb{C}^K$  of the local problem (82) using (86) and (87).

Matrix  $\mathbb{A}^K$ :

$\mathbb{A}^K$	$\underline{p}^K$	$\underline{q}_x^K$	$\underline{q}_y^K$	$\underline{q}_z^K$
$\underline{p}^K$	$-\rho_0 \omega^2 \mathbb{M}^K + 2i\omega \sum_u b_{0,u} \mathbb{D}_u^K$	$(\mathbb{D}_x^K)^T$	$(\mathbb{D}_y^K)^T$	$(\mathbb{D}_z^K)^T$
$\underline{q}_x^K$	$+2i\omega \sum_\ell (\tau_\ell - t_{\text{upw},\ell}) \mathbb{E}_\ell^K$	$-\mathbb{D}_x^K + \sum_\ell n_x^{K,\ell} \mathbb{E}_\ell^K$	$-\mathbb{D}_y^K + \sum_\ell n_y^{K,\ell} \mathbb{E}_\ell^K$	$-\mathbb{D}_z^K + \sum_\ell n_z^{K,\ell} \mathbb{E}_\ell^K$
$\underline{q}_y^K$	$-\mathbb{D}_x^K$	$\overline{W}_{0,11}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,12}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,13}^K \overline{\mathbb{M}}^K$
$\underline{q}_z^K$	$-\mathbb{D}_y^K$	$\overline{W}_{0,21}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,22}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,23}^K \overline{\mathbb{M}}^K$
$\underline{q}_z^K$	$-\mathbb{D}_z^K$	$\overline{W}_{0,31}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,32}^K \overline{\mathbb{M}}^K$	$\overline{W}_{0,33}^K \overline{\mathbb{M}}^K$

The red terms can be used instead of the black ones, depending on whether or not an integration by parts is performed.

Matrix  $\mathbb{C}^K$ :

$\mathbb{C}^K$	$\widehat{\underline{p}}^{g(K,1)}$	$\widehat{\underline{p}}^{g(K,2)}$	$\widehat{\underline{p}}^{g(K,3)}$	$\widehat{\underline{p}}^{g(K,4)}$
$\underline{p}^K$	$-2i\omega(\tau_1 + \alpha_1) \mathbb{F}_1^K$	$-2i\omega(\tau_2 + \alpha_2) \mathbb{F}_2^K$	$-2i\omega(\tau_3 + \alpha_3) \mathbb{F}_3^K$	$-2i\omega(\tau_4 + \alpha_4) \mathbb{F}_4^K$
$\underline{q}_x^K$	$n_x^{K,1} \mathbb{F}_1^K$	$n_x^{K,2} \mathbb{F}_2^K$	$n_x^{K,3} \mathbb{F}_3^K$	$n_x^{K,4} \mathbb{F}_4^K$
$\underline{q}_y^K$	$n_y^{K,1} \mathbb{F}_1^K$	$n_y^{K,2} \mathbb{F}_2^K$	$n_y^{K,3} \mathbb{F}_3^K$	$n_y^{K,4} \mathbb{F}_4^K$
$\underline{q}_z^K$	$n_z^{K,1} \mathbb{F}_1^K$	$n_z^{K,2} \mathbb{F}_2^K$	$n_z^{K,3} \mathbb{F}_3^K$	$n_z^{K,4} \mathbb{F}_4^K$

**Remark 5.3:** We chose to write the matrix form of the method with the equation for  $p^K$  first which is the opposite of the system (46a)–(46b)–(46c). This convention is interesting because the matrices for the method in dimension  $n$  are submatrices of the ones for dimension  $n + 1$ .

### 5.2.3 Global problem

The transmission condition (46c) can be written as

$$\sum_{K,\ell} \left( \sum_u n_u^{K,\ell} (\mathbb{F}_\ell^K)^T \underline{q}_u^K + 2i\omega \left[ -(\alpha_\ell \mathbb{G}^{g(K,\ell)} \widehat{p}^{g(K,\ell)} + t_{\text{upw},\ell} (\mathbb{F}_\ell^K)^T \underline{p}^K) + \tau \left( (\mathbb{F}_\ell^K)^T \underline{p}^K - \mathbb{G}^{g(K,\ell)} \widehat{p}^{g(K,\ell)} \right) \right] \right) = 0 \quad (88)$$

**Matrix form:** We will now use (88) to construct the matrices  $\mathbb{B}^K$  and  $\mathbb{L}^K$  of (83).

Matrix  $\mathbb{B}^K$ :

$$\mathbb{B}^K \parallel \begin{array}{c|ccc|c} & \underline{p}^K & \underline{q}_x^K & \underline{q}_y^K & \underline{q}_z^K \\ \hline e_1 & 2i\omega(\tau_1 - t_{\text{upw},1}) (\mathbb{F}_1^K)^T & n_x^{K,1} (\mathbb{F}_1^K)^T & n_y^{K,1} (\mathbb{F}_1^K)^T & n_z^{K,1} (\mathbb{F}_1^K)^T \\ e_2 & 2i\omega(\tau_2 - t_{\text{upw},2}) (\mathbb{F}_2^K)^T & n_x^{K,2} (\mathbb{F}_2^K)^T & n_y^{K,2} (\mathbb{F}_2^K)^T & n_z^{K,2} (\mathbb{F}_2^K)^T \\ e_3 & 2i\omega(\tau_3 - t_{\text{upw},3}) (\mathbb{F}_3^K)^T & n_x^{K,3} (\mathbb{F}_3^K)^T & n_y^{K,3} (\mathbb{F}_3^K)^T & n_z^{K,3} (\mathbb{F}_3^K)^T \\ e_4 & 2i\omega(\tau_4 - t_{\text{upw},4}) (\mathbb{F}_4^K)^T & n_x^{K,4} (\mathbb{F}_4^K)^T & n_y^{K,4} (\mathbb{F}_4^K)^T & n_z^{K,4} (\mathbb{F}_4^K)^T \end{array}$$

Matrix  $\mathbb{L}^K$ :

$$\mathbb{L}^K \parallel \begin{array}{c|cccc} & \widehat{p}^{g(K,1)} & \widehat{p}^{g(K,2)} & \widehat{p}^{g(K,3)} & \widehat{p}^{g(K,4)} \\ \hline e_1 & -2i\omega(\tau_1 + \alpha_1) \mathbb{G}^{g(K,1)} & & & \\ e_2 & & -2i\omega(\tau_2 + \alpha_2) \mathbb{G}^{g(K,2)} & & \\ e_3 & & & -2i\omega(\tau_3 + \alpha_3) \mathbb{G}^{g(K,3)} & \\ e_4 & & & & -2i\omega(\tau_4 + \alpha_4) \mathbb{G}^{g(K,4)} \end{array}$$

### 5.2.4 Boundary conditions:

The transmission condition (46c) was also used to weakly enforce the boundary conditions

$$\begin{aligned} \mathbf{q} \cdot \mathbf{n} - 2i\omega(\mathbf{b}_0 \cdot \mathbf{n})p &= g_N && \text{on } \Gamma_N \\ p &= g_D && \text{on } \Gamma_D \end{aligned}$$

**Neumann boundary condition:** If we want to enforce the Neumann boundary condition on edge  $\ell$  of element  $K$  we only need to add a right-hand side to (88) :

$$\mathbb{G}^{g(K,\ell)} \underline{g}_N$$

where we assumed that

$$g_N = \sum_i g_{N,i} \Psi_i.$$

**Weak Dirichlet boundary conditions:** However if we want to enforce the Dirichlet boundary condition on edge  $\ell$  of element  $K$  we need to change the matrices  $\mathbb{B}^K$  and  $\mathbb{L}^K$ . This first method of implementation corresponds to weakly enforcing the boundary conditions, *ie* using the expression

$$\langle \widehat{p}_h - g_D, \mu \rangle_{\Gamma_D} = 0.$$

This method is used for a *modal* implementation of the HDG method. It can be implemented as follows described in [ALGORITHM 4](#).

---

**Algorithm 4:** Implementation of weak Dirichlet BC
 

---

- 1 The corresponding row of  $\mathbb{B}^K$  is set to 0 :  $\mathbb{B}^K[e_\ell, :] = 0$
  - 2 The corresponding entry of  $\mathbb{L}^K$  is changed :  $\mathbb{L}^K[e_\ell, \widehat{p}^{g(K,\ell)}] = \mathbb{G}^{g(K,\ell)}$
  - 3 The following right-hand side is added :  $\mathbb{G}^{g(K,\ell)} \underline{g}_D$ , where  $\underline{g}_D = \sum_i g_{D,i} \Psi_i$ .
- 

**Strong Dirichlet boundary conditions:** For a nodal implementation of the HDG method, it is also possible to strongly enforce the boundary condition on  $\Gamma_D$  :

$$\widehat{p}_h = g_D,$$

on edge  $\ell$  of element  $K$ . We need to change the local contribution

$$\mathbb{K}^K = \mathbb{L}^K - \mathbb{B}^K (\mathbb{A}^K)^{-1} \mathbb{C}^K,$$

to the global problem as described in [ALGORITHM 5](#). We recall that  $\mathbb{K}^K$  has the same shape as  $\mathbb{L}^K$ .

---

**Algorithm 5:** Implementation of strong Dirichlet BC
 

---

- 1 The corresponding entries of  $\mathbb{K}^K$  are replaced with an identity block :  $\mathbb{K}^K[e_\ell, \widehat{p}^{g(K,\ell)}] = \mathbf{Id}$
  - 2 The following right-hand side is added :  $[g_D(\mathbf{x}_r)]_r^T$ , where  $\mathbf{x}_r \in \text{dof}(e^{g(K,\ell)})$
- 

### 5.3 Implementation of the total flux HDG method

Using the same notations as in the previous section, we can write the discrete form of the *total flux* HDG method (15a)–(15b)–(15c).

#### 5.3.1 Local problem

The local problem (12a)–(12b) can be written as

$$\begin{aligned} \sum_u W_{0,vu} \underline{\sigma}_u^K - \mathbb{D}_v^K \underline{p}^K + 2i\omega \left( \sum_u W_{0,vu} b_{0,u} \right) \mathbb{M}^K \underline{p}^K + \sum_\ell n_v^\ell \mathbb{F}_\ell^K \widehat{p}^{g(K,\ell)} &= 0, \quad \forall v \in \{x, y, e\} \\ -\omega^2 \rho_0 \mathbb{M}^K \underline{p}^K + \sum_u (\mathbb{D}_u^K)^T \underline{\sigma}_u^K + i\omega \sum_\ell \tau^\ell \left( \mathbb{E}_\ell^K \underline{p}^K - \mathbb{F}_\ell^K \widehat{p}^{g(K,\ell)} \right) &= \mathbb{S}^K \end{aligned}$$

**Matrix form:** Matrix  $\mathbb{A}^K$  : We introduce  $\beta_v := \sum_u W_{0,vu} b_{0,u}$ .

$$\begin{array}{c|cccc} \mathbb{A}^K & \underline{p}^K & \underline{\sigma}_x^K & \underline{\sigma}_y^K & \underline{\sigma}_z^K \\ \hline \underline{p}^K & -\rho_0 \omega^2 \mathbb{M}^K + i\omega \sum_\ell \tau_\ell \mathbb{E}_\ell^K & (\mathbb{D}_x^K)^T & (\mathbb{D}_y^K)^T & (\mathbb{D}_z^K)^T \\ \hline \underline{\sigma}_x^K & -\mathbb{D}_x^K + 2i\omega \beta_x \mathbb{M}^K & \overline{W}_{0,11}^K \overline{\mathbb{M}}^K & \overline{W}_{0,12}^K \overline{\mathbb{M}}^K & \overline{W}_{0,13}^K \overline{\mathbb{M}}^K \\ \underline{\sigma}_y^K & -\mathbb{D}_y^K + 2i\omega \beta_y \mathbb{M}^K & \overline{W}_{0,21}^K \overline{\mathbb{M}}^K & \overline{W}_{0,22}^K \overline{\mathbb{M}}^K & \overline{W}_{0,23}^K \overline{\mathbb{M}}^K \\ \underline{\sigma}_z^K & -\mathbb{D}_z^K + 2i\omega \beta_z \mathbb{M}^K & \overline{W}_{0,31}^K \overline{\mathbb{M}}^K & \overline{W}_{0,32}^K \overline{\mathbb{M}}^K & \overline{W}_{0,33}^K \overline{\mathbb{M}}^K \end{array}$$

Matrix  $\mathbb{C}^K$ :

$$\begin{array}{c|cccc} \mathbb{C}^K & \widehat{p}^{g(K,1)} & \widehat{p}^{g(K,2)} & \widehat{p}^{g(K,3)} & \widehat{p}^{g(K,4)} \\ \hline \underline{p}^K & -i\omega \tau_1 \mathbb{F}_1^K & -i\omega \tau_2 \mathbb{F}_2^K & -i\omega \tau_3 \mathbb{F}_3^K & -i\omega \tau_4 \mathbb{F}_4^K \\ \underline{\sigma}_x^K & n_x^{K,1} \overline{\mathbb{F}}_1^K & n_x^{K,2} \overline{\mathbb{F}}_2^K & n_x^{K,3} \overline{\mathbb{F}}_3^K & n_x^{K,4} \overline{\mathbb{F}}_4^K \\ \underline{\sigma}_y^K & n_y^{K,1} \overline{\mathbb{F}}_1^K & n_y^{K,2} \overline{\mathbb{F}}_2^K & n_y^{K,3} \overline{\mathbb{F}}_3^K & n_y^{K,4} \overline{\mathbb{F}}_4^K \\ \underline{\sigma}_z^K & n_z^{K,1} \overline{\mathbb{F}}_1^K & n_z^{K,2} \overline{\mathbb{F}}_2^K & n_z^{K,3} \overline{\mathbb{F}}_3^K & n_z^{K,4} \overline{\mathbb{F}}_4^K \end{array}$$



### 5.3.2 Global problem

The discrete transmission condition (14) can be written

$$\sum_{K,\ell} \left[ n_u^{K,\ell} (\mathbb{F}_\ell^K)^T \underline{\sigma}_u^K + i\omega\tau^\ell \left( (\mathbb{F}_\ell^K) \underline{p}^K - \mathbb{G}^{g(K,\ell)} \widehat{\underline{p}}^{g(K,\ell)} \right) \right] = 0$$

**Matrix form:** Matrix  $\mathbb{B}^K$ :

$$\mathbb{B}^K \parallel \begin{array}{c|ccc|c} & \underline{p}^K & \underline{\sigma}_x^K & \underline{\sigma}_y^K & \underline{\sigma}_z^K \\ \hline e_1 & i\omega\tau_1 (\mathbb{F}_1^K)^T & n_x^{K,1} (\mathbb{F}_1^K)^T & n_y^{K,1} (\mathbb{F}_1^K)^T & n_z^{K,1} (\mathbb{F}_1^K)^T \\ e_2 & i\omega\tau_2 (\mathbb{F}_2^K)^T & n_x^{K,2} (\mathbb{F}_2^K)^T & n_y^{K,2} (\mathbb{F}_2^K)^T & n_z^{K,2} (\mathbb{F}_2^K)^T \\ e_3 & i\omega\tau_3 (\mathbb{F}_3^K)^T & n_x^{K,3} (\mathbb{F}_3^K)^T & n_y^{K,3} (\mathbb{F}_3^K)^T & n_z^{K,3} (\mathbb{F}_3^K)^T \\ e_4 & i\omega\tau_4 (\mathbb{F}_4^K)^T & n_x^{K,4} (\mathbb{F}_4^K)^T & n_y^{K,4} (\mathbb{F}_4^K)^T & n_z^{K,4} (\mathbb{F}_4^K)^T \end{array}$$

Matrix  $\mathbb{L}^K$ :

$$\mathbb{L}^K \parallel \begin{array}{c|ccc|c} & \widehat{\underline{p}}^{g(K,1)} & \widehat{\underline{p}}^{g(K,2)} & \widehat{\underline{p}}^{g(K,3)} & \widehat{\underline{p}}^{g(K,4)} \\ \hline e_1 & -i\omega\tau_1 \mathbb{G}^{g(K,1)} & & & \\ e_2 & & -i\omega\tau_2 \mathbb{G}^{g(K,2)} & & \\ e_3 & & & -i\omega\tau_3 \mathbb{G}^{g(K,3)} & \\ e_4 & & & & -i\omega\tau_4 \mathbb{G}^{g(K,4)} \end{array}$$

## 5.4 Implementation of the HDG+ method

As  $p^K \in \mathcal{P}_{k+1}(K)$ ,  $\mathbf{q}^K \in \mathcal{P}_k(K)$  and  $\widehat{p}^e \in \mathcal{P}_k(e)$ , one should be very careful while writing the matrix form of the system, indeed we have

$$\underline{p}^K \in \mathbb{C}^{d_n(k+1)} \quad \text{and} \quad \underline{\mathbf{q}}^K \in \mathbb{C}^{n \times d_n(k)} \quad (\iff \underline{q}_u^K \in \mathbb{C}^{d_n(k)}, u \in \{x, y, z\})$$

so some of the elementary matrices will be rectangular.

### 5.4.1 Definition of the elementary matrices

Let  $(\Phi_j^{K,k})_{j=1}^{d_n(k)}$  be the basis for  $\mathcal{P}_k(K)$  and  $(\Psi_j^{e,k})_{j=1}^{d_{n-1}(k)}$  be the basis for  $\mathcal{P}_k(e)$ .

The unknowns are therefore decomposed in the following way

$$p^K = \sum_{j=1}^{d_n(k+1)} p_j^K \Phi_j^{K,k+1} ; \quad q_u^K = \sum_{j=1}^{d_n(k)} q_{u,j}^K \Phi_j^{K,k} ; \quad \widehat{p}^e = \sum_{j=1}^{d_{n-1}(k)} \widehat{p}_j^e \Psi_j^{e,k}$$

for  $u \in \{x, y, z\}$ .

We define the following elementary matrices :

$$\begin{aligned} \mathbb{M}_{ij}^{K,k} &= \int_K \Phi_i^{K,k} \Phi_j^{K,k} \, d\mathbf{x} \\ \mathbb{D}_{u,ij}^{K,k} &= \int_K \Phi_j^{K,k} \partial_u \Phi_i^{K,k} \, d\mathbf{x} & \mathbb{D}_{u,ij}^{K,k,k+1} &= \int_K \Phi_j^{K,k+1} \partial_u \Phi_i^{K,k} \, d\mathbf{x} \\ \mathbb{F}_{\ell,ij}^{K,k} &= \int_{\partial K^\ell} \Phi_i^{K,k} \Psi_j^{g(K,\ell),k} \, d\sigma & \mathbb{F}_{\ell,ij}^{K,k+1,k} &= \int_{\partial K^\ell} \Phi_i^{K,k+1} \Psi_j^{g(K,\ell),k} \, d\sigma \\ \mathbb{E}_{\ell,ij}^{K,k} &= \int_{\partial K^\ell} \Phi_i^{K,k} \Phi_j^{K,k} \, d\sigma \\ \mathbb{G}_{i,j}^m &= \int_{e_m} \Psi_i^{m,k} \Psi_j^{m,k} \, d\sigma \end{aligned} \tag{89}$$

where  $g(K, \ell)$  is the global number of the  $\ell$ -th face of element  $K$  and  $m$  is also the global number of the edge  $e_m$ .

Most of those matrices are not square matrices and their sizes are recalled in [TABLE 3](#).

Matrix	Rows	Columns
$\mathbb{M}^{K,k+1}$	$d_n(k+1)$	$d_n(k+1)$
$\mathbb{D}_u^{K,k+1}$	$d_n(k+1)$	$d_n(k+1)$
$\mathbb{D}_u^{K,k,k+1}$	$d_n(k)$	$d_n(k+1)$
$\mathbb{F}_\ell^{K,k+1,k}$	$d_n(k+1)$	$d_{n-1}(k)$
$\mathbb{T}_\ell^{K,k+1}$	$d_n(k+1)$	$d_n(k+1)$
$\mathbb{E}_\ell^{K,k+1}$	$d_n(k+1)$	$d_n(k+1)$

(a) Matrices needed for [\(46b\)](#)

Matrix	Rows	Columns
$\mathbb{M}^{K,k}$	$d_n(k)$	$d_n(k)$
$\mathbb{D}_u^{K,k,k+1}$	$d_n(k)$	$d_n(k+1)$
$\mathbb{F}_\ell^{K,k}$	$d_n(k)$	$d_{n-1}(k)$

(b) Matrices needed for [\(46a\)](#)

Matrix	Rows	Columns
$\mathbb{F}_\ell^{K,k}$	$d_n(k)$	$d_{n-1}(k)$
$\mathbb{F}_\ell^{K,k+1,k}$	$d_n(k+1)$	$d_{n-1}(k)$
$\mathbb{G}_\ell^{m,k}$	$d_{n-1}(k)$	$d_{n-1}(k)$

(c) Matrices needed for [\(46c\)](#)

Table 3: Summary of the dimensions of the local matrices

The matrix  $\mathbb{T}_\ell^{K,k+1}$  will be used to evaluate the projection  $P_M p$  and will be defined in [\(93\)](#). It was added to [TABLE 3](#) for completeness.

**Remark 5.4:** To compute the integrals involving polynomials of different degrees in a nodal framework, it is possible to directly use the expressions given in [\(89\)](#) or to use the following trick. Let  $(\mathbf{x}_j^{k+1})_{j=1}^{d_n(k+1)}$  be the degrees of freedom associated with  $\mathcal{P}_{k+1}(K)$ . We define the *projection matrix*  $\mathbb{P}$  by

$$\mathbb{P}_{ij}^{K,k,k+1} = \Phi_i^{K,k}(\mathbf{x}_j^{k+1}).$$

Seeing  $\Phi_i^{K,k}$  as a polynomial of degree  $k+1$ , we can express it in the basis  $(\Phi_j^{K,k+1})_{j=1}^{d_n(k+1)}$ . Due to the nodal nature of the basis functions, we have

$$\Phi_i^{K,k} = \sum_{r=1}^{d_n(k+1)} \Phi_i^{K,k}(\mathbf{x}_r^{k+1}) \Phi_r^{K,k+1} = \sum_{r=1}^{d_n(k+1)} \mathbb{P}_{ir}^{K,k,k+1} \Phi_r^{K,k+1}.$$

The matrix  $\mathbb{D}^{K,k,k+1}$  is therefore given by

$$\mathbb{D}^{K,k,k+1} = \mathbb{P}^{K,k,k+1} \mathbb{D}^{K,k+1}.$$

#### 5.4.2 Evaluating the projection:

The evaluation of the term

$$\langle \tau P_M p_h, w_h \rangle_{\partial K} \tag{90}$$

where  $w_h \in \mathcal{P}_{k+1}(K)$  is a difficult part of the implementation of the HDG+ method.

Notice however that the implementation of the term

$$\langle P_M p_h, \mu_h \rangle_{\partial K}, \quad \mu \in M_h(\partial K)$$

arising in the discretization of the transmission condition [\(46c\)](#) is straightforward. Indeed as  $\mu_h \in \mathcal{P}_k(e)$ , by definition of  $P_M$  we have

$$\langle P_M p_h, \mu_h \rangle_{\partial K} = \langle p_h, \mu_h \rangle_{\partial K}.$$

In this section we present three techniques to compute [\(90\)](#).

**Efficient implementation using a hierarchical basis:** The easiest way to compute (90) is to use a hierarchical and orthogonal basis for  $\mathcal{P}_{k+1}(K)$ . Let  $(\Phi_j^K)_{j=1}^{d_n(k+1)}$  be such a basis, then  $(\Phi_j^K)_{j=1}^{d_n(k)}$  is a basis for  $\mathcal{P}_k(K)$ . Therefore if

$$p^K = \sum_{j=1}^{d_n(k+1)} p_j^K \Phi_j^K$$

then

$$P_M p^K = \sum_{j=1}^{d_n(k)} p_j^K \Phi_j^K.$$

A good choice for such a basis would probably be Dubiner's one which is  $L^2$ -orthogonal if the reference element is

- *in 2D*: the triangle with vertices

$$\mathbf{x}_1 = (-1, -1) ; \quad \mathbf{x}_2 = (1, -1) ; \quad \mathbf{x}_3 = (-1, 1),$$

- *in 3D*: the tetrahedron with vertices

$$\mathbf{x}_1 = (-1, -1, -1) ; \quad \mathbf{x}_2 = (1, -1, -1) ; \quad \mathbf{x}_3 = (-1, 1, -1) ; \quad \mathbf{x}_4 = (-1, -1, 1).$$

See [Dub91, Kir04] for more details.

As our solver is developed in the framework of nodal discontinuous Galerkin methods, we will not use this method for our implementation.

**Efficient implementation through Gauss-Legendre quadrature:** As  $\tau$  is constant on each edge, we have

$$\langle \tau P_M p_h, u \rangle_{\partial K} = \langle \tau P_M p_h, P_M u \rangle_{\partial K}.$$

This integral can be efficiently computed in two dimensions by using a Gauss-Legendre quadrature. This trick is described in [Oik14, Sec. 3.4]. In two dimensions, the edges of a triangle is a one-dimensional interval  $I$ . For simplicity, we assume that  $I = [-1, 1]$ .

Denote by  $\mathcal{G}_k$  the  $k$  points Gauss-Legendre quadrature rule. For a function  $f$ , it is given by

$$\mathcal{G}_k[f] = \sum_{i=1}^k w_i f(a_i)$$

where  $(w_i)_{i \in \llbracket 1, k \rrbracket}$  are the quadrature weights and  $(a_i)_{i \in \llbracket 1, k \rrbracket}$  are the quadrature points. If the function  $f$  is regular enough, then  $\mathcal{G}_k$  approximates the integral of  $f$  over  $I$ . In particular  $\mathcal{G}_{k+1}$  is exact for polynomials of degree up to  $2k+1$ . The usual stabilization term with  $p, u \in \mathcal{P}_{k+1}(K)$ , can be exactly computed with the  $k+2$  quadrature rule  $\mathcal{G}_{k+2}$  (which is exact for polynomials of degree up to  $2k+3$ ), *ie.*

$$\mathcal{G}_{k+2}[\tau p u] = \langle \tau p, u \rangle_{\partial K}.$$

If the  $k+1$  points rule  $\mathcal{G}_{k+1}$  is used instead, the HDG+ stabilization term is obtained

$$\mathcal{G}_{k+1}[\tau p u] = \langle \tau P_M p, P_M u \rangle_{\partial K}.$$

### Lemma 5.1:

$$\forall u, v \in \mathcal{P}_{k+1}(I), \quad \mathcal{G}_{k+1}[uv] = \int_I P_M u P_M v dx$$

**Proof:** Let  $\varphi_m$  be the Legendre polynomial of order  $m \geq 0$ . On  $I$ , we can write

$$u = \sum_{i=0}^{k+1} u_i \varphi_i \quad \text{and} \quad v = \sum_{i=0}^{k+1} v_i \varphi_i.$$

Due to the hierarchical nature of the Legendre polynomials, we know that

$$P_M u = \sum_{i=0}^k u_i \varphi_i \quad \text{and} \quad P_M v = \sum_{i=0}^k v_i \varphi_i.$$

We recall that  $\mathcal{G}_{k+1}$  is exact for polynomials of degree up to  $2k+1$  and that  $\varphi_{k+1}$  vanishes at the quadrature points, *ie.*

$$\forall i \in \llbracket 1, k+1 \rrbracket, \quad \varphi_{k+1}(a_i) = 0.$$

Then we have

$$\begin{aligned} \mathcal{G}_{k+1}[uv] &= \sum_{i=1}^{k+1} \left[ w_i \sum_{\ell=0}^{k+1} u_\ell \varphi_\ell(a_i) \sum_{m=0}^{k+1} v_m \varphi_m(a_i) \right] \\ &= \sum_{i=1}^{k+1} \left[ w_i \sum_{\ell=0}^k u_\ell \varphi_\ell(a_i) \sum_{m=0}^k v_m \varphi_m(a_i) \right] \\ &= \sum_{i=1}^{k+1} w_i (P_M u)(a_i) (P_M v)(a_i) \\ &= \mathcal{G}_{k+1}[P_M u P_M v] \\ &= \int_I P_M u P_M v dx, \end{aligned}$$

as  $\deg(P_M u P_M v) = 2k$ .

**Remark 5.5:** the decomposition onto the Legendre basis is needed to prove the result, but another basis (*eg.* Lagrange's nodal basis) can be used in the implementation of the method.

**Remark 5.6:** unfortunately this result is not yet generalized for three-dimensional problems.

**Generic implementation:** We will now present a way to compute (90) that works for every choice of polynomial basis in every space dimension. As this technique requires to precompute the projections for all the basis functions of order  $k+1$ , the nodal HDG+ method may be less attractive when  $p$ -adaptivity is needed.

We assume that

$$p^K = \sum_{j=1}^{d_n(k+1)} p_j^K \Phi_j^{K,k+1} \quad \text{and} \quad w_h = \Phi_i^{K,k+1}$$

(90) therefore becomes

$$\langle P_M p_h^K, w_h \rangle_{\partial K^\ell} = \sum_{j=1}^{d_n(k+1)} p_j^K \langle P_M \Phi_j^{K,k+1}, \Phi_i^{K,k+1} \rangle_{\partial K^\ell}. \quad (91)$$

The next step is to decompose  $P_M \Phi_j^{K,k+1}$  onto the basis of  $M_h(\partial K)$

$$P_M \Phi_j^{K,k+1} = \sum_{r=1}^{d_{n-1}(k)} \theta_{j,r}^{K,\ell} \Psi_r^{e,k},$$

and we can construct a linear system for  $\underline{\theta_j^{K,\ell}}$  using the definition of  $P_M$

$$\sum_{r=1}^{d_{n-1}(k)} \langle \Psi_r^{e,k}, \Psi_i^{e,k} \rangle_{\partial K^\ell} \theta_{j,r}^{K,\ell} = \langle \Phi_j^{K,k+1}, \Psi_i^{e,k} \rangle_{\partial K^\ell}, \quad \forall i \in \llbracket 1, d_{n-1}(k) \rrbracket$$

or in matrix form

$$\mathbb{G}^{g(K,\ell)} \underline{\theta_j^{K,\ell}} = \underline{s_j^{K,\ell}}$$

where

$$s_{j,i}^{K,\ell} := \langle \Phi_j^{K,k+1}, \Psi_i^{e,k} \rangle_{\partial K^\ell}.$$

Plugging this result into (91), we have

$$\langle P_M p_h^K, \Phi_i^{K,k+1} \rangle_{\partial K^\ell} = \sum_{j=1}^{d_n(k+1)} \sum_{r=1}^{d_{n-1}(k)} p_j^K \theta_{j,r}^{K,\ell} \langle \Psi_r^{e,k}, \Phi_i^{K,k+1} \rangle_{\partial K^\ell}. \quad (92)$$

We define the following matrix

$$\mathbb{T}_{\ell,ij}^{K,k+1} := \sum_{r=1}^{d_{n-1}(k)} \mathbb{F}_{\ell,ir}^{K,k+1,k} \theta_{j,r}^{K,\ell}, \quad (93)$$

where we recall that

$$\mathbb{F}_{\ell,ij}^{K,k+1,k} := \int_{\partial K^\ell} \Phi_i^{K,k+1} \Psi_j^{g(K,\ell),k} d\sigma,$$

and we can rewrite (92) in matrix form

$$\langle P_M p_h^K, \Phi_i^{K,k+1} \rangle_{\partial K} = \sum_{\ell=1}^4 \sum_{j=1}^{d_n(k+1)} \mathbb{T}_{\ell,ij}^{K,k+1} p_j^K.$$

### 5.4.3 Local problem

Using the elementary matrices introduced in (89) we can write (46b) as

$$\begin{aligned} -\rho_0 \omega^2 \mathbb{M}^{K,k+1} \underline{p}^K + 2i\omega \sum_u b_0^u \mathbb{D}_u^{K,k+1} \underline{p}^K + \sum_u (\mathbb{D}_u^{K,k,k+1})^T \underline{q}_u^K + 2i\omega \sum_\ell \tau_\ell \left( \mathbb{T}_\ell^{K,k+1} \underline{p}^K - \mathbb{F}^{K,k+1,k} \widehat{\underline{p}}^{g(K,\ell)} \right) \\ - 2i\omega \sum_\ell \left[ \alpha_\ell \mathbb{F}^{K,k+1,k} \widehat{\underline{p}}^{g(K,\ell)} + t_{\text{upw},\ell} \mathbb{E}_\ell^{K,k+1} \underline{p}^K \right] = \mathbb{S}^K \end{aligned} \quad (94)$$

and (46a) as

$$\sum_u W_{0,vu} \mathbb{M}^{K,k} \underline{q}_u^K - \mathbb{D}_v^{K,k,k+1} \underline{p}^K + \sum_\ell n_v^{K,\ell} \mathbb{F}_\ell^{K,k} \widehat{\underline{p}}^{g(K,\ell)} = 0 \quad \forall v \in \{x, y, z\} \quad (95)$$

Using TABLE 3 it is easy (and important) to check that (94) has  $d_n(k+1)$  equations and (95) has  $d_n(k)$  equations.

**Matrix form:** We can now construct the matrices  $\mathbb{A}^K$  and  $\mathbb{C}^K$  for the local problem (82).  
 Matrix  $\mathbb{A}^K$ : (size:  $(d_n(k+1) + 3d_n(k))^2$ )

$$\mathbb{A}^K \parallel \begin{array}{c|ccc} & \underline{p^K} & \underline{q_x^K} & \underline{q_y^K} & \underline{q_z^K} \\ \hline \underline{p^K} & -\rho_0\omega^2\mathbb{M}^{K,k+1} + 2i\omega \sum b_{0,u}\mathbb{D}_u^{K,k+1} & & & \\ & -2i\omega \sum_{\ell} t_{\text{upw},\ell} \mathbb{E}_{\ell}^{K,k+1} & (\mathbb{D}_x^{K,k,k+1})^T & (\mathbb{D}_y^{K,k,k+1})^T & (\mathbb{D}_z^{K,k,k+1})^T \\ & + 2i\omega \sum_{\ell} \tau_{\ell} \mathbb{T}_{\ell}^{K,k+1} & & & \\ \hline \underline{q_x^K} & -\mathbb{D}_x^{K,k,k+1} & W_{0,11}^K \mathbb{M}^{K,k} & W_{0,12}^K \mathbb{M}^{K,k} & W_{0,13}^K \mathbb{M}^{K,k} \\ \underline{q_y^K} & -\mathbb{D}_y^{K,k,k+1} & W_{0,21}^K \mathbb{M}^{K,k} & W_{0,22}^K \mathbb{M}^{K,k} & W_{0,23}^K \mathbb{M}^{K,k} \\ \underline{q_z^K} & -\mathbb{D}_z^{K,k,k+1} & W_{0,31}^K \mathbb{M}^{K,k} & W_{0,32}^K \mathbb{M}^{K,k} & W_{0,33}^K \mathbb{M}^{K,k} \end{array}$$

Matrix  $\mathbb{C}^K$ : (size:  $(d_n(k+1) + 3d_n(k)) \times (4d_{n-1}(k))$ )

$$\mathbb{C}^K \parallel \begin{array}{c|cccc} & \underline{\hat{p}^{g(K,1)}} & \underline{\hat{p}^{g(K,2)}} & \underline{\hat{p}^{g(K,3)}} & \underline{\hat{p}^{g(K,4)}} \\ \hline \underline{p^K} & -2i\omega(\tau_1 + \alpha_1)\mathbb{F}_1^{K,k+1,k} & -2i\omega(\tau_2 + \alpha_2)\mathbb{F}_2^{K,k+1,k} & -2i\omega(\tau_3 + \alpha_3)\mathbb{F}_3^{K,k+1,k} & -2i\omega(\tau_4 + \alpha_4)\mathbb{F}_4^{K,k+1,k} \\ \underline{q_x^K} & n_x^{K,1}\mathbb{F}_1^{K,k} & n_x^{K,2}\mathbb{F}_2^{K,k} & n_x^{K,3}\mathbb{F}_3^{K,k} & n_x^{K,4}\mathbb{F}_4^{K,k} \\ \underline{q_y^K} & n_y^{K,1}\mathbb{F}_1^{K,k} & n_y^{K,2}\mathbb{F}_2^{K,k} & n_y^{K,3}\mathbb{F}_3^{K,k} & n_y^{K,4}\mathbb{F}_4^{K,k} \\ \underline{q_z^K} & n_z^{K,1}\mathbb{F}_1^{K,k} & n_z^{K,2}\mathbb{F}_2^{K,k} & n_z^{K,3}\mathbb{F}_3^{K,k} & n_z^{K,4}\mathbb{F}_4^{K,k} \end{array}$$

#### 5.4.4 Global problem

Using the elementary matrices introduced in (89) we can write (46c) as

$$\sum_{K,\ell} \left[ \sum_u n_u^{K,\ell} (\mathbb{F}_{\ell}^{K,k})^T \underline{q_u^K} + 2i\omega \sum_{\ell} \tau_{\ell} \left( (\mathbb{F}_{\ell}^{K,k+1,k})^T \underline{p^K} - \mathbb{G}^{g(K,\ell),k} \underline{\hat{p}^{g(K,\ell)}} \right) \right] - 2i\omega \sum_{K,\ell} \left[ \alpha_{\ell} \mathbb{G}^{g(K,\ell),k} \underline{\hat{p}^{g(K,\ell)}} + t_{\text{upw},\ell} (\mathbb{F}_{\ell}^{K,k+1,k})^T \underline{p^K} \right] = 0 \quad (96)$$

Once again, TABLE 3 should be used to check the dimensions of the matrices involved in (96).

**Matrix form:** We will now use (96) to construct the matrices  $\mathbb{B}^K$  and  $\mathbb{L}^K$  of (83).

Matrix  $\mathbb{B}^K$ : (size:  $4d_{n-1}(k) \times (d_n(k+1) + 3d_n(k))$ )

$$\mathbb{B}^K \parallel \begin{array}{c|ccc} & \underline{p^K} & \underline{q_x^K} & \underline{q_y^K} & \underline{q_z^K} \\ \hline e_1 & 2i\omega(\tau_1 - t_{\text{upw},1})(\mathbb{F}_1^{K,k+1,k})^T & n_x^{K,1}(\mathbb{F}_1^{K,k})^T & n_y^{K,1}(\mathbb{F}_1^{K,k})^T & n_z^{K,1}(\mathbb{F}_1^{K,k})^T \\ e_2 & 2i\omega(\tau_2 - t_{\text{upw},2})(\mathbb{F}_2^{K,k+1,k})^T & n_x^{K,2}(\mathbb{F}_2^{K,k})^T & n_y^{K,2}(\mathbb{F}_2^{K,k})^T & n_z^{K,2}(\mathbb{F}_2^{K,k})^T \\ e_3 & 2i\omega(\tau_3 - t_{\text{upw},3})(\mathbb{F}_3^{K,k+1,k})^T & n_x^{K,3}(\mathbb{F}_3^{K,k})^T & n_y^{K,3}(\mathbb{F}_3^{K,k})^T & n_z^{K,3}(\mathbb{F}_3^{K,k})^T \\ e_4 & 2i\omega(\tau_4 - t_{\text{upw},4})(\mathbb{F}_4^{K,k+1,k})^T & n_x^{K,4}(\mathbb{F}_4^{K,k})^T & n_y^{K,4}(\mathbb{F}_4^{K,k})^T & n_z^{K,4}(\mathbb{F}_4^{K,k})^T \end{array}$$

Matrix  $\mathbb{L}^K$ : (size:  $(4d_{n-1}(k))^2$ )

$$\mathbb{L}^K \parallel \begin{array}{c|cccc} & \underline{\hat{p}^{g(K,1)}} & \underline{\hat{p}^{g(K,2)}} & \underline{\hat{p}^{g(K,3)}} & \underline{\hat{p}^{g(K,4)}} \\ \hline e_1 & -2i\omega(\tau_1 + \alpha_1)\mathbb{G}^{g(K,1),k} & & & \\ e_2 & & -2i\omega(\tau_2 + \alpha_2)\mathbb{G}^{g(K,2),k} & & \\ e_3 & & & -2i\omega(\tau_3 + \alpha_3)\mathbb{G}^{g(K,3),k} & \\ e_4 & & & & -2i\omega(\tau_4 + \alpha_4)\mathbb{G}^{g(K,4),k} \end{array}$$

## 5.5 Comparison of the cost of the HDG and HDG+ methods

Now that we have written down the discrete systems for the HDG and HDG+ methods, it is interesting to compare the sizes of systems that we need to solve.

In TABLE 4 we have written down the sizes of the matrices for both methods. We can see that the global problem  $\mathbb{K}$  of the HDG+ has the same dimension as the one of the HDG method of degree  $k$ . As the resolution of the global problem is the most expensive step in the resolution of the method, the cost of the HDG+ method is therefore similar to the cost of the HDG method of degree  $k$ , while yielding to an order of convergence of  $k + 2$  instead of  $k + 1$ . We can also notice that the cost of the local problems of the HDG+ method is intermediate between the cost of the local problems of the HDG methods of degree  $k$  and  $k + 1$ . However this has a really limited impact on the computational cost.

We denote by  $N_{\text{elt}}$  the number of elements in the mesh, *ie*  $N_{\text{elt}} = \text{card}(\mathcal{T}_h)$ .

Matrix	HDG $k$	HDG+	HDG $k + 1$
$\mathbb{A}^K$	$(4d_n(k))^2$	$(d_n(k+1) + 3d_n(k))^2$	$(4d_n(k+1))^2$
$\mathbb{C}^K$	$4d_n(k) \times 4d_{n-1}(k)$	$(d_n(k+1) + 3d_n(k)) \times (4d_{n-1}(k))$	$4d_n(k+1) \times 4d_{n-1}(k+1)$
$\mathbb{B}^K$	$4d_{n-1}(k) \times 4d_n(k)$	$4d_{n-1}(k) \times (d_n(k+1) + 3d_n(k))$	$4d_{n-1}(k+1) \times 4d_n(k+1)$
$\mathbb{L}^K$	$(4d_{n-1}(k))^2$	$(4d_{n-1}(k))^2$	$(4d_{n-1}(k+1))^2$
$\mathbb{K}$	$\sim (N_{\text{elt}}d_{n-1}(k))^2$	$\sim (N_{\text{elt}}d_{n-1}(k))^2$	$\sim (N_{\text{elt}}d_{n-1}(k+1))^2$

Table 4: Size of the matrices for the HDG and HDG+ methods in 3D

In FIGURE 5, we have plotted the number of degrees of freedom for local problems ( $N_{\text{dof}}^{\text{loc}}$ ) and for the global problem ( $N_{\text{dof}}^{\text{glob}}$ ) with  $N_{\text{elt}} = 10^3$  for several polynomial degrees in 3D. Notice that  $\mathbb{A}^K$  has dimension  $(N_{\text{dof}}^{\text{loc}})^2$  and  $\mathbb{K}$  has dimension  $(N_{\text{dof}}^{\text{glob}})^2$ .

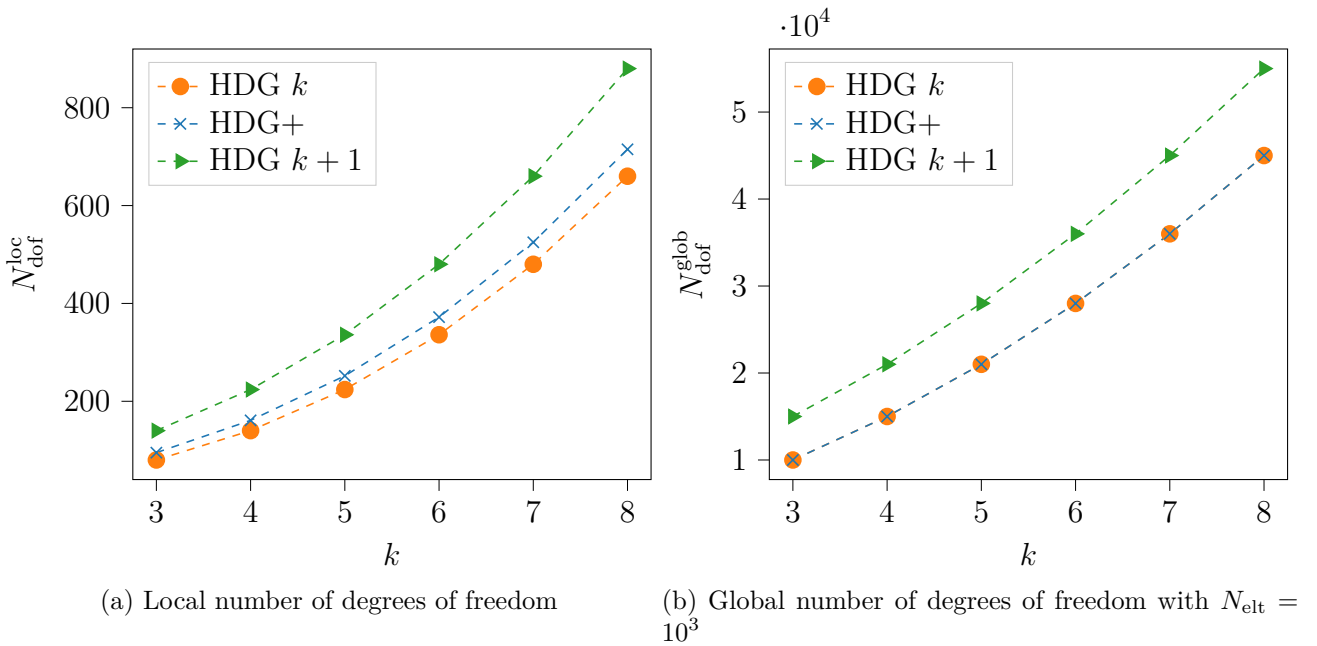


Figure 5: Local and global number of degrees of freedom for the three HDG methods in 3D

To put this different costs in perspective, we can also compare them to the cost of a standard DG method. Using TABLE 4, we know that the number of degrees of freedom of the HDG

method of degree  $k$  is

$$N_{\text{dof}}^{\text{HDG-}k} \sim N_{\text{elt}} d_{n-1}(k) = N_{\text{elt}} \binom{k+n-1}{n-1} \sim k^{n-1} N_{\text{elt}},$$

whereas the number of degrees of freedom of a DG method of degree  $k$  is

$$N_{\text{dof}}^{\text{DG-}k} \sim N_{\text{elt}} d_n(k) = N_{\text{elt}} \binom{k+n}{n} \sim k^n N_{\text{elt}},$$

which becomes much larger as  $k$  increases. We therefore retrieved the fact that HDG in dimension  $n$  has the same cost as DG in dimension  $n - 1$ , this is depicted in [FIGURE 6](#).

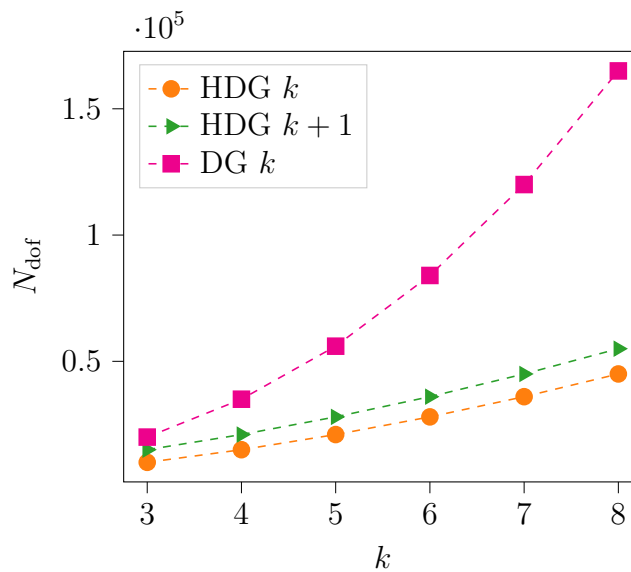


Figure 6:  $N_{\text{dof}}^{\text{HDG-}k}$ ,  $N_{\text{dof}}^{\text{HDG-}k+1}$  and  $N_{\text{dof}}^{\text{DG-}k}$  in 3D with  $N_{\text{elt}} = 10^3$

## 6 Numerical experiments

In this section we will provide some numerical experiments for the three HDG methods that were described and analysed in this paper. We will first focus on the simple case of duct modes propagating in a waveguide to obtain convergence curves and validate the theoretical results of the previous sections. We will then provide some illustrative examples to show that those methods can be used in more realistic cases.

### 6.1 Convergence rate

In this subsection, we will present some numerical experiments to illustrate our theoretical results. As most of the estimates obtained in our analysis involve projection errors of the form

$$\|p_h - \pi_W p\|_{\mathcal{T}_h} \quad \text{or} \quad \|\mathbf{q}_h - \boldsymbol{\pi}_V \mathbf{q}\|_{\mathcal{T}_h},$$

we will need to evaluate those projections before actually computing errors. In [TABLE 5](#) we recall the different projections used for the analysis of the three different variants of the HDG method that we considered.



Method	$\pi_W$	$\pi_V$	$P_M$
HDG $\mathbf{q}$	$L^2$	$L^2$	$L^2$
HDG +	$L^2$	$L^2$	$L^2$
HDG $\boldsymbol{\sigma}$	HDG	HDG	$L^2$

Table 5: Summary of the different projections used for the analysis of the HDG methods

**Geometric settings:** As depicted on [FIGURE 7](#) we consider a uniform directional flow  $\mathbf{v}_0 = Mc_0\mathbf{e}_x$ , where  $M$  is the *Mach number*.

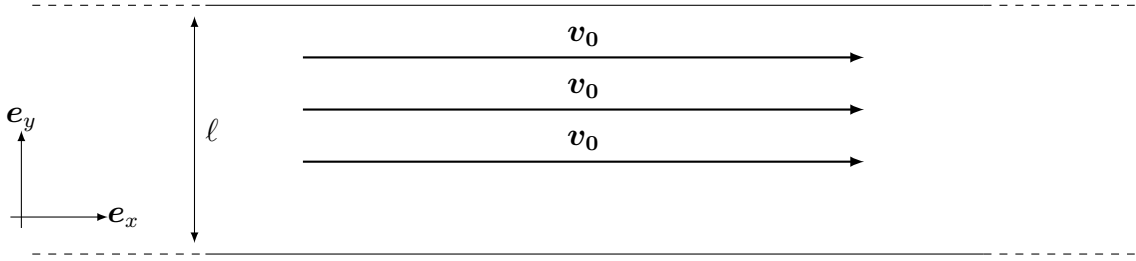


Figure 7: Sketch of the geometric configuration

Unless stated otherwise, we will always use the following parameters for the convergence tests

$$\mathcal{O} = (0, 2) \times (0, 1) \quad ; \quad \rho_0, c_0 \equiv 1 \quad ; \quad \omega = 5.55\pi,$$

and the choice of  $M$  will be specified for each numerical experiment.

**Analytic solution:** The duct modes are a family of analytic solutions of (1) in a waveguide, see [\[BDL03\]](#). They are given by

$$p_n^\pm(x, y) = e^{i\beta_n^\pm x} \varphi_n(y)$$

where

$$\begin{aligned} n < N_0 : & \quad \beta_n^\pm = \frac{-\kappa M \pm \sqrt{\kappa^2 - \frac{n^2\pi^2}{\ell^2}(1 - M^2)}}{1 - M^2} \\ n > N_0 : & \quad \beta_n^\pm = \frac{-\kappa M \pm i\sqrt{\frac{n^2\pi^2}{\ell^2}(1 - M^2) - \kappa^2}}{1 - M^2} \end{aligned}$$

with

$$\kappa = \frac{\omega}{c_0} \quad \text{and} \quad M = \frac{v_0}{c_0}$$

$$N_0 = \left\lfloor \frac{\kappa\ell}{\pi\sqrt{1 - M^2}} \right\rfloor$$

and

$$\begin{aligned} \varphi_0(y) &:= \sqrt{\ell^{-1}} \\ \varphi_n(y) &:= \sqrt{2\ell^{-1}} \cos\left(\frac{n\pi y}{\ell}\right), \quad n \in \mathbb{N}^* \end{aligned}$$

The choice of  $n$  will be specified for each numerical experiment.

**Evaluating the projections:** Here we give the details for evaluating the  $L^2$  projection onto  $W_h$ , the process is very similar for the other projections.

We recall that the dimension of the polynomial spaces is given by

$$\dim \mathcal{P}_k(K) = \binom{n+k}{n} =: d_n(k), \quad \text{for } K \subset \mathbb{R}^n.$$

All the projections considered are local to an element  $K$ , so evaluating them amounts to solving a linear system on each  $K$ . Indeed, the definition of the  $L^2$  projection onto  $W_h(K) = \mathcal{P}_k(K)$  gives

$$\forall i \in \llbracket 1, d_n(k) \rrbracket, \quad (\pi_W p, \Phi_i^{K,k})_K = (p, \Phi_i^{K,k})_K.$$

As  $\pi_W p \in W_h(K)$ , we can write

$$\pi_W p = \sum_{j=1}^{d_n(k)} \pi_j \Phi_j^{K,k},$$

where  $\underline{\pi} = (\pi_j)_j$  is the vector of the coordinates of  $\pi_W p$  in the basis  $(\Phi_j^K)_j$  of  $\mathcal{P}_k(K)$ . We therefore obtain the following system

$$\forall i \in \llbracket 1, d_n(k) \rrbracket, \quad \sum_{j=1}^{d_n(k)} \pi_j (\Phi_j^{K,k}, \Phi_i^{K,k})_K = (p, \Phi_i^{K,k})_K,$$

or in matrix form

$$\mathbb{M}^{K,k} \underline{\pi} = \underline{s}, \quad \text{where } s_i = (p, \Phi_i^{K,k})_K.$$

The integral in the right-hand side is evaluated using a 91 points Gauss-Lobatto quadrature formula and the linear system can be solved using `lapack`.

**Evaluating the  $L^2$ -error:** Now that the projections can be evaluated, it remains to compute the  $L^2$  norms. As the numerical solution and the projections are polynomial quantities, this can be done using the mass matrix. Indeed, for  $u \in \mathcal{P}_k(K)$  we have

$$\begin{aligned} \|u\|_K^2 &= \left( \sum_{j=1}^{d_n(k)} u_j \Phi_j^{K,k}, \sum_{i=1}^{d_n(k)} u_i \Phi_i^{K,k} \right)_K \\ &= \sum_{i,j=1}^{d_n(k)} u_j u_i^* (\Phi_j^{K,k}, \Phi_i^{K,k})_K \\ &= \sum_{i=1}^{d_n(k)} u_i^* \sum_{j=1}^{d_n(k)} \mathbb{M}_{ij}^{K,k} u_j \\ &= \underline{u}^* \mathbb{M}^{K,k} \underline{u}. \end{aligned}$$

To obtain more meaningful results, we will use relative errors instead of the standard  $L^2$ -error. This choice allows us to compare the errors computed on different meshes without any pollution coming from the different number of elements. The relative error for  $p_h$  is given by

$$\mathcal{E}_p := \frac{\|p_h - \pi_W p\|_{\mathcal{T}_h}}{\|\pi_W p\|_{\mathcal{T}_h}},$$

and similar expressions will be used for the other volumetric quantities.

**Some notation:** To allow the comparison between the HDG+ method and the two HDG methods, we would like to emphasize that  $k$  always denotes the polynomial degree used to approximate the *trace unknown*  $\widehat{p}_h$ . We have chosen to plot the relative errors against the quantity  $\frac{k}{h}$  which is proportional to the number of degrees of freedom per wavelength. We would like to point out that all the plots in the next sections will use a log-log scale.

**Hardware configuration:** Those numerical experiments were carried out on a *miriel* node on the *plafrim* cluster<sup>8</sup>. This node is equipped with a 2 dodeca-core Haswell Intel Xeon E5-2680 v3 with a clock rate of 2.5 GHz and 128 Go of memory.

### 6.1.1 Acoustic case without flow

A first important test to validate our implementation in the *acoustic case without flow*. In this case, the convected Helmholtz equation reduces to the standard Helmholtz equation and the HDG methods, either with the diffusive or total flux, are the same when there is no convection. We should therefore be able to reproduce the *super-convergence* of the HDG method for the Helmholtz equation, *ie.*

$$\mathcal{E}_p = \mathcal{O}(h^{k+2}),$$

when an approximation of degree  $k$  is used for all the unknowns.

Here we have chosen to use the following parameters

$$n = 0, \quad \text{and} \quad M = 0,$$

which correspond to a plane-wave.

The resulting convergence curve is depicted on [FIGURE 8](#) and we can clearly see that the super-convergence is obtained.

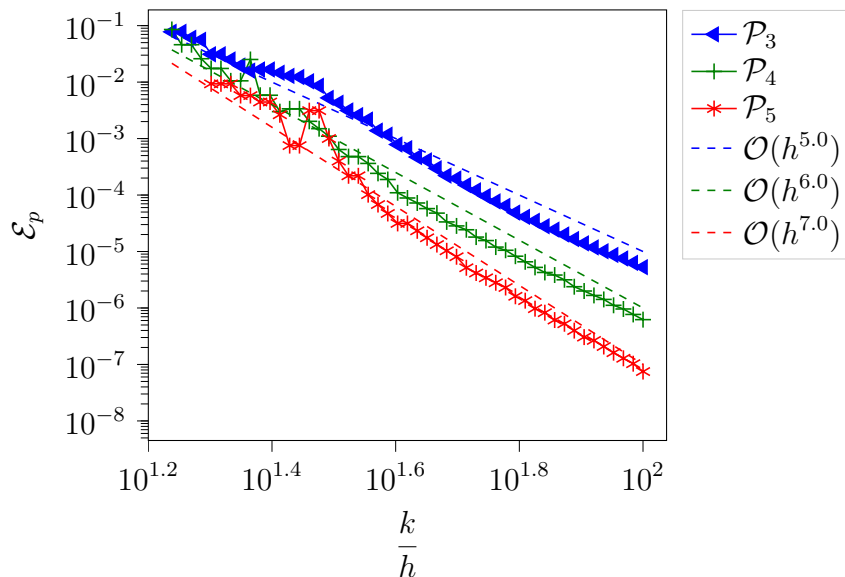


Figure 8: Convergence history for the HDG method without flow for the volumetric unknown  $p_h$

<sup>8</sup>See <http://www.plafrim.fr>.

### 6.1.2 Low Mach

We then move to a flow with a low Mach number. In this case we have used the following parameters

$$n = 3, \quad \text{and} \quad M = 0.2.$$

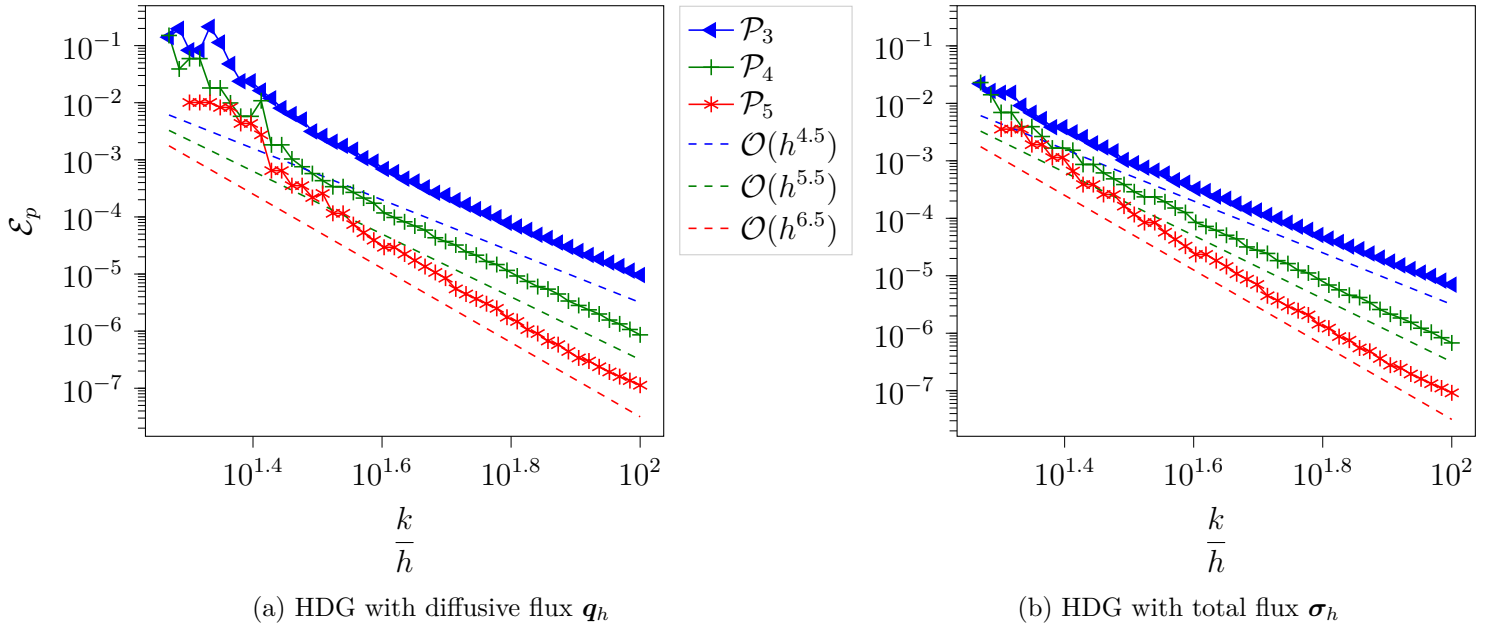


Figure 9: Low Mach convergence history for the volumetric unknown  $p_h$  for both the diffusive and total flux HDG method

The convergence history for the volumetric unknown  $p_h$  is displayed on [FIGURE 9](#) for both of the HDG methods. We can see that the *diffusive flux* formulation achieves an order of convergence of  $k + 3/2$  as expected. On the other hand the *total flux* formulation also achieves an order of convergence of  $k + 3/2$  which is better than the expected order  $k + 1$ . However for uniform flows and upon well choosing the penalization parameters both of those methods are algebraically equivalent and we therefore expect to obtain the same order of convergence for  $p_h$ . On [FIGURE 10](#) the convergence rate for the volumetric unknown  $p_h$  for the HDG+ method is displayed. As expected the optimal convergence rate of  $k + 2$  is obtained. As discussed in [TABLE 6](#), it is clear that the use of the HDG+ method is less expensive than the use of the HDG methods.

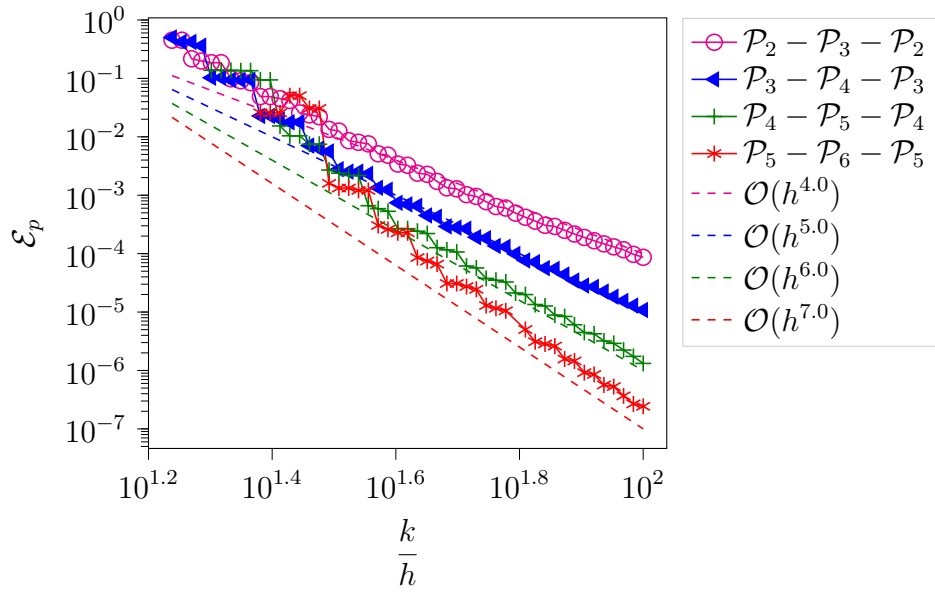


Figure 10: Low Mach convergence history for the volumetric unknown  $p_h$  for the HDG+ method

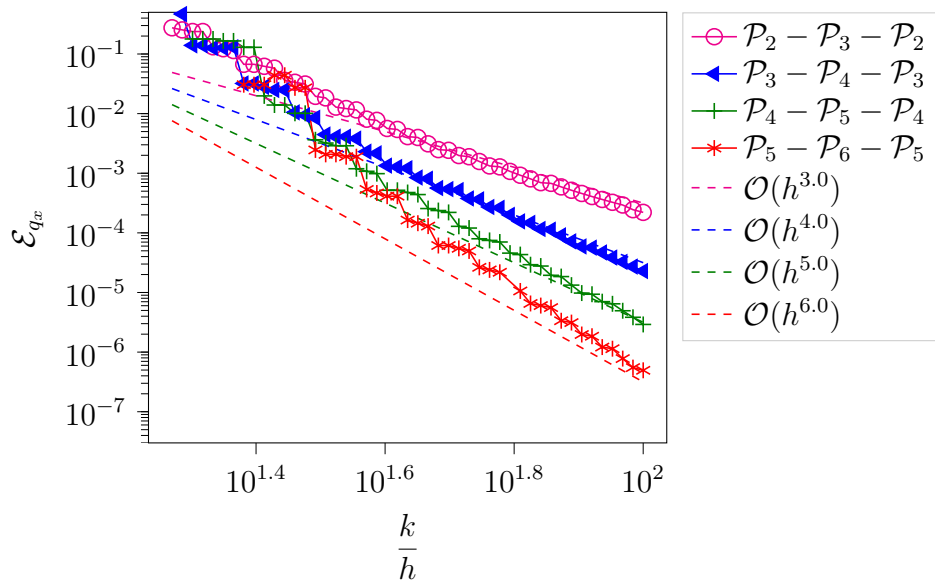


Figure 11: Low Mach convergence history for the first component of volumetric unknown  $\mathbf{q}_h$  for the HDG+ method

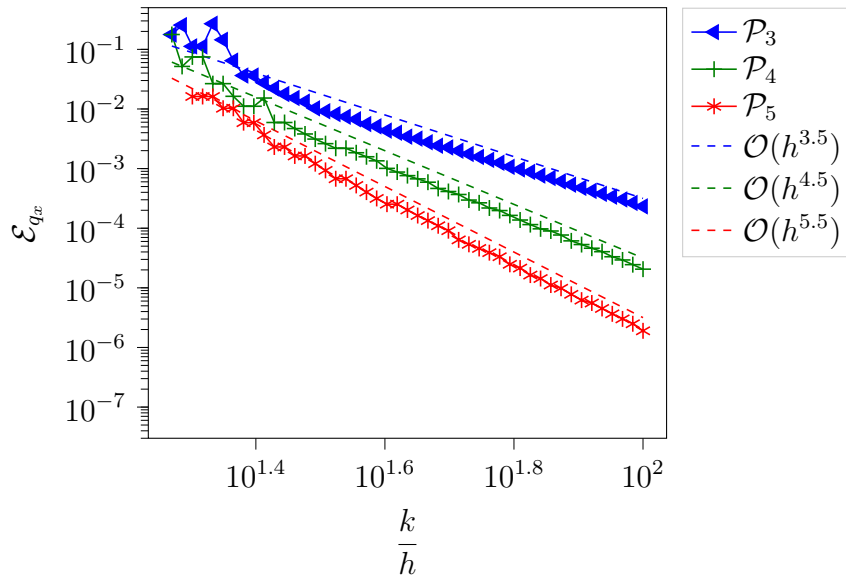
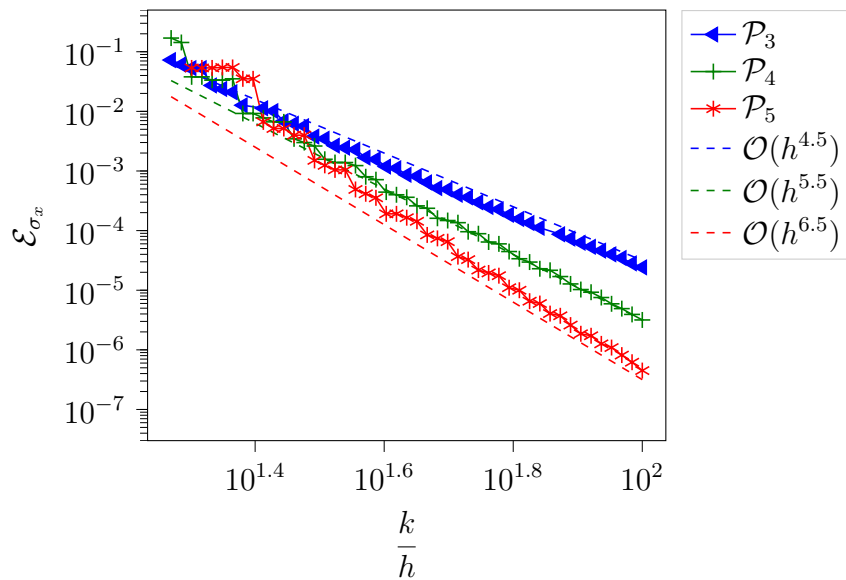

 (a) HDG with diffusive flux  $\mathbf{q}_h$ 

 (b) HDG with total flux  $\sigma_h$ 

Figure 12: Low Mach convergence history for the first component of the volumetric flux unknown

If we now move to the flux unknown, we can see on [FIGURE 11](#) that  $\mathbf{q}_h$  converges with the optimal order  $k + 1$  as expected. For the HDG methods, the convergence histories for  $\mathbf{q}_h$  and  $\sigma_h$  are depicted on [FIGURE 12](#). The *diffusive flux* formulation achieves a convergence order of  $k + 1/2$  as expected, whereas the *total flux* formulation outperforms the theoretical results and achieves an order of  $k + 3/2$ . For this formulation it was however expected that both  $p_h$  and  $\sigma_h$  converge with the same order, which is actually the case. Finally we would like to point out that the choice between the formulations with  $\mathbf{q}_h$  or  $\sigma_h$  is application dependent.

Finally information regarding the size of the global problems for a fixed number of degrees of freedom per wavelength are given in [TABLE 6](#). We have included the HDG+ method with  $(k, k-1)$  and the HDG+ method with  $(k+1, k)$ . When using the HDG+ method with  $(k, k-1)$ , a smaller global problem is solved to obtain the same convergence rate as the HDG methods. When using the HDG+ method with  $(k+1, k)$ , larger elements can be used to obtain the same number of degrees of freedom per wavelength because of the higher polynomial interpolation degree. We can clearly see that the HDG+ methods are computationally less expensive than

both of the HDG methods.

$k$	$h/k$		HDG- $\mathbf{q}_h$	HDG- $\boldsymbol{\sigma}_h$	HDG+ ( $k, k-1$ )	HDG+ ( $k+1, k$ )
3	$10^{-1}$	nnz	20 384	20 384	10 968	8 720
		nnz LÜ	36 872	36 872	15 216	9 800
		MUMPS time	$2.5 \cdot 10^{-2}$	$1.6 \cdot 10^{-2}$	$1.8 \cdot 10^{-2}$	$1.7 \cdot 10^{-2}$
3	$10^{-2}$	nnz	872 216	872 216	481 528	467 280
		nnz LÜ	3 572 752	3 572 752	1 835 148	1 673 560
		MUMPS time	0.27	0.23	0.15	0.10
4	$10^{-1}$	nnz	17 130	17 130	8 720	4 980
		nnz LÜ	24 125	24 125	9 800	4 644
		MUMPS time	$1.9 \cdot 10^{-2}$	$2.0 \cdot 10^{-2}$	$1.7 \cdot 10^{-2}$	$9.5 \cdot 10^{-3}$
4	$10^{-2}$	nnz	752 330	752 330	467 280	423 600
		nnz LÜ	2 874 075	2 874 075	1 673 560	1 397 418
		MUMPS time	0.19	0.17	0.10	$8.6 \cdot 10^{-2}$
5	$10^{-1}$	nnz	12 552	12 552	4 980	4 332
		nnz LÜ	14 112	14 112	4 644	4 068
		MUMPS time	$1.7 \cdot 10^{-2}$	$2.7 \cdot 10^{-2}$	$9.5 \cdot 10^{-3}$	$9.9 \cdot 10^{-3}$
5	$10^{-2}$	nnz	672 840	672 840	432 600	414 456
		nnz LÜ	2 409 006	2 409 006	1 397 418	1 333 368
		MUMPS time	0.12	0.11	$8.6 \cdot 10^{-2}$	$6.8 \cdot 10^{-2}$

Table 6: Size of the global systems and `mumps` elapsed time for different interpolation degree  $k$  for the low Mach case with a fixed number of degrees of freedom per wavelength

Notice that the numerical experiments performed here lead to relatively small linear system. Indeed our goal was to validate the numerical method rather than showing the ability of `hawen` to handle large numerical simulations.

In [TABLE 7](#), we review the size of the global system for various interpolation degrees  $k$  with a fixed error threshold for  $p_h$ . As we did before, we use the relative  $L^2$ -error defined by

$$\mathcal{E}_p := \frac{\|p_h - \pi_{WP}\|_{\mathcal{T}_h}}{\|\pi_{WP}\|_{\mathcal{T}_h}}.$$

We can clearly see that, when the desired error is smaller than  $10^{-2}$ , increasing the interpolation degree and thus the order of the method leads to solving a smaller linear system to obtain the same error level. This is less visible for an error threshold of  $10^{-2}$ , as this accuracy can be obtained with relatively large elements even with a low interpolation degree. For this error level, we can also see that there is no real difference in the size of the linear systems of HDG and HDG+ methods. This is expected as there is no real difference between the different interpolation degrees for this error threshold. However for smaller error thresholds, we can see that the HDG- $\mathbf{q}_h$ , HDG- $\boldsymbol{\sigma}_h$  and HDG+ with  $(k, k-1)$  lead to linear systems with similar sizes, which is expected as they share the same order. The HDG+ method with  $(k+1, k)$  also has a higher convergence rate and therefore lead to significantly smaller linear system for the same error level.

$k$	$\mathcal{E}_p$		HDG- $\mathbf{q}_h$	HDG- $\boldsymbol{\sigma}_h$	HDG+ ( $k, k-1$ )	HDG+ ( $k+1, k$ )
3	$10^{-2}$	nnz	49 784	39 216	37 200	26 920
		nnz LU	116 208	85 416	97 371	53 088
		MUMPS time	$3.1 \cdot 10^{-2}$	$1.8 \cdot 10^{-2}$	$2.5 \cdot 10^{-2}$	$1.3 \cdot 10^{-2}$
3	$10^{-3}$	nnz	99 920	102 800	92 040	49 784
		nnz LU	293 328	293 328	283 155	116 272
		MUMPS time	$3.0 \cdot 10^{-2}$	$2.8 \cdot 10^{-2}$	$3.2 \cdot 10^{-2}$	$1.9 \cdot 10^{-2}$
3	$10^{-4}$	nnz	210 360	267 784	255 552	140 640
		nnz LU	689 568	926 352	944 826	440 584
		MUMPS time	$4.9 \cdot 10^{-2}$	$6.2 \cdot 10^{-2}$	$8.1 \cdot 10^{-2}$	$3.63 \cdot 10^{-2}$
4	$10^{-2}$	nnz	28 240	28 240	26 920	28 240
		nnz LU	51 050	51 150	53 088	51 050
		MUMPS time	$2.3 \cdot 10^{-2}$	$1.4 \cdot 10^{-2}$	$1.3 \cdot 10^{-2}$	$1.4 \cdot 10^{-2}$
4	$10^{-3}$	nnz	58 560	58 560	49 784	42 050
		nnz LU	126 350	126 450	116 272	84 075
		MUMPS time	$1.9 \cdot 10^{-2}$	$2.0 \cdot 10^{-2}$	$1.9 \cdot 10^{-2}$	$2.2 \cdot 10^{-2}$
4	$10^{-4}$	nnz	106 880	151 600	140 640	77 770
		nnz LU	284 050	433 450	440 584	182 875
		MUMPS time	$2.5 \cdot 10^{-2}$	$3.2 \cdot 10^{-2}$	$3.63 \cdot 10^{-2}$	$2.1 \cdot 10^{-2}$
4	$10^{-5}$	nnz	255 580	255 580	372 640	156 100
		nnz LU	787 925	787 925	1 383 296	446 855
		MUMPS time	$3.9 \cdot 10^{-2}$	$4.1 \cdot 10^{-2}$	0.10	$3.4 \cdot 10^{-2}$
5	$10^{-2}$	nnz	24 660	24 660	28 240	24 660
		nnz LU	34 740	34 740	51 050	34 740
		MUMPS time	$1.1 \cdot 10^{-2}$	$1.4 \cdot 10^{-2}$	$1.4 \cdot 10^{-2}$	$1.9 \cdot 10^{-2}$
5	$10^{-3}$	nnz	40 656	40 656	42 050	40 656
		nnz LU	73 440	73 296	84 075	73 512
		MUMPS time	$1.6 \cdot 10^{-2}$	$1.4 \cdot 10^{-2}$	$2.2 \cdot 10^{-2}$	$1.3 \cdot 10^{-2}$
5	$10^{-4}$	nnz	84 312	88 200	77 770	60 540
		nnz LU	181 944	192 996	182 875	120 636
		MUMPS time	$2.6 \cdot 10^{-2}$	$1.9 \cdot 10^{-2}$	$2.1 \cdot 10^{-2}$	$1.5 \cdot 10^{-2}$
5	$10^{-5}$	nnz	148 704	153 888	156 100	111 972
		nnz LU	392 580	409 752	446 855	263 340
		MUMPS time	$3.1 \cdot 10^{-2}$	$2.9 \cdot 10^{-2}$	$3.4 \cdot 10^{-2}$	$2.5 \cdot 10^{-2}$

Table 7: Size of the global systems and `mumps` elapsed time for different interpolation degree  $k$  for the low Mach case with a fixed error threshold  $\mathcal{E}_p$

### 6.1.3 Large Mach

Finally we also considered a flow with a large March number. In this case, we used the following parameters

$$n = 3, \quad \text{and} \quad M = 0.8.$$

As the simulations of acoustic wave propagation in flows with large Mach numbers is known to be more challenging, we expect to see worse performances than in the previous subsection.



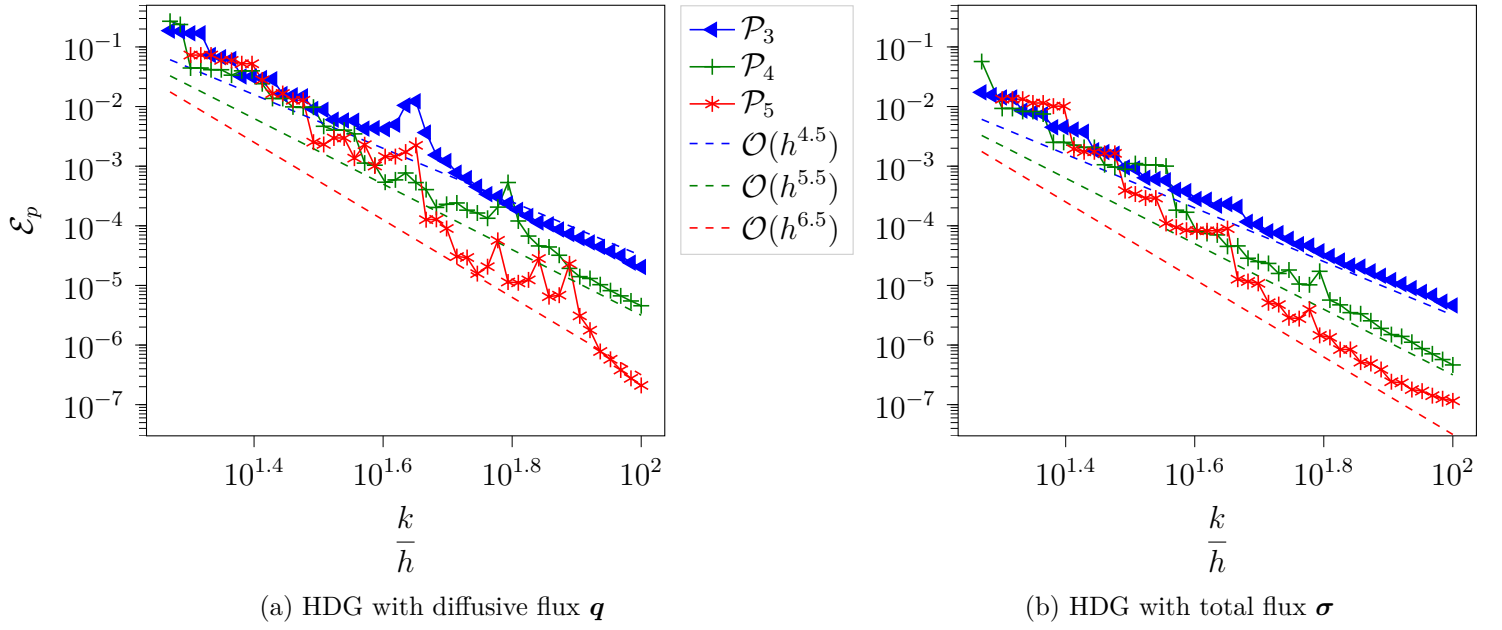


Figure 13: Large Mach convergence history for the volumetric unknown  $p_h$  for both of the HDG methods

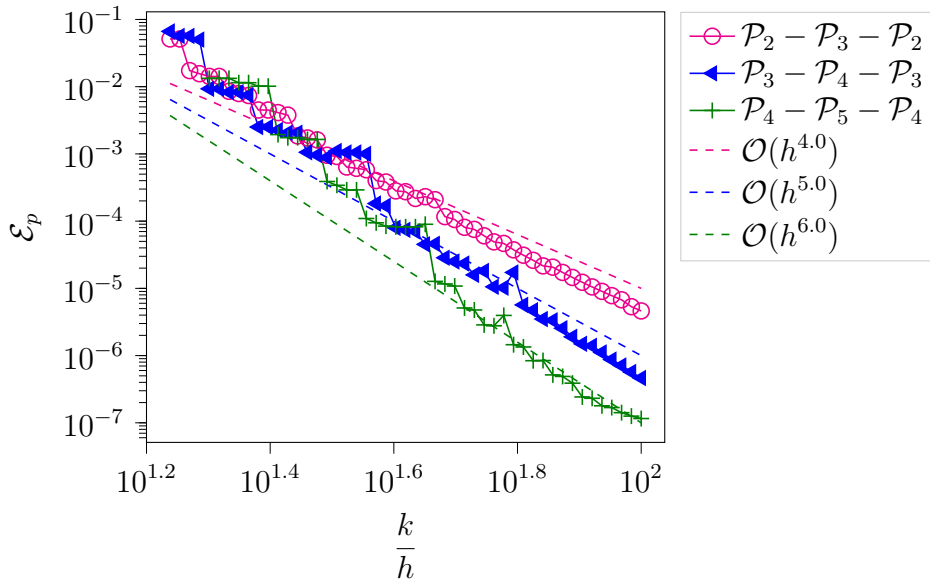


Figure 14: Large Mach convergence history for the volumetric unknown  $p_h$  for the HDG+ method

The convergence history for the volumetric unknown is depicted in [FIGURE 13](#) for both of the HDG methods. We can see that the *total flux* formulation still achieves the convergence order of  $k + 3/2$  whereas the behaviour of the *diffusive flux* formulation seems less robust. The same convergence history for HDG+ method is displayed on [FIGURE 14](#) and we can see that it still achieves the optimal convergence rate of  $k + 2$ .

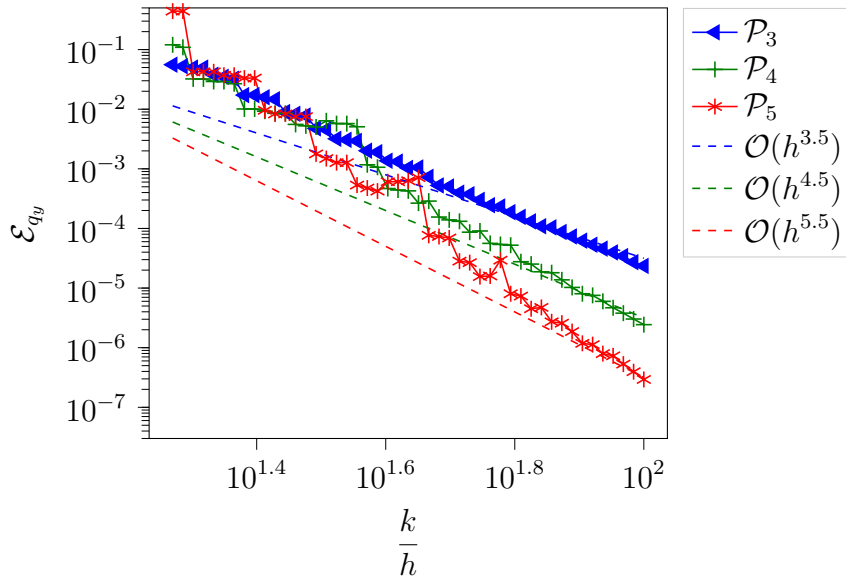
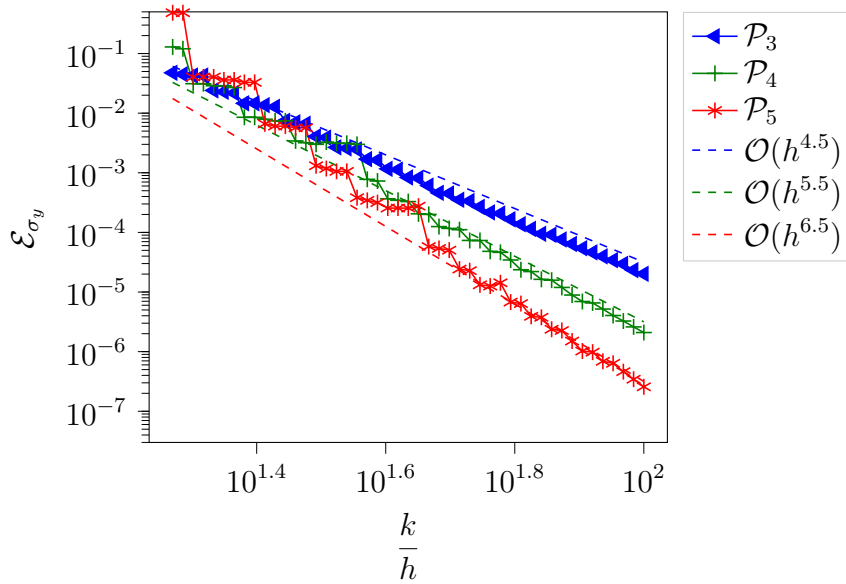

 (a) HDG with diffusive flux  $\mathbf{q}_h$ 

 (b) HDG with total flux  $\sigma_h$ 

Figure 15: Large Mach convergence history for the second component of the volumetric flux unknown

The convergence history for the volumetric flux unknown  $\mathbf{q}_h$  or  $\sigma_h$  is depicted in [FIGURE 15](#) for the HDG methods and in [FIGURE 16](#) for the HDG+ method. As in the low-Mach case, the HDG methods have a convergence rate of  $k + 3/2$  and the HDG+ method has a convergence rate of  $k + 1$ . Notice that the HDG- $\sigma_h$  method seems to be the most robust method for the approximation of the flux unknown for high Mach numbers.

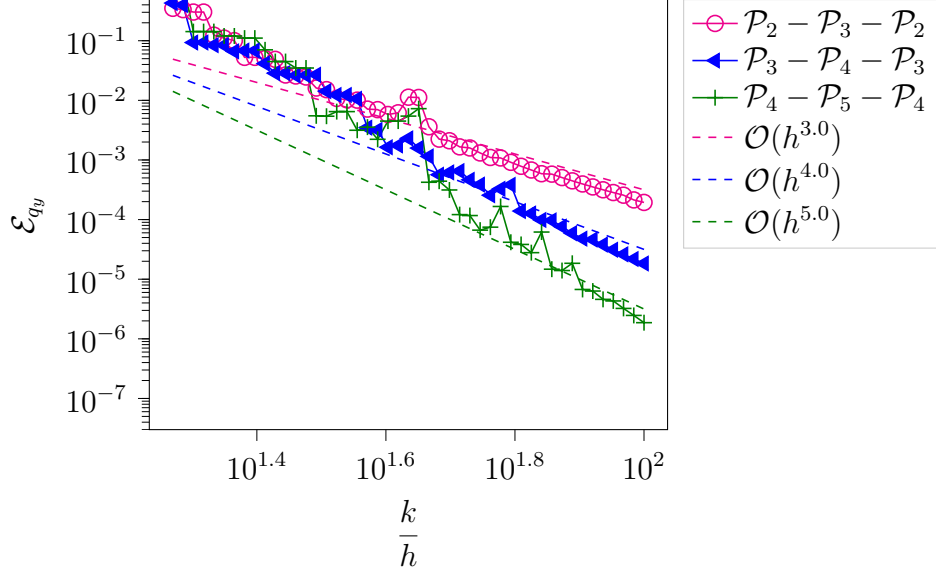


Figure 16: Large Mach convergence history for the second component of the volumetric unknown  $\mathbf{q}_h$  for the HDG+ method

## 6.2 A posteriori error estimate

In this section, we will show that it is possible to compute a simple *a posteriori* error indicator. For more complete approach to *a posteriori* error analysis for HDG methods, we refer to [CZ12, CZ13]. We introduce the *relative jump error*

$$\mathcal{E}_{\text{jump}} = \frac{\sqrt{\sum_{K,\ell} \|\widehat{p}_h - p_h\|_{\partial K^\ell}^2}}{\sqrt{\sum_{K,\ell} \|\widehat{p}_h\|_{\partial K^\ell}^2}},$$

which measures the jump between  $p_h$  and  $\widehat{p}_h$ . For the HDG+ method,  $p_h$  should be replaced with  $P_M p_h$ .

**Jump error and residuals for the HDG methods:** To understand the importance of this quantity, we introduce the residuals

$$\delta_h := -\omega^2 \rho_0 p_h + \text{div}(\boldsymbol{\sigma}_h) - s, \quad \text{and} \quad \boldsymbol{\Delta}_h := \mathbf{W}_0 \boldsymbol{\sigma}_h + \nabla p_h + 2i\omega p_h \mathbf{W}_0 \mathbf{b}_0.$$

We have chosen to work with the HDG- $\boldsymbol{\sigma}_h$  formulation as we recommend it over the HDG- $\mathbf{q}_h$  one, but the adaptation to the HDG- $\mathbf{q}_h$  formulation is immediate.

Reverting the integrations by parts in (12a)–(12b), we have

$$\begin{aligned} (\boldsymbol{\Delta}_h, \mathbf{r}_h)_K &= \langle p_h - \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial K}, \\ (\delta_h, w_h)_K &= -i\omega \langle \tau(p_h - \widehat{p}_h), w_h \rangle_{\partial K}, \end{aligned}$$

as  $\delta_h$  and  $\boldsymbol{\Delta}_h$  are usually not polynomial quantities this means that

$$\begin{aligned} (\mathbf{P}_V \boldsymbol{\Delta}_h, \mathbf{r}_h)_K &= \langle p_h - \widehat{p}_h, \mathbf{r}_h \cdot \mathbf{n} \rangle_{\partial K}, \\ (P_W \delta_h, w_h)_K &= -i\omega \langle \tau(p_h - \widehat{p}_h), w_h \rangle_{\partial K}, \end{aligned}$$

where  $\mathbf{P}_V$  and  $P_W$  are the  $L^2$ -orthogonal projections onto  $\mathbf{V}_h(K)$  and  $W_h(K)$  respectively. Notice that if  $\delta_h$  and  $\boldsymbol{\Delta}_h$  are actually polynomials things are even simpler as  $P_W \delta_h = \delta_h$  and  $\mathbf{P}_V \boldsymbol{\Delta}_h = \boldsymbol{\Delta}_h$ . Taking  $\mathbf{r}_h = \mathbf{P}_V \boldsymbol{\Delta}_h$  and  $w_h = P_W \delta_h$ , and using the following inverse inequality

$$\|w_h\|_{\partial K} \leq Ch_K^{-\frac{1}{2}} \|w_h\|_K, \quad \forall w_h \in W_h,$$

leads to

$$\begin{aligned}\|\mathbf{P}_V \Delta_h\|_K &\leq C_\Delta h_K^{-\frac{1}{2}} \|p_h - \widehat{p}_h\|_{\partial K}, \\ \|P_W \delta_h\|_K &\leq C_\delta h_K^{-\frac{1}{2}} \|p_h - \widehat{p}_h\|_{\partial K}.\end{aligned}$$

As

$$\|\delta_h\|_K \leq \|P_W \delta_h\|_K + \|(\text{Id} - P_W) \delta_h\|_K, \quad (97)$$

we can see that the size of the residuals and hence the quality of the approximation depends only on the size of the jump  $p_h - \widehat{p}_h$  (first term of the rhs) and the approximation properties of  $W_h(K)$  (second term of the rhs).

When  $p \in H^s(K)$  with  $s \in \llbracket 0, k+1 \rrbracket$ , we have

$$\|p - P_W p\|_K \leq C h_K^s \|p\|_{s,K},$$

so the penalization term  $\tau(p_h - \widehat{p}_h)$  will ensure the stability of the method. Indeed (97) will therefore lead to

$$\|\delta_h\|_K \leq C_\delta h_K^{-\frac{1}{2}} \|p_h - \widehat{p}_h\|_{\partial K} + C h_K^s \|\delta_h\|_{s,K},$$

and assuming that  $h_K$  is small enough, we have

$$\|\delta_h\|_K \leq C h_K^{-\frac{1}{2}} \|p_h - \widehat{p}_h\|_{\partial K},$$

as the last term of the right-hand side can be absorbed by the left-hand side. A similar result can be obtained for  $\Delta_h$ .

To illustrate this, we have depicted  $\mathcal{E}_{\text{jump}}$  for the HDG methods for a low-Mach flow in [FIGURE 17](#) and for a large-Mach flow in [FIGURE 18](#).

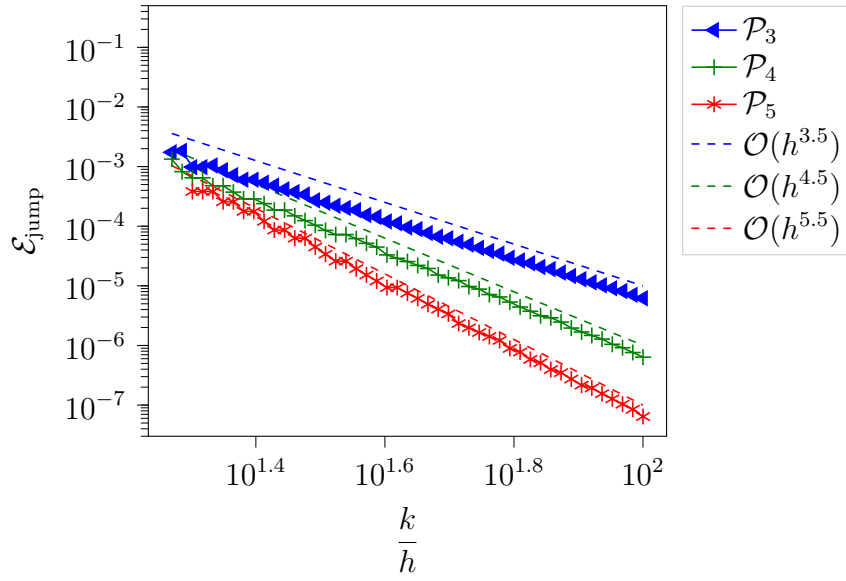
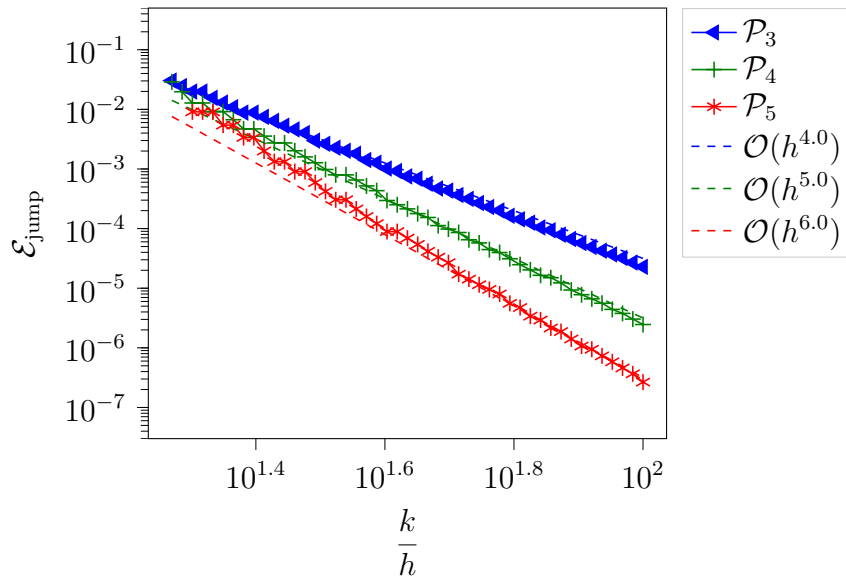
(a) HDG with diffusive flux  $\mathbf{q}_h$ (b) HDG with total flux  $\boldsymbol{\sigma}_h$ 

Figure 17: Low Mach convergence history for the jump error

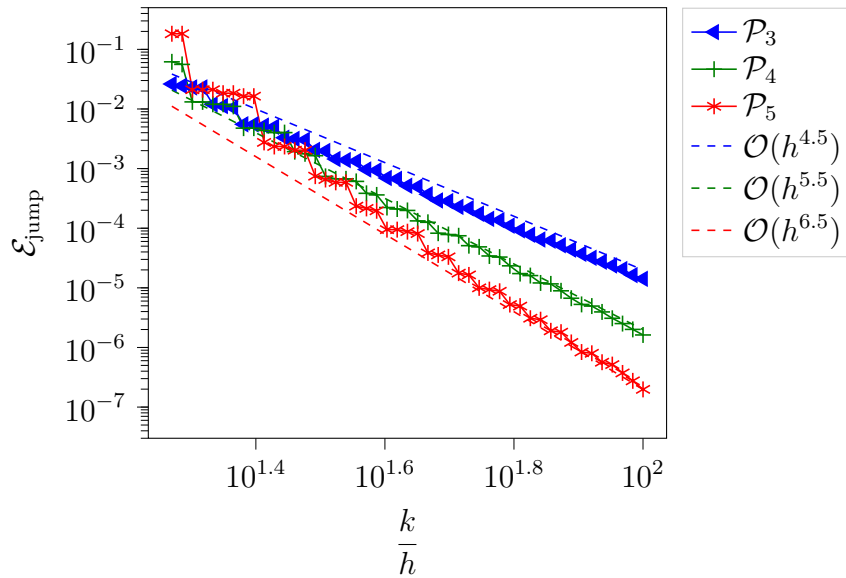
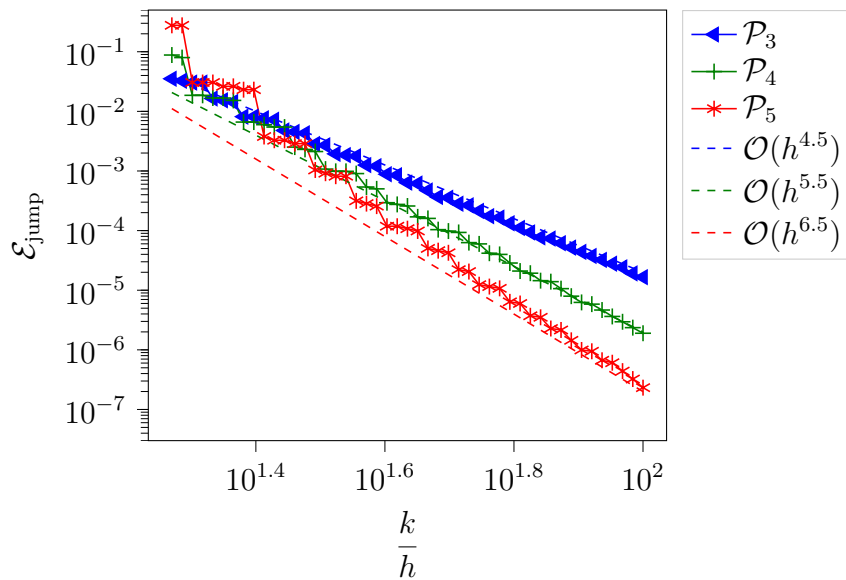

 (a) HDG with diffusive flux  $\mathbf{q}_h$ 

 (b) HDG with total flux  $\sigma_h$ 

Figure 18: Large Mach convergence history for the jump error

**Adaptation to the HDG+ method:** For the HDG+ method, nothing changes for the vectorial residual  $\Delta_h$  but things are a little different for the scalar residual  $\delta_h$ . Indeed this time we have

$$\begin{aligned} \|P_W \delta_h\|_K &\lesssim h_K^{-\frac{1}{2}} \left( \|\tau\| (P_M p_h - \widehat{p}_h)_{\partial K} + \|\tau_{\text{upw}}\| (p_h - \widehat{p}_h)_{\partial K} \right), \\ &\lesssim h_K^{-\frac{1}{2}} \|\tau\| (P_M p_h - \widehat{p}_h)_{\partial K}, \\ (\text{as } \tau = \mathcal{O}(h_K^{-1})) &\lesssim h_K^{-\frac{3}{2}} \|P_M p_h - \widehat{p}_h\|_{\partial K}. \end{aligned}$$

To illustrate this, we have depicted  $\mathcal{E}_{\text{jump}}$  for the HDG+ method for a low-Mach flow in [FIGURE 19](#) and for a large-Mach flow in [FIGURE 20](#). Once again, we can see that these quantities are clearly decreasing.

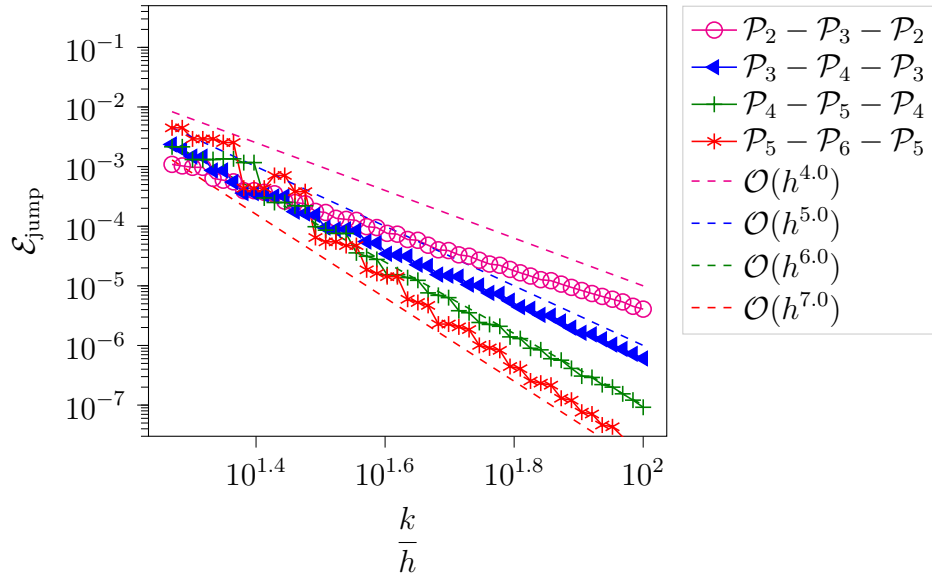


Figure 19: Low Mach convergence history for the jump error for the HDG+ method

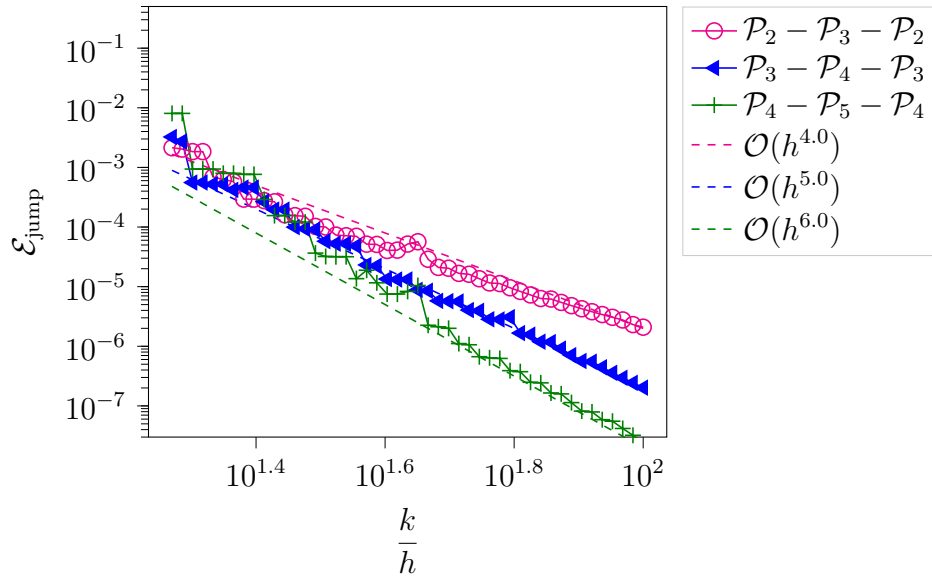


Figure 20: Large Mach convergence history for the jump error for the HDG+ method

### 6.3 Is the upwinding mechanism necessary ?

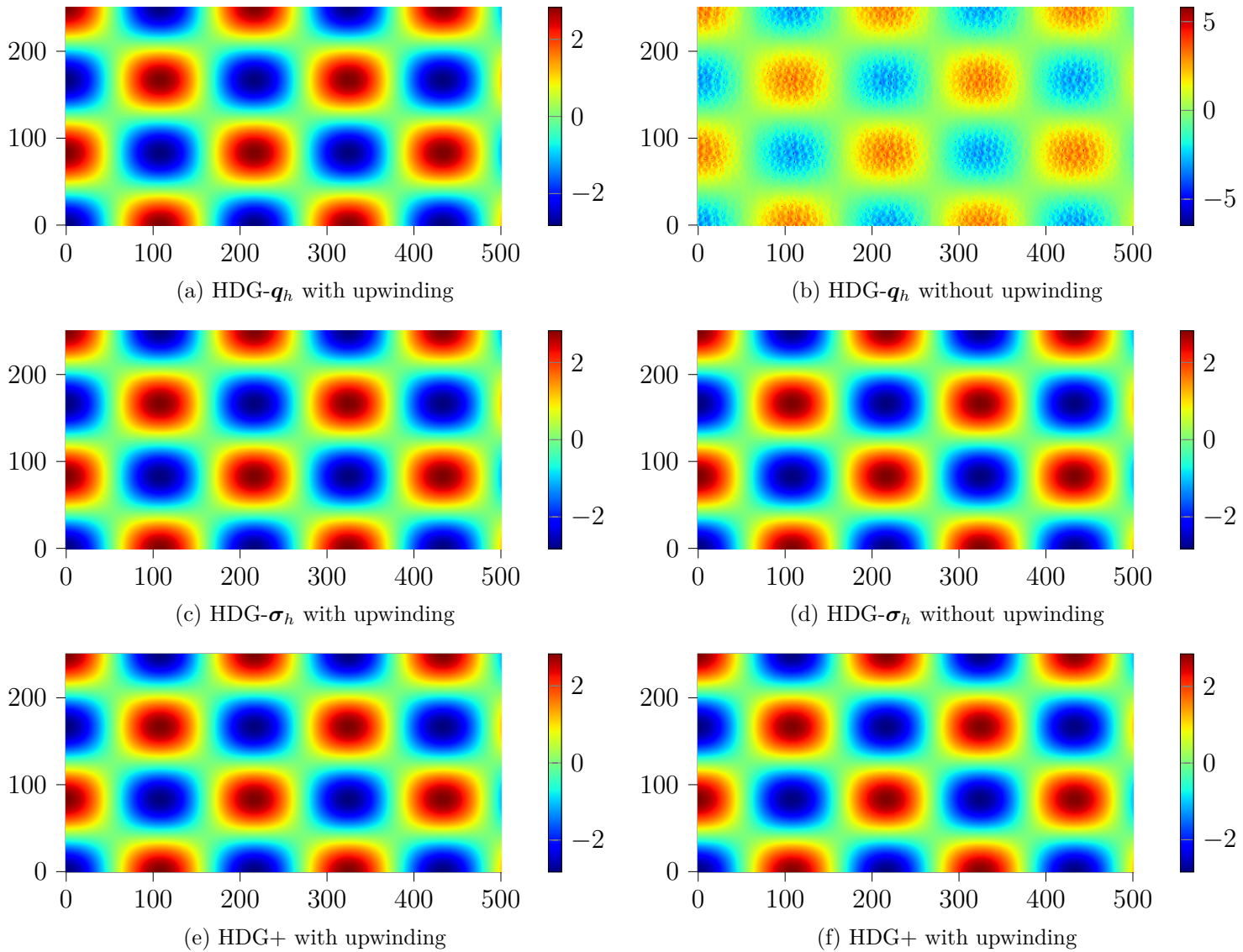
In this section we investigate the numerical impact of upwinding mechanisms. We have performed the same numerical simulations with the three methods with and without those mechanisms. For the HDG- $\mathbf{q}_h$  and the HDG+ methods, deactivating the upwinding mechanism corresponds to taking  $\tau_{\text{upw}} \equiv 0$ . For the HDG- $\boldsymbol{\sigma}_h$  it corresponds to taking  $\tau = 1$  instead of the value given by the Riemann solver.

As it can be seen on [FIGURE 21](#) the HDG- $\mathbf{q}_h$  without upwinding leads to poor numerical results, whereas the two other methods seem to perform well.

Method	With upwinding	Without upwinding
HDG- $\mathbf{q}_h$	$4.3 \cdot 10^{-7}$	0.53
HDG- $\boldsymbol{\sigma}_h$	$4.9 \cdot 10^{-7}$	$5.8 \cdot 10^{-7}$
HDG+	$6.3 \cdot 10^{-8}$	$5.77 \cdot 10^{-7}$

 Table 8: Jump error  $\mathcal{E}_{\text{jump}}$  for the different method

We have also computed the jump error  $\mathcal{E}_{\text{jump}}$  as quality indicator of the numerical solution. The values are given in TABLE 8. Even if the error is higher without upwinding for the HDG+ and HDG- $\boldsymbol{\sigma}_h$ , it seems to remain at a reasonable level. For the HDG+ method, we can conclude that the  $\tau_{\text{upw}}$  penalization seems optional. This can be understood as we need to choose  $\tau = \mathcal{O}(h_K^{-1})$  which is large, so the method seems less sensitive to changes in the penalization. On the other hand, this  $\tau_{\text{upw}}$  penalization is mandatory to make the HDG- $\mathbf{q}_h$  formulation work. Finally, for the HDG- $\boldsymbol{\sigma}_h$  we still recommend to use the upwind penalization parameter  $\tau$  as it leads to a method with no arbitrary choice to make.


 Figure 21:  $p_h$  for the different HDG methods with and without upwind



## 6.4 Point-sources in a uniform flow

For many practical applications it is necessary to consider point-sources. In this section, we show that our method can be used for such computations.

**Settings:** We consider a uniform flow in an infinite plane. We use a Dirac point source  $s = \delta_{(0,0)}$  located at the origin. We write  $\mathbf{b}_0$  as

$$\mathbf{b}_0 = M (\cos \alpha \mathbf{e}_x + \sin \alpha \mathbf{e}_y),$$

where  $M$  is the Mach number,  $\alpha$  the angle between  $\mathbf{b}_0$  and the horizontal axis, and we normalize  $\rho_0, c_0 \equiv 1$ .

The physical domain  $\mathcal{O}_{\text{phys}}$  is surrounded with PMLs, see [FIGURE 22](#) for a sketch of the geometric settings.

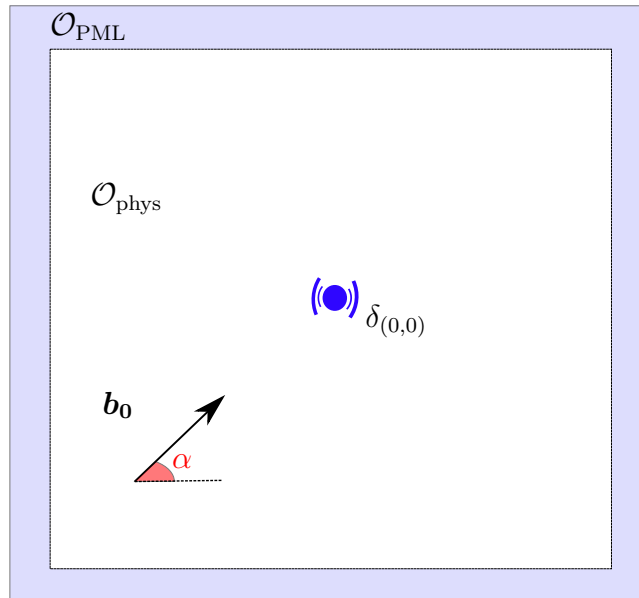


Figure 22: Geometric settings

We have depicted the solution for a large Mach number ( $M = 0.8$ ) in [FIGURE 23](#) and for a low Mach number ( $M = 0.4$ ) in [FIGURE 24](#). The presence of convection leads to a clearly visible Doppler effect : indeed because of the convection, the apparent frequency changes in the domain. The apparent frequency is higher in the bottom-left part of the domain than in the top-right part. The artifacts in the top-right part of the domain are due to the PMLs. Notice that the artifacts are more visible for the large Mach flow.

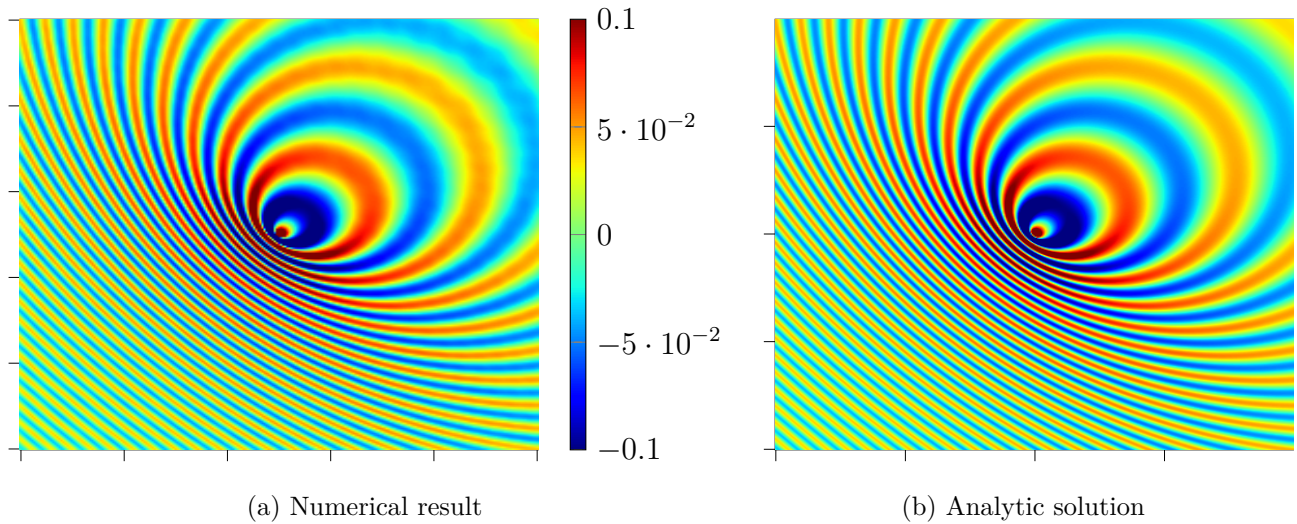


Figure 23:  $p_h$  computed with the HDG- $\sigma_h$  method for  $\omega = 6\pi$ ,  $M = 0.8$  and  $\alpha = \pi/4$ . PMLs are not displayed. The colorbar is the same for the two pictures.

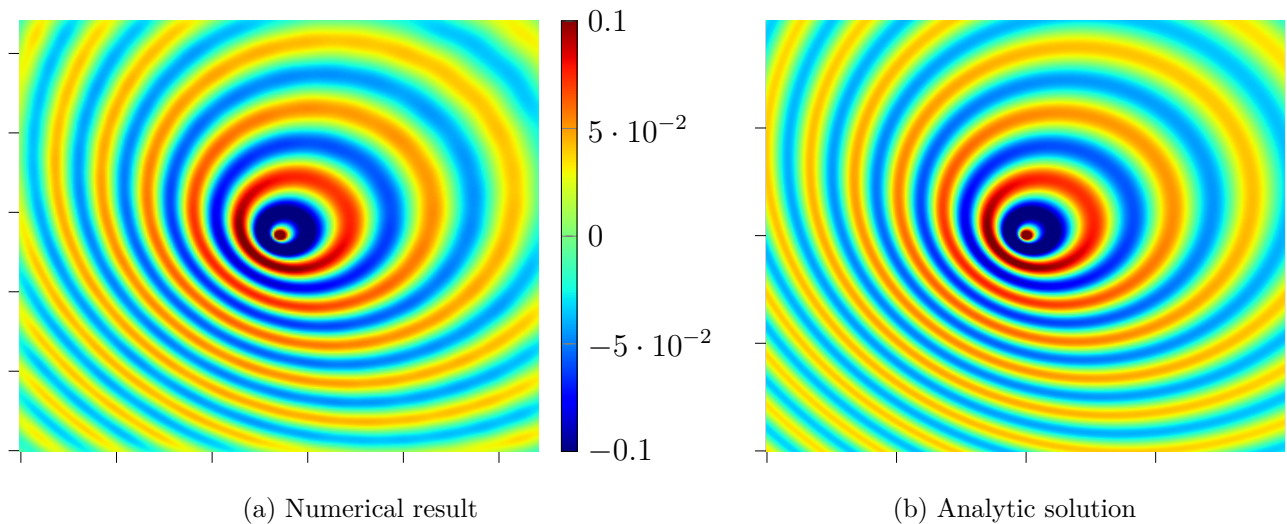


Figure 24:  $p_h$  computed with the HDG- $\sigma_h$  method for  $\omega = 6\pi$ ,  $M = 0.4$  and  $\alpha = \pi/4$ . PMLs are not displayed. The colorbar is the same for the two pictures.

**Analytic solution:** In this context it is possible to write an analytic solution for (1). Following [HPN19], we first need to introduce the so-called *Prandtl-Glauert-Lorentz transformation* in the frequency domain. This transformation maps  $\mathbf{x} = (x, y)$  to  $\tilde{\mathbf{x}} = (\tilde{x}, \tilde{y})$  and  $\omega$  to  $\tilde{\omega}$ , where

$$\begin{cases} \tilde{x} = \left(1 + M^2 \frac{\cos^2 \alpha}{\beta(1 + \beta)}\right) x + M^2 \frac{\cos \alpha \sin \alpha}{\beta(1 + \beta)} y \\ \tilde{y} = M^2 \frac{\cos \alpha \sin \alpha}{\beta(1 + \beta)} x + \left(1 + M^2 \frac{\sin^2 \alpha}{\beta(1 + \beta)}\right) y, \\ \tilde{\omega} = \frac{\omega}{\beta} \end{cases}$$

with  $\beta := \sqrt{1 - M^2}$ . As a shorthand, we write this transformation

$$\tilde{\mathbf{x}} = \mathbf{A}\mathbf{x}, \quad \text{with } \mathbf{A} := \begin{bmatrix} \left(1 + M^2 \frac{\cos^2 \alpha}{\beta(1 + \beta)}\right) & M^2 \frac{\cos \alpha \sin \alpha}{\beta(1 + \beta)} \\ M^2 \frac{\cos \alpha \sin \alpha}{\beta(1 + \beta)} & \left(1 + M^2 \frac{\sin^2 \alpha}{\beta(1 + \beta)}\right) \end{bmatrix}.$$

This Prandtl-Glauert-Lorentz transformation is closely related to the Lorentz transform arising in special relativity. It is well-known that, when the background coefficients are uniform, these transformation maps the convected Helmholtz equation to a standard Helmholtz equation, see *eg.* [MBAG20, HPN19]. This can be understood as the Lorentz transformation was introduced as the transformation between two inertial frames that preserves wave equations, for a deeper insight of the connection between flow acoustics and Lorentzian geometry we refer to [Vis98]. The analytic solution is given by

$$p_{\text{exact}}(\mathbf{x}, \omega) = -\frac{i}{4\beta} H_0^{(1)}\left(\frac{\omega}{\beta} |\mathbf{A}\mathbf{x}|\right) \exp\left[\frac{i\omega}{\beta} \mathbf{A}\mathbf{x} \cdot \mathbf{b}_0\right]$$

where  $H_0^{(1)}$  is the Hankel function of the first kind of order 0 and

$$\tilde{r} := \sqrt{\tilde{x}^2 + \tilde{y}^2}.$$

Even if there is an analytic solution in this case, we were not able to obtain meaningful convergence plots due to the bad quality of the PMLs.

**Computational cost:** In TABLE 9 we have written down the sizes of the linear systems to solve for the different HDG methods using an interpolation of order 5 for the trace variable. To give a reference, we also added the size of the system obtained when solving the convected Helmholtz equation (1) with a continuous finite element method (CG) with same interpolation degree using the `montjoie` solver<sup>9</sup>. We have also added the size of the system obtained when solving the standard Helmholtz equation with a Local Discontinuous Galerkin method (LDG), which is a first-order DG method, using the `montjoie` solver (as LDG methods are not implemented for the convected Helmholtz equation in `montjoie`).

Method	HDG- $\mathbf{q}_h$	HDG- $\boldsymbol{\sigma}_h$	HDG+	CG	LDG
$k$	5	5	4	5	5
nnz	12 237 750	12 237 750	6 243 750	1 353 750	46 862 808
nnz LU	67 515 161	67 515 161	34 244 845	39 008 186	259 272 979
$\mathcal{E}_{\text{jump}}$	$1.4 \cdot 10^{-3}$	$1.0 \cdot 10^{-3}$	$8.9 \cdot 10^{-4}$		

Table 9: Size of the linear system to solve for the Dirac in a uniform flow

We can see that using HDG instead of LDG leads to significantly smaller linear systems. For the convected Helmholtz equation, it is possible to use the CG method which less expensive than the HDG ones, see [CD16]. However, we would like to point out that the CG method is known to give bad numerical results for more realistic aeroacoustic models such as Galbrun's equation, see [CD18].

<sup>9</sup>`montjoie` is a versatile and well-tested high-order finite element solver. For more informations about the numerical method used to solve the convected Helmholtz equation, see <http://montjoie.gforge.inria.fr/helmholtz.php>.

**Local refinement:** To handle the point-sources, a mesh with a local refinement around the source should be used, otherwise artifacts could be present in the numerical solution. With the HDG method, those artifacts seem really limited when no local refinement is used, see [FIGURE 25](#).

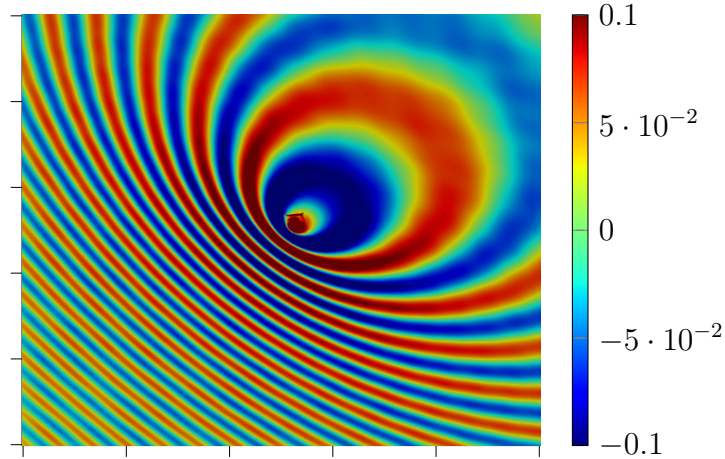


Figure 25: Point-source without local refinement

**Illustrative example:** To show that our method can handle more complex simulations, we consider the same test-case as before but with two point-sources located near the origin. On [FIGURE 26](#), we can still see the changes in the apparent frequency due to the Doppler effect, but also interference patterns due to the interactions between the two sources.

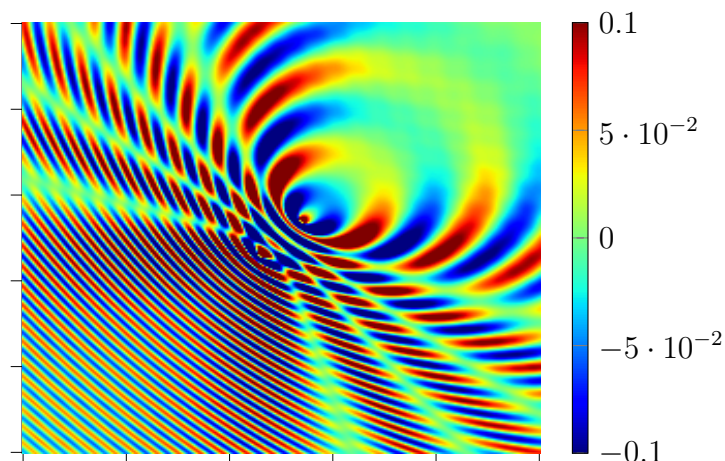


Figure 26: Interference between two point-sources

## 6.5 Gaussian jet

In this section we focus on a space-varying flow. We will use a *gaussian jet*, which is a common test-case in the literature, see *eg.* [\[MMMP17, Sec. 3\]](#). To work with parameters and unknowns without dimension, we choose

$$\rho_0 = c_0 \equiv 1,$$

and we consider the following gaussian jet flow

$$\mathbf{v}_0 = M_0(y)\mathbf{e}_x, \quad \text{where} \quad M_0(y) := M_\infty + \mu \exp\left(-\frac{y^2}{R^2}\right),$$

which is a gaussian perturbation of the uniform flow  $\mathbf{v}_0 = M_\infty \mathbf{e}_x$ .

For this simulation, we work with the following values for the parameters

$$\mathcal{O}_{\text{phys}} = (0, 3) \times (-1, 1) \ ; \ \omega = 6\pi \ ; \ s = \delta_{(1,0.5)} \ ;$$

and

$$M_\infty = 0.1 \ ; \ \mu = 0.3 \ ; \ R = 0.35.$$

The physical domain  $\mathcal{O}_{\text{phys}}$  is surrounded by PMLs.

Notice that these kind of jet-flows are not potential flows and should therefore not be used with the convected Helmholtz equation. However, as discussed in [MMMP17], with this choice of parameters the vorticity of the flow stays small and can be neglected. The convected Helmholtz equation is therefore a good approximation of the more realistic models.

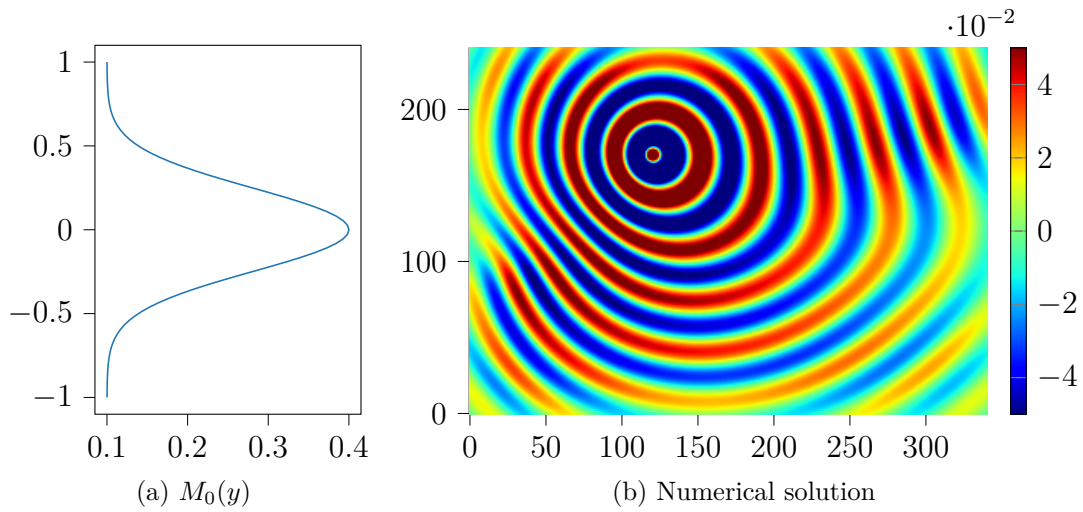


Figure 27:  $\Re p_h$  for the gaussian jet, obtained with the HDG+ method

On [FIGURE 27](#), we can see that there is a phase-shift inside the jet, as it is expected : this refraction-like effect can be seen in the center-left and center-right parts of the domain at the limit between the jet and uniform flows. There is also a small Doppler effect as the apparent frequency changes in the domain : it is different in the the bottom-left part of the domain and in the center-right one. However this effect is less obvious than in the previous example as the Mach number is significantly lower.

## Conclusion

In this paper, we have introduced three HDG methods to solve the convected Helmholtz equation. Two of them are standard HDG methods that use the same polynomial degree for the approximation of all the unknowns. The third one is less standard and uses a higher polynomial degree for the scalar unknown and a reduced stabilization process.

For all of those methods, detailed theoretical results on convergence and well-posedness are provided. It is important to note that we could not obtain the super-convergence property for the two standard HDG methods because of convection. We also provided numerical experiments that are consistent with the absence of super-convergence. The HDG+ method achieves optimal convergence. Due to the reduced stabilization this leads to a "super-convergence like behaviour" as we obtain a convergence rate of  $k + 2$  for the cost of a HDG method of degree  $k$  without post-processing.

During the numerical experiments, it occurred to us that the HDG+ and HDG- $\sigma_h$  methods seemed more robust than the HDG- $q_h$  method. In particular, they are less sensitive to the choice of penalization parameter.

**Future work:** Now that the use of HDG method has been validated for the convected Helmholtz equation, we will generalize this work to more realistic aeroacoustic models such as Galbrun's equation or Goldstein's equation.

We also noticed that the use of PMLs may lead to bad numerical results, and we will therefore address the construction of absorbing boundary conditions for the convected Helmholtz equation in a future paper.

**Open-source implementation:** An open-source implementation of those three methods in the `hawn` solver will be released soon.

**Acknowledgement:** The authors would like to thank Florian Faucher for his help with the numerical implementation.

Experiments presented in this paper were carried out using the PlaFRIM experimental testbed, supported by Inria, CNRS (LABRI and IMB), Université de Bordeaux, Bordeaux INP and Conseil Régional d'Aquitaine (see <https://www.plafrim.fr>).

Nathan Rouxelin acknowledges financial support from e2s-UPPA (see <https://e2s-uppa.eu>).

## A Intermediate results for the error analysis of the HDG method with diffusive flux

In this appendix we state the intermediate results allowing to adapt the analysis conducted in [SUBSECTION 4.3](#) to prove the convergence of the HDG method stated in [THEOREM 6](#). The proofs are omitted but are very similar to the ones given in [SUBSECTION 4.3](#).

The main differences to keep in mind are :  $s \in [1, k + 1]$  instead of  $s \in [1, k + 2]$  and  $\tau = \mathcal{O}(1)$  instead of  $\tau = \mathcal{O}(h^{-1})$ .

**Gradient estimate:** this corresponds to [LEMMA 4.3](#)

### Lemma A.1:

The following estimate holds

$$\|\nabla \varepsilon_h^p\|_{\mathcal{T}_h} \lesssim \|\varepsilon_h^q\|_{\mathbf{w}_0, \mathcal{T}_h} + \|\delta_h^q\|_{\mathcal{T}_h} + h^{-\frac{1}{2}} \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial \mathcal{T}_h}$$

**Energy-like equality:** this corresponds to [LEMMA 4.5](#).

**Lemma A.2:**

The following energy-like equality holds

$$\begin{aligned} & \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h}^2 - \omega^2 \|\varepsilon_h^p\|_{\rho_0, K}^2 - 2i\omega \left( \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h}^2 + \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h}^2 \right) \\ &= -\omega^2 (\rho_0 \delta_h^p, \varepsilon_h^p)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} + (\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\delta}_h^q)_{\mathcal{T}_h} + 2\omega \Im (\varepsilon_h^p \mathbf{b}_0, \nabla \varepsilon_h^p)_{\mathcal{T}_h} \\ &+ \left\langle \boldsymbol{\delta}_h^q \cdot \mathbf{n} + 2i\omega (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p - 2i\omega \tau (\delta_h^p - \widehat{\delta}_h^p) + 2i\omega \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \varepsilon_h^p - \widehat{\varepsilon}_h^p \right\rangle_{\partial\mathcal{T}_h} \end{aligned}$$

Furthermore if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s, t \in [1, k+1]$  then the following estimate holds

$$\begin{aligned} & \left| \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h}^2 - 2i\omega \left( \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h}^2 + \left\| \left( \frac{1}{2} |\mathbf{b}_0 \cdot \mathbf{n}| \right)^{\frac{1}{2}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p) \right\|_{\partial\mathcal{T}_h}^2 \right) \right| \\ & \lesssim \omega^2 \|\varepsilon_h^p\|_{\mathcal{T}_h}^2 + \omega \|\varepsilon_h^p\|_{\mathcal{T}_h} \left( \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} + h^{-\frac{1}{2}} \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h} + \omega h^s \|p\|_{s, \mathcal{O}} \right) \\ & + \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h} \left( h^t \|\mathbf{q}\|_{t, \mathcal{O}} + \omega h^s \|p\|_{s, \mathcal{O}} \right) + h^{2t} \|\mathbf{q}\|_{t, \mathcal{O}}^2 \\ & + h^{-\frac{1}{2}} \|\varepsilon_h^p - \widehat{\varepsilon}_h^p\|_{\partial\mathcal{T}_h} \left( \omega h^s \|p\|_{s, \mathcal{O}} + h^t \|\mathbf{q}\|_{t, \mathcal{O}} \right) \end{aligned}$$

**Dual identity:** this corresponds to [LEMMA 4.6](#).

**Lemma A.3:**

The following dual identity holds

$$\begin{aligned} \|\varepsilon_h^p\|_{\mathcal{T}_h}^2 &= -(\mathbf{W}_0 \boldsymbol{\varepsilon}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi} - \boldsymbol{\xi})_{\mathcal{T}_h} + \omega^2 (\rho_0 \varepsilon_h^p, \pi_W \theta - \theta)_{\mathcal{T}_h} + 2i\omega (\nabla \varepsilon_h^p, (\pi_W \theta - \theta) \mathbf{b}_0)_{\mathcal{T}_h} \\ & - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \varepsilon_h^p, \pi_W \theta - \theta \rangle_{\partial\mathcal{T}_h} \\ & + (\mathbf{W}_0 \boldsymbol{\delta}_h^q, \boldsymbol{\pi}_V \boldsymbol{\xi})_{\mathcal{T}_h} - \omega^2 (\rho_0 \delta_h^p, \pi_W \theta)_{\mathcal{T}_h} + 2i\omega (\delta_h^p \mathbf{b}_0, \nabla (\pi_W \theta))_{\mathcal{T}_h} \\ & + 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\varepsilon}_h^p - \tau (\varepsilon_h^p - \widehat{\varepsilon}_h^p) + \tau_{\text{upw}} (\varepsilon_h^p - \widehat{\varepsilon}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial\mathcal{T}_h} \\ & - 2i\omega \langle (\mathbf{b}_0 \cdot \mathbf{n}) \widehat{\delta}_h^p - \tau (\delta_h^p - \widehat{\delta}_h^p) + \tau_{\text{upw}} (\delta_h^p - \widehat{\delta}_h^p), \pi_W \theta - P_M \theta \rangle_{\partial\mathcal{T}_h} + \langle \boldsymbol{\delta}_h^q \cdot \mathbf{n}, \pi_W \theta - P_M \theta \rangle_{\partial\mathcal{T}_h} \end{aligned}$$

**Dual estimate:** this corresponds to [LEMMA 4.7](#).

**Lemma A.4:**

Assuming that the regularity assumption (72) holds and that  $\omega^2 h^2 \|\rho_0\|_{\infty} C_{\text{reg}} C$  (where  $C$  is the constant of [THEOREM 4](#)) is small enough, if  $p \in H^s(\mathcal{O})$  and  $\mathbf{q} \in \mathbf{H}^t(\mathcal{O})$  where  $s, t \in [1, k+1]$  then

$$\|\varepsilon_h^p\|_{\mathcal{T}_h} \lesssim h^{t+1} (1 + \omega) \|\mathbf{q}\|_{t, \mathcal{O}} + h^s (1 + \omega + \omega^2) \|p\|_{s, \mathcal{O}} + h(1 + \omega) \|\boldsymbol{\varepsilon}_h^q\|_{\mathbf{W}_0, \mathcal{T}_h}$$

where

$$h := \max_{K \in \mathcal{T}_h} h_K$$

## References

- [ALA13] Natalia C. B. Arruda, Abimael F. D. Loula, and Regina C. Almeida. Locally discontinuous but globally continuous Galerkin methods for elliptic problems. *Com-*

*puter Methods in Applied Mechanics and Engineering*, 255:104–120, March 2013. Cited on page 9.

- [BCDL15] Marie Bonnasse-Gahot, Henri Calandra, Julien Diaz, and Stéphane Lanteri. Hybridizable Discontinuous Galerkin method for the simulation of the propagation of the elastic wave equations in the frequency domain. Research Report RR-8990, INRIA Bordeaux ; INRIA Sophia Antipolis - Méditerranée, June 2015. Cited on pages 4 and 49.
- [BDL03] Eliane Bécache, Anne-Sophie Bonnet-Ben Dhia, and Guillaume Legendre. Perfectly matched layers for the convected Helmholtz equation. Technical report, 2003. Cited on page 62.
- [BDMP21] Hélène Barucq, Julien Diaz, Rose-Cloé Meyer, and Ha Pham. Implementation of hybridizable discontinuous Galerkin method for time-harmonic anisotropic poroelasticity in two dimensions. *International Journal for Numerical Methods in Engineering*, 122(12):3015–3043, 2021. Cited on page 4.
- [CC12] Yanlai Chen and Bernardo Cockburn. Analysis of variable-degree HDG methods for convection–diffusion equations. Part I: General nonconforming meshes. *IMA Journal of Numerical Analysis*, 32(4):1267–1293, October 2012. Cited on page 4.
- [CC14] Yanlai Chen and Bernardo Cockburn. Analysis of variable-degree HDG methods for convection-diffusion equations. Part II: Semimatching nonconforming meshes. *Mathematics of Computation*, 83(285):87–111, January 2014. Cited on page 4.
- [CD16] Juliette Chabassier and Marc Duruffé. High Order Finite Element Method for solving Convected Helmholtz equation in radial and axisymmetric domains. Application to Helioseismology. Research Report RR-8893, Inria Bordeaux Sud-Ouest, March 2016. Cited on page 80.
- [CD18] Juliette Chabassier and Marc Duruffé. Solving time-harmonic Galbrun’s equation with an arbitrary flow. Application to Helioseismology. *Rapport de recherche Inria*, 2018. Cited on page 80.
- [CDG<sup>+</sup>09] Bernardo Cockburn, Bo Dong, Johnny Guzmán, Marco Restelli, and Riccardo Sacco. A Hybridizable Discontinuous Galerkin Method for Steady-State Convection-Diffusion-Reaction Problems. *SIAM Journal on Scientific Computing*, 31(5):3827–3846, January 2009. Cited on page 4.
- [CDPE16] Bernardo Cockburn, Daniele A. Di Pietro, and Alexandre Ern. Bridging the hybrid high-order and hybridizable discontinuous Galerkin methods. *ESAIM: Mathematical Modelling and Numerical Analysis*, 50(3):635–650, May 2016. Cited on page 4.
- [CGL09] Bernardo Cockburn, Jayadeep Gopalakrishnan, and Raytcho Lazarov. Unified Hybridization of Discontinuous Galerkin, Mixed, and Continuous Galerkin Methods for Second Order Elliptic Problems. *SIAM J. Numer. Anal.*, 47:1319–1365, August 2009. Cited on page 4.
- [CGS10] Bernardo Cockburn, Jayadeep Gopalakrishnan, and Francisco-Javier Sayas. A projection-based error analysis of HDG methods. *Mathematics of Computation*, 79(271):1351–1367, March 2010. Cited on pages 4, 13, 26, and 48.



- [CLOS20] Liliana Camargo, Bibiana López-Rodríguez, Mauricio Osorio, and Manuel Solano. An HDG method for Maxwell’s equations in heterogeneous media. *Computer Methods in Applied Mechanics and Engineering*, 368:113178, August 2020. Cited on page 4.
- [Coc14] Bernardo Cockburn. Static Condensation, Hybridization, and the Devising of the HDG Methods. *springerprofessional.de*, 2014. Cited on page 4.
- [CQS18] Huangxin Chen, Weifeng Qiu, and Ke Shi. A priori and computable a posteriori error estimates for an HDG method for the coercive Maxwell equations. *Computer Methods in Applied Mechanics and Engineering*, 333:287–310, May 2018. Cited on page 4.
- [CQSS17] Huangxin Chen, Weifeng Qiu, Ke Shi, and Manuel Solano. A Superconvergent HDG Method for the Maxwell Equations. *Journal of Scientific Computing*, 70(3):1010–1029, March 2017. Cited on page 4.
- [CS13] Bernardo Cockburn and Ke Shi. Superconvergent HDG methods for linear elasticity with weakly symmetric stresses | IMA Journal of Numerical Analysis | Oxford Academic. *IMA Journal of Numerical Analysis*, 2013. Cited on page 4.
- [CZ12] Bernardo Cockburn and Wujun Zhang. A Posteriori Error Estimates for HDG Methods. *Journal of Scientific Computing*, 51(3):582–607, June 2012. Cited on page 72.
- [CZ13] Bernardo Cockburn and Wujun Zhang. A Posteriori Error Analysis for Hybridizable Discontinuous Galerkin Methods for Second Order Elliptic Problems. *SIAM Journal on Numerical Analysis*, 51(1):676–693, January 2013. Cited on page 72.
- [DS19] Shukai Du and Francisco-Javier Sayas. *An Invitation to the Theory of the Hybridizable Discontinuous Galerkin Method: Projections, Estimates, Tools*. SpringerBriefs in Mathematics. Springer International Publishing, Cham, 2019. Cited on pages 4, 21, 24, 25, 26, 34, and 41.
- [Dub91] Moshe Dubiner. Spectral methods on triangles and other domains. *Journal of Scientific Computing*, 6(4):345–390, December 1991. Cited on page 56.
- [EG04] Alexandre Ern and Jean-Luc Guermond. *Theory and Practice of Finite Elements*. Applied Mathematical Sciences. Springer-Verlag, New York, 2004. Cited on pages 8, 12, 23, 33, 35, and 41.
- [Fau21] Florian Faucher. ‘hawen’: Time-harmonic wave modeling and inversion using hybridizable discontinuous Galerkin discretization. *Journal of Open Source Software*, 6(57):2699, January 2021. Cited on page 4.
- [FCS15] G. Fu, B. Cockburn, and H. Stolarski. Analysis of an HDG method for linear elasticity. *International Journal for Numerical Methods in Engineering*, 102(3-4):551–575, 2015. Cited on page 4.
- [FLd14] Cristiane O. Faria, Abimael F. D. Loula, and Antônio J. B. dos Santos. Primal stabilized hybrid and DG finite element methods for the linear elasticity problem. *Computers & Mathematics with Applications*, 68(4):486–507, August 2014. Cited on page 9.

- [FS20] Florian Faucher and Otmar Scherzer. Adjoint-state method for Hybridizable Discontinuous Galerkin discretization, application to the inverse acoustic wave problem. *Computer Methods in Applied Mechanics and Engineering*, 372:113406, December 2020. Cited on pages 4 and 49.
- [GM11] Roland Griesmaier and Peter Monk. Error Analysis for a Hybridizable Discontinuous Galerkin Method for the Helmholtz Equation. *Journal of Scientific Computing*, 49(3):291–310, December 2011. Cited on page 4.
- [GSV18] Jay Gopalakrishnan, Manuel Solano, and Felipe Vargas. Dispersion Analysis of HDG Methods. *Journal of Scientific Computing*, 77(3):1703–1735, December 2018. Cited on page 4.
- [HPN19] Fang Q. Hu, Michelle E. Pizzo, and Douglas M. Nark. On the use of a Prandtl-Glauert-Lorentz transformation for acoustic scattering by rigid bodies with a uniform flow. *Journal of Sound and Vibration*, 443:198–211, March 2019. Cited on pages 79 and 80.
- [HPS17] Allan Hungria, Daniele Prada, and Francisco-Javier Sayas. HDG methods for elastodynamics. *Computers & Mathematics with Applications*, 74(11):2671–2690, December 2017. Cited on page 4.
- [Hun19] Allan Hungria. *Using HDG+ to Compute Solutions of the 3D Linear Elastic and Poroelastic Wave Equations*. PhD thesis, University of Delaware, 2019. Cited on pages 4, 34, and 41.
- [HW08] Jan S. Hesthaven and Tim Warburton. *Nodal Discontinuous Galerkin Methods: Algorithms, Analysis, and Applications*. Texts in Applied Mathematics. Springer-Verlag, New York, 2008. Cited on page 16.
- [Kir04] Robert C. Kirby. Algorithm 839: FIAT, a new paradigm for computing finite element basis functions. *ACM Transactions on Mathematical Software (TOMS)*, 30(4):502–516, December 2004. Cited on page 56.
- [KSC12] Robert M. Kirby, Spencer J. Sherwin, and Bernardo Cockburn. To CG or to HDG: A Comparative Study. *Journal of Scientific Computing*, 51(1):183–212, April 2012. Cited on page 4.
- [Leh10] Christoph Lehrenfeld. *Hybrid Discontinuous Galerkin Methods for Solving Incompressible Flow Problems*. PhD thesis, RWTH Aachen, 2010. Cited on pages 4 and 28.
- [MBAG20] Philippe Marchner, Hadrien Beriot, Xavier Antoine, and Christophe Geuzaine. Stable Perfectly Matched Layers with Lorentz transformation for the convected Helmholtz equation. Technical report, 2020. Cited on page 80.
- [MMMP17] Jean-François Mercier, Colin Mietka, Florence Millot, and Vincent Pagneux. Acoustic propagation in a vortical homentropic flow. page 19, 2017. Cited on pages 81 and 82.
- [NPRC15] N. C. Nguyen, J. Peraire, F. Reitich, and B. Cockburn. A phase-based hybridizable discontinuous Galerkin method for the numerical solution of the Helmholtz equation. *Journal of Computational Physics*, 290:318–335, June 2015. Cited on page 4.

- [Oik14] Issei Oikawa. A hybridized discontinuous Galerkin method with reduced stabilization. *arXiv:1405.2491 [math]*, November 2014. Cited on pages 4 and 56.
- [Oik16] Issei Oikawa. Analysis of a Reduced-Order HDG Method for the Stokes Equations. *Journal of Scientific Computing*, 67(2):475–492, May 2016. Cited on page 4.
- [Oik18] Issei Oikawa. An HDG Method with Orthogonal Projections in Facet Integrals. *Journal of Scientific Computing*, 76(2):1044–1054, August 2018. Cited on page 4.
- [PE12] Daniele Antonio Di Pietro and Alexandre Ern. *Mathematical Aspects of Discontinuous Galerkin Methods*. Mathématiques et Applications. Springer-Verlag, Berlin Heidelberg, 2012. Cited on page 13.
- [Pie90] Allan D. Pierce. Wave equation for sound in fluids with unsteady inhomogeneous flow. *The Journal of the Acoustical Society of America*, 87(6):2292–2299, June 1990. Cited on page 5.
- [QS16a] Weifeng Qiu and Ke Shi. An HDG Method for Convection Diffusion Equation. *Journal of Scientific Computing*, 66(1):346–357, January 2016. Cited on pages 4, 28, and 41.
- [QS16b] Weifeng Qiu and Ke Shi. A superconvergent HDG method for the incompressible Navier–Stokes equations on general polyhedral meshes. *IMA Journal of Numerical Analysis*, 36(4):1943–1967, October 2016. Cited on page 4.
- [QSS16] Weifeng Qiu, Jiguang Shen, and Ke Shi. An HDG method for linear elasticity with strong symmetric stresses. *arXiv:1312.1407 [math]*, February 2016. Cited on pages 4 and 41.
- [SM50] Jack Sherman and Winifred J. Morrison. Adjustment of an Inverse Matrix Corresponding to a Change in One Element of a Given Matrix. *Annals of Mathematical Statistics*, 21(1):124–127, March 1950. Cited on page 6.
- [Ste91] Rolf Stenberg. Postprocessing schemes for some mixed finite elements. *ESAIM: Mathematical Modelling and Numerical Analysis*, 25(1):151–167, 1991. Cited on pages 26 and 48.
- [Vis98] Matt Visser. Acoustic black holes: Horizons, ergospheres, and Hawking radiation. *Classical and Quantum Gravity*, 15(6):1767–1791, June 1998. Cited on page 80.
- [YMKS16] Sergey Yakovlev, David Moxey, Robert M. Kirby, and Spencer J. Sherwin. To CG or to HDG: A Comparative Study in 3D. *Journal of Scientific Computing*, 67(1):192–220, April 2016. Cited on page 4.



**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

200 avenue de la Vieille Tour  
33405 Talence Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399