



**HAL**  
open science

# A unified approach for handling convection terms in finite volumes and mimetic discretization methods for elliptic problems

Lourenco Beirao da Veiga, Jerome Droniou, Gianmarco Manzini

## ► To cite this version:

Lourenco Beirao da Veiga, Jerome Droniou, Gianmarco Manzini. A unified approach for handling convection terms in finite volumes and mimetic discretization methods for elliptic problems. IMA Journal of Numerical Analysis, 2011, 31 (4), pp.1357-1401. 10.1093/imanum/drq018 . hal-00808695

**HAL Id: hal-00808695**

**<https://hal.science/hal-00808695>**

Submitted on 8 Apr 2013

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

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

# A unified approach to handle convection terms in Finite Volumes and Mimetic Discretization Methods for elliptic problems

Lourenço Beirão da Veiga<sup>1</sup>, Jérôme Droniou<sup>2 3</sup>, Gianmarco Manzini<sup>4</sup>.  
February 23, 2013

**Abstract** We study the numerical approximation to the solution of the steady convection-diffusion equation. The diffusion term is discretized by using the Hybrid Mimetic Method (HMM), which is the unified formulation for the Hybrid Finite Volume Method, the Mixed Finite Volume Method and the Mimetic Finite Difference Method recently proposed in [33]. In such a setting, we discuss several techniques to discretize the convection term, which are mainly adapted from the literature of Finite Volume or Finite Element schemes. For this family of schemes, we provide a full proof of convergence under very general regularity conditions of the solution field, and derive an error estimate when the scalar solution is in  $H^2(\Omega)$ . Finally, we compare the performance of these schemes on a set of test cases selected from the literature in order to document the accuracy of the numerical approximation in both diffusion and convection-dominated regimes. Moreover, we numerically investigate the behavior of these methods in the approximation of solutions with boundary layers or internal regions with strong gradients.

## 1 Introduction

Many physical models of fluid flows involve partial differential equations (PDEs) with both convection and diffusion terms such as the Navier-Stokes equations, flows in porous media, etc. Analytical solutions are not normally available for real applications and numerical approximations must be devised in some way. To this purpose, efficient numerical schemes based on Finite and Mixed Finite Element and 2-points Finite Volumes have been developed for the numerical treatment of the diffusive part of the equation. In such a framework, a great amount of work has been done to investigate the connections between the lowest-order Raviart-Thomas Mixed Finite Element ( $RT_0 - P_0$ ) methods and various cell-centered Finite Volume and Finite Difference numerical formulations on meshes of simplexes and quadrilaterals/hexahedrons. The relationship between the Mixed Finite Element method and cell-centered Finite Difference on rectangular meshes was first established in [45], and further developed in subsequent papers, see for example [3]. Basically, it can be shown that applying appropriate quadrature rules to the numerical formulation in the  $RT_0$  space on rectangles, the vector variable (the velocity) is eliminated thus reducing the method to a positive definite cell-centered Finite Difference method for the scalar variable (the pressure). Using this approach, classical cell-centered Finite Difference methods on rectangular meshes are easily retrieved based on a 9-points stencil for full tensor coefficients and a 5-points stencil for scalar (diagonal) tensor. Similar results are also obtainable on regular hexahedron meshes. These developments led to the formulation of enhanced cell-centered Finite Differences, cf. [3], that can handle general shape elements (triangles, quadrilaterals and hexahedra) and are suitable to full tensor coefficients. A similar relationship exists between the  $RT_0 - P_0$  scheme and the 2-point Finite Volume formulation on triangular meshes using triangle circumcenters. This relationship was originally established by [6] for two-dimensional diffusion problems with scalar coefficients. This approach has been further developed by [49], which investigates the case of a full diffusion tensor in two and three dimensions on meshes of simplexes.

Nonetheless, practical situations, such as those encountered in petroleum engineering, require computational grids that are not structured or simple enough to make use of the methods mentioned above.

---

<sup>1</sup>Dipartimento di Matematica “F. Enriques”, Università degli Studi di Milano, via Saldini 50, I – 20133 Milano, Italy; email: [lourenco.beirao@unimi.it](mailto:lourenco.beirao@unimi.it)

<sup>2</sup>Université Montpellier 2, Institut de Mathématiques et de Modélisation de Montpellier, CC 051, Place Eugène Bataillon, 34095 Montpellier cedex 5, France; email: [droniou@math.univ-montp2.fr](mailto:droniou@math.univ-montp2.fr)

<sup>3</sup>Funded by GDR MOMAS CNRS/PACEN and Project VFSitCom (ANR-08-BLAN-0275-01).

<sup>4</sup>Istituto di Matematica Applicata e Tecnologie Informatiche (IMATI) – CNR, via Ferrata 1, I – 27100 Pavia, Italy; email: [manzini@imati.cnr.it](mailto:manzini@imati.cnr.it)

Thus, alternative and more sophisticated techniques have been developed in the last decade to approximate the solution to diffusive equations on general grids. In this framework, we mention, for instance, the Discontinuous Galerkin Method, [5, 44] and references therein, the Multi-Point Flux Approximation [1, 2, 48], the Mimetic Finite Difference Method, [8, 11–14, 18, 20–23, 37, 40, 41] and references therein, the Hybrid Finite Volume Method [36], and the Mixed Finite Volume Method [25, 31, 32]. Strict correlations also exist among these numerical approximations and with respect to the lowest order Mixed Finite Element Method, and it is not surprising that sometimes the lowest-order schemes may belong to more than one of these families of methods. For example, the first-order Discontinuous Galerkin scheme can be easily re-interpreted as a Finite Volume method. The lowest-order Raviart-Thomas scheme on grids of simplexes (triangles in 2-D, tetrahedrons in 3-D) is a member of the family of MFD methods, cf. [23]. Note, however, that on meshes of quadrilaterals and hexaedrons no connection has been established yet between the MFD method in mixed form and the Mixed Finite Element method. Again, we mention [47] that outlined the relationship existing between the Multi-Point Flux Approximation and the Mixed Finite Element Method.

A remarkable fact has been recently discovered in [33]: after some generalization a unified formulation exists for three of the methods cited above, i.e., the Hybrid Finite Volume method, the Mixed Finite Volume method and the Mimetic Finite Difference method. Consequently, these three methods are members of the same family of discretization techniques. Following [33], we will refer to such a family of numerical methods by the wording *Hybrid Mimetic Mixed* methods, or by the abbreviation HMM.

Since the HMM method is at the juncture of two different frameworks, namely the Mimetic/Finite Element and the Finite Volume ones, the convective term can be naturally discretized using quite different techniques depending on the adopted point of view on the scheme. There are, indeed, two possible approaches: either the diffusive flux is approximated and, then, some form of centered or upwind approximation of the convection term is considered in the discretization of the divergence equation, or the total flux, which includes both diffusive and convective terms, is approximated, which leads to a centered-type approximation of the convection terms. The first approach is, perhaps, more popular in the Finite Difference and Finite Volume practitioner community, cf. [26, 32], while the second approach seems to be more popular in the Finite Element practitioner community. Nevertheless, it is worth mentioning that both approaches have been considered in the framework of Mixed Finite Element methods, see [29, 30, 39].

In the Mimetic Finite Difference setting, a numerical discretization of the full diffusion and convection flux has been proposed by [24]. A proper reformulation of the mimetic scheme as a conforming method, using the finite dimensional subspace of  $H(\text{div}, \Omega)$  given by the lifting of the degrees of freedom of the vector variable, makes it possible to perform the convergence analysis in a very similar way to that presented in [30].

From this overview, we can conclude that several numerical discretizations of the convection-diffusion equations that may fit in the HMM setting have been proposed in the literature. However, no systematic study has been carried out so far on the possible ways, and related advantages and drawbacks, in which a convective term can be treated numerically by using the more general HMM formulation. It is our main goal in this work to perform such an investigation in order to assess the behavior of such methods both theoretically and numerically.

The plan of the paper is as follows. In Section 2, we recall the principles of the HMM schemes for the pure diffusion equation, and we discuss how to discretize the convection term, using some centered, upwind or exponential fitting-like choice in accordance with a 2-point Finite Volume flux formula (or, from the point of view of Finite Elements, see [28, 38]). We also show that the numerical approximation proposed in [24], possibly with a stabilization term, is an HMM method, to which the theoretical analysis of the present paper apply. In Section 3 we provide full proofs of convergence under very general regularity conditions when the mesh size tends to zero and derive error estimates in suitable mesh-dependent norms when the scalar solution is in  $H^2(\Omega)$ . Section 4 is devoted to present and discuss how various instances of the HMM discretizations perform when applied to a set of standard test cases for the convection-diffusion equations including the approximation of solutions with boundary and internal layers. Finally, conclusions are given in Section 5.

## 2 The HMM formulation for convection-diffusion problems

### 2.1 The mathematical model

Let us consider the steady convection diffusion equation:

$$-\operatorname{div}(\Lambda \nabla p) + \operatorname{div}(Vp) = f \quad \text{in } \Omega, \quad (2.1)$$

$$p = g^D \quad \text{on } \partial\Omega \quad (2.2)$$

under the hypotheses:

- (H1)  $\Omega$  is a bounded, open, polygonal subset of  $\mathbf{R}^d$  with  $d \geq 1$ ;
- (H2)  $\Lambda : \Omega \rightarrow M_d(\mathbf{R})$  is a bounded, measurable, symmetric and uniformly elliptic tensor;
- (H3)  $f \in L^2(\Omega)$ ;
- (H4)  $V \in C^1(\overline{\Omega})^d$  is such that  $\operatorname{div}(V) \geq 0$ .

Moreover, let us introduce the diffusive flux and the total flux:

$$F = -\Lambda \nabla p \quad \text{and} \quad \tilde{F} = F + Vp. \quad (2.3)$$

For simplicity, we will restrict the presentation of the methods and the theoretical analysis to the case of homogeneous Dirichlet boundary condition by setting  $g^D = 0$  in (2.2) and we will consider the non-homogeneous case in the numerical experiments of Section 4.

Under Assumptions (H1)-(H4), the existence and uniqueness of a weak solution in  $H_0^1(\Omega)$  to (2.1)-(2.2) with  $g^D = 0$  is completely standard since the bilinear form associated with this problem is continuous and coercive.

**Remark 2.1** *The  $C^1$  regularity assumption on  $V$  in (H4) can be weakened for the convergence study (see Section 3.1.3). We assume the smoothness of the convection field in order to simplify a little bit some (already lengthy) technical arguments, and also to prove error estimates.*

### 2.2 Mesh notation and regularity

Let us begin with the definition of an admissible discretization of  $\Omega$  and the related notation.

**Definition 2.2 [Admissible discretization]** *An admissible discretization of  $\Omega$  is given by the triplet  $\mathcal{D}_h = (\Omega_h, \mathcal{E}_h, \mathcal{P}_h)$ , where the mesh size  $h$  will be defined in the following and where:*

- $\Omega_h$  is a finite family of non-empty open polygonal disjoint subsets  $E$  of  $\Omega$ , the cells of the mesh, such that  $\overline{\Omega} = \cup_{E \in \Omega_h} \overline{E}$ ;
- $\mathcal{E}_h$  is a finite family of non-empty open disjoint subsets  $e$  of  $\overline{\Omega}$ , the faces of the mesh, such that for all  $e \in \mathcal{E}_h$  there exists an affine hyperplane  $\mathcal{A}$  of  $\mathbf{R}^d$  and a cell  $E \in \Omega_h$  such that  $e \subset (\overline{E} \setminus E) \cap \mathcal{A}$ . We also assume that:
  - for all  $E \in \Omega_h$  there exists a subset  $\partial E$  of  $\mathcal{E}_h$  such that  $\overline{E} \setminus E = \cup_{e \in \partial E} \overline{e}$ ;
  - for all  $e \in \mathcal{E}_h$  either we have that  $e \subset \partial\Omega$  or we have that  $e \subset \overline{E} \cap \overline{E'}$  for some pair of elements  $E, E' \in \Omega_h$  with  $E \neq E'$ ;
- $\mathcal{P}_h$  is a family of points of  $\Omega$  indexed by  $E$ , i.e.,  $\mathcal{P}_h = (x_E)_{E \in \Omega_h}$ , and such that each mesh cell  $E$  is star-shaped with respect to  $x_E$ .

**Remark 2.3** *When all the mesh cells are convex shaped, a convenient choice for the points  $(x_E)_{E \in \Omega_h}$  is given, for instance, by the centers of gravity of the cells.*

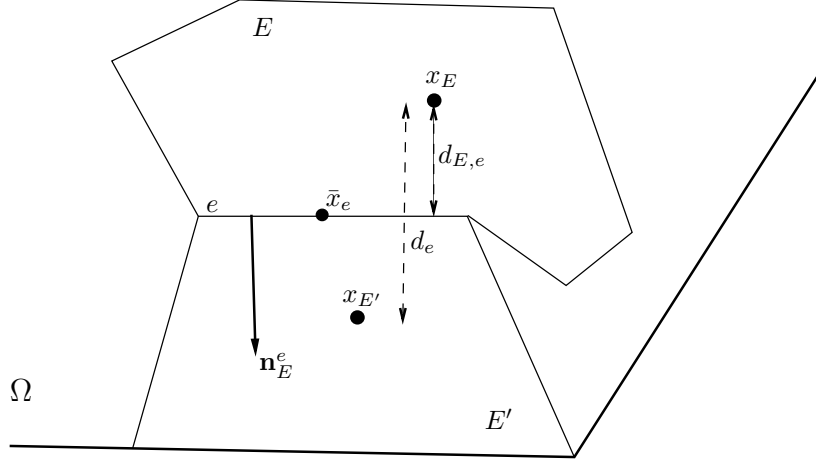


Figure 1: Mesh notations.

The  $d$ -dimensional measure of each cell  $E$  is denoted by  $|E|$  and the cell size by  $h_E$ . As usual, the mesh size is given by  $h = \sup_{E \in \Omega_h} h_E$ . For notation's consistency,  $|e|$  and  $h_e$  denote the  $(d-1)$ -dimensional measure of face  $e$  and the face diameter. For each face  $e \in \mathcal{E}_h$ ,  $\bar{x}_e$  denotes the barycenter of  $e$  and  $\mathbf{n}_E^e$  its normal direction pointing out of  $E$ . Moreover, to each face  $e$  we associate the unit normal vector  $\mathbf{n}^e$ , whose orientation is arbitrarily chosen when  $e$  is an internal face, and assumed pointing out of  $\Omega$  when  $e$  is a boundary face. We denote the set of the internal faces by  $\mathcal{E}_{h,\text{int}}$ , i.e.,  $\mathcal{E}_{h,\text{int}} = \{e \in \mathcal{E}_h \text{ for } e \not\subset \partial\Omega\}$ , and the set of the boundary faces by  $\mathcal{E}_{h,\text{ext}}$ , i.e.,  $\mathcal{E}_{h,\text{ext}} = \{e \in \mathcal{E}_h \text{ for } e \subset \partial\Omega\}$ . We will find it convenient to denote the two cells that share an internal face  $e$  by  $E$  and  $E'$ , and, where required, to fix the orientation of  $e$  so that  $\mathbf{n}_E^e \cdot \mathbf{n}^e = 1$ . Finally, we introduce the following geometric quantities that will be useful in the definition of the numerical convection flux in subsection 2.4.1:

$$d_{E,e} = \text{distance between } x_E \text{ and the hyperplane containing } e,$$

and

$$d_e = \begin{cases} d_{E,e} + d_{E',e} & \text{for any internal face } e \in \mathcal{E}_{h,\text{int}}, \\ d_{E,e} & \text{for any boundary face } e \in \mathcal{E}_{h,\text{ext}}. \end{cases}$$

Figure 1 illustrates some of these notations.

The proof of convergence for  $h \rightarrow 0$  that we present in Section 3 requires the following very mild geometrical assumptions on the meshes of  $\mathcal{D}_h$ .

(G1) Every mesh cell  $E$  is star-shaped with respect to the corresponding point  $x_E$ .

(G2) For any internal face  $e \in \mathcal{E}_{h,\text{int}}$ , let us introduce  $\mathcal{M}_e = \{E, E'\}$ , i.e., the cells on the opposite side of  $e$ ; then, the quantity

$$\text{regul}(\mathcal{D}_h) = \max \left( \max_{e \in \mathcal{E}_{h,\text{int}}, (E,E') \in \mathcal{M}_e} \frac{d_{E,e}}{d_{E',e}}, \max_{E \in \Omega_h, e \in \partial E} \frac{h_E}{d_{E,e}}, \max_{E \in \Omega_h} \text{Card}(\partial E) \right),$$

which expresses the mesh regularity, is *uniformly* bounded from above for  $h \rightarrow 0$ .

In the mimetic framework it is often used a similar condition, that we state as follows.

(ME) [*Star-shaped elements*] There exists a positive number  $\tau^*$  such that each element  $E$  is star-shaped with respect to *all* the points of a ball of radius  $\tau^* h_E$  centered at  $x_E$ .

Stronger conditions on the mesh regularity are required to derive an error estimate for the HMM approximations to the exact solution and flux. We formulate these mesh regularity conditions for  $d = 3$ ; the restriction to other dimensions is straightforward.

(HG)[*Shape-regularity*] There exist two positive real numbers  $N_s$  and  $\rho_s$  independent of  $h$  such that every mesh  $\Omega_h$  of the sequence admits a sub-partition into tetrahedrons  $\mathcal{S}_h$  such that:

(HG1) the decomposition of every polyhedron  $E \in \Omega_h$  denoted by  $\mathcal{S}_h|_E$  is formed by at most  $N_s$  tetrahedrons, and each vertex of  $\Omega_h$  is a vertex of  $\mathcal{S}_h$ ;

(HG2) every tetrahedron of  $\mathcal{S}_h$  is *shape-regular* in the sense that the ratio between  $r_T$ , the radius of its inscribed sphere, and  $h_T$ , its diameter, is bounded from below by  $\rho_s$ ; formally, we have that

$$\forall T \in \mathcal{S}_h : \frac{r_T}{h_T} \geq \rho_s > 0.$$

From the above assumptions several properties of the mesh, which are useful in the error analysis of the mimetic formulation, can be derived. For the sake of the reader's convenience, we list them below for future reference in the paper.

(M1) There exist two positive integers  $N_E$  and  $N_e$  that are independent of  $h$ ,  $E \in \Omega_h$  and  $e \in \mathcal{E}_h$  and such that every element  $E$  has  $\text{Card}(\partial E) \leq N_E$  faces, and every face  $e$  has  $\text{Card}(\partial e) \leq N_e$  edges.

(M2) For any mesh element  $E \in \Omega_h$ , the quantities  $|E|$ ,  $|e|$  for  $e \in \partial E$ , and  $|l|$  for each edge  $l \in \partial e$  properly scale with respect to  $h_E$ ; in particular, there exists a positive constant  $a^*$  such that

$$a^* h_E^{d-1} \leq |e|, \quad a^* h_E \leq h_e, \quad a^* h_E^{d-2} \leq |l|.$$

(M3) There exists a constant  $C^{\text{Ag}}$  independent of  $h_E$  and such that [20]:

$$\sum_{e \in \partial E} \|\phi\|_{L^2(e)}^2 \leq C^{\text{Ag}} \left( h_E^{-1} \|\phi\|_{L^2(E)}^2 + h_E |\phi|_{H^1(E)}^2 \right) \quad (2.4)$$

for any function  $\phi \in H^1(E)$ . We will refer to (2.4) as *the Agmon inequality*.

(M4) For any function  $q \in H^2(E)$ , there exists a *linear polynomial*  $\mathcal{L}_1(q)$  interpolating  $q$  and a constant  $C$ , independent of  $h_E$ , such that [17]:

$$\|q - \mathcal{L}_1(q)\|_{L^2(E)} + h_E |q - \mathcal{L}_1(q)|_{H^1(E)} \leq C h_E^2 |q|_{H^2(E)}. \quad (2.5)$$

### 2.3 Discretization of the diffusion term

To approximate (2.1)-(2.2), we introduce the space of the *discrete scalar fields*  $Q_h$  and the space of the *discrete flux fields*  $X_h$ . The discrete scalars  $q \in Q_h$  are defined by taking one degree of freedom per cell denoted by  $q_E$ , i.e.,  $q = (q_E)_{E \in \Omega_h}$ . Therefore, the space  $Q_h$  can be identified with the space of the piecewise constant polynomials defined on  $\Omega_h$ . Similarly, the *discrete fluxes* are defined by taking one degree of freedom per face per element denoted by  $F_E^e$ , i.e.,  $F = (F_E^e)_{E \in \Omega_h}^{e \in \partial E}$ , which represents the normal flux across the face  $e$  along the direction  $\mathbf{n}_E^e$ . We require that every flux  $F \in X_h$  satisfies the *flux conservation property* at any internal face:

$$\forall e \in \mathcal{E}_{h,\text{int}}, e \subset \partial E \cup \partial E' : F_E^e + F_{E'}^e = 0, \quad (2.6)$$

so that the elements of  $X_h$  only possess one degree-of-freedom per face, and the sign of  $F_E^e$  depends on the orientation of the face  $e$  with respect to  $E$ . The restriction of  $F$  to cell  $E \in \Omega_h$  is denoted by  $F_E = (F_E^e)_{e \in \partial E}$  and represents the collection of the normal fluxes along the directions  $\mathbf{n}_E^e$  for  $e \in \partial E$ . The set of these vector fields forms the linear space  $X_E$ . Throughout the paper we will also make use of

the symbol  $\widehat{X}_h$  to denote the linear space of the *discontinuous fluxes*, i.e., of the vectors having the same form  $F = (F_E^e)_{E \in \Omega_h}^{e \in \partial E}$  but that do not satisfy condition (2.6). Note that  $X_h$  is a linear sub-space of  $\widehat{X}_h$ .

The next ingredient of the HMM formulation is the *discrete divergence operator*  $\text{div}_h : \widehat{X}_h \rightarrow Q_h$ , which is defined as

$$\forall G \in \widehat{X}_h, \forall E \in \Omega_h : \quad (\text{div}_h(G))_E = \frac{1}{|E|} \sum_{e \in \partial E} |e| G_E^e. \quad (2.7)$$

To any sufficiently regular vector field  $G$  and scalar field  $q$ , we associate the interpolated fields  $G^I \in X_h$  and  $q^I \in Q_h$  that are given by

$$\forall e \in \mathcal{E}_h : (G^I)^e = \frac{1}{|e|} \int_e G \cdot \mathbf{n}^e \quad \text{and} \quad \forall E \in \Omega_h : (q^I)_E = \frac{1}{|E|} \int_E q. \quad (2.8)$$

**Remark 2.4** *The definition of the discrete divergence operator in (2.7) is consistent with the Gauss divergence theorem for the interpolations of (2.8), so that the following commutation property holds*

$$(\text{div}(G))^I = \text{div}_h(G^I). \quad (2.9)$$

We endow  $Q_h$  with the usual  $L^2(\Omega)$  scalar product for piecewise constant functions, i.e.,  $[\cdot, \cdot]_{Q_h} := [\cdot, \cdot]_{L^2}$ . On the other hand,  $X_h$  and  $\widehat{X}_h$  are equipped with the scalar product

$$[F, G]_{\widehat{X}_h} = \sum_{E \in \Omega_h} [F_E, G_E]_E, \quad (2.10)$$

that assembles the locally defined scalar products  $[\cdot, \cdot]_E$ . The local scalar products  $([\cdot, \cdot]_E)$  satisfy the coercivity and consistency assumptions:

(S1) there exist two positive constants  $\sigma_*$  and  $\sigma^*$  independent of the mesh size  $h$  such that for every mesh cell  $E$

$$\sigma_* |E| \sum_{e \in \partial E} (G_E^e)^2 \leq [G, G]_E \leq \sigma^* |E| \sum_{e \in \partial E} (G_E^e)^2 \quad \forall G \in X_h ;$$

(S2) for every element  $E$ , we have that

$$[(\Lambda_E \nabla q^1)^I, G]_E = - [\text{div}_h(G), q^1]_{L^2(E)} + \sum_{e \in \partial E} G_E^e \int_e q^1$$

for all  $G \in X_h$  and all linear polynomials  $q^1$ , and where  $\Lambda_E$  is the cell average of  $\Lambda$ .

**Remark 2.5**  $\Lambda_E$  is, actually, an approximation of  $\Lambda|_E$ , the restriction of the diffusion tensor  $\Lambda$  to cell  $E$ . To prove the convergence of the numerical solution in subsection 3.1, we only require that the diffusion tensor  $\Lambda$  satisfy the regularity assumption (H2), while we need a stronger regularity condition to derive the error estimates of subsection 3.2. In this latter case, we will find it convenient to assume (H2) and also that  $\Lambda$  be locally Lipschitz continuous on  $\Omega_h$ , i.e., for all  $E \in \Omega_h$ , the components of  $\Lambda|_E$  are Lipschitz continuous functions on  $E$ . Consequently,  $\Lambda_E$  can be any constant approximation of  $\Lambda|_E$  such that the estimate

$$\max_{i,j=1,d} \sup_{x \in E} |(\Lambda_E)_{ij} - \Lambda_{ij}(x)| = \mathcal{O}(h)$$

holds.

The construction of a family of scalar products satisfying the above assumptions when  $x_E$  is the center of gravity of  $E$  can be found in [22]. Moreover, in this case it has been proved in [33] that (S1)-(S2) lead necessarily to the following form:

$$\forall (F_E, G_E) \in X_E : [F_E, G_E]_E = |E| \Lambda_E \mathbf{v}_E(F_E) \cdot \mathbf{v}_E(G_E) + T_E(G_E)^T \mathbb{B}_E T_E(F_E) \quad (2.11)$$

where

$$\mathbf{v}_E(F_E) = -\frac{1}{|E|}\Lambda_E^{-1} \sum_{e \in \partial E} |e|F_E^e(\bar{x}_e - x_E) \quad (2.12)$$

is a constant approximation of  $\nabla p$  on cell  $E$ ,  $T_E(F_E) = (T_{E,e}(F^e))_{e \in \partial E}$  is given by

$$T_{E,e}(F_E) = F_E^e + \Lambda_E \mathbf{v}_E(F_E) \cdot \mathbf{n}_E^e, \quad (2.13)$$

and  $\mathbb{B}_E$  is a symmetric positive definite matrix of size  $\text{Card}(\partial E)$ . More precisely, it turns out that the matrix  $\mathbb{B}_E$  satisfies the following coercivity condition, which is directly related to (S1).

(C) There exists a positive constant  $\alpha$ , which is independent of the mesh size, such that for all  $E \in \Omega_h$  and  $G_E \in X_E$  there holds that:

$$\alpha \sum_{e \in \partial E} |e|d_{E,e}(T_{E,e}(G_E))^2 \leq T_E(G_E)^T \mathbb{B}_E T_E(G_E) \leq \frac{1}{\alpha} \sum_{e \in \partial E} |e|d_{E,e}(T_{E,e}(G_E))^2.$$

If  $x_E$  is not the barycenter of  $E$ , the same construction (2.11)-(2.13) still holds provided that (S2) be modified by introducing a suitable integration weight, see [33].

The HMM discretization to problem (2.1)-(2.2) with  $V = 0$ , which provides us the desired approximation of the diffusion operator, takes the form:

*find*  $(p_h, F_h) \in Q_h \times X_h$  *such that:*

$$\forall G \in X_h : \quad [F_h, G]_{\widehat{X}_h} = [\text{div}_h(G), p_h]_{Q_h} \quad (2.14)$$

$$\forall q \in Q_h : \quad [\text{div}_h(F_h), q]_{Q_h} = [f, q]_{Q_h}. \quad (2.15)$$

$p_h \in Q_h$  and  $F_h \in X_h$  are the approximations to  $p^J$  and  $F^I$ , the interpolations of the exact scalar solution  $p$  and its flux  $F = -\Lambda \nabla p$ .

The HMM method can be easily hybridized through the introduction of  $H(\mathcal{E}_h)$ , the space of face values  $q_{\mathcal{E}_h} = (q_e)_{e \in \mathcal{E}_h} \in \mathbf{R}^{\text{Card}(\mathcal{E}_h)}$  with  $q_e = 0$  for  $e \in \mathcal{E}_{h,\text{ext}}$  and imposing explicitly the flux conservation property (2.6). The discrete variational form (2.14)-(2.15) with  $[\cdot, \cdot]_E$  satisfying (2.11)-(2.13) is equivalent to:

*find*  $(p_h, F_h, p_{\mathcal{E}_h}) \in Q_h \times \widehat{X}_h \times H(\mathcal{E}_h)$  *such that:*

$$\forall E \in \Omega_h, \forall G_E \in X_E : \quad [F_E, G_E]_E = \sum_{e \in \partial E} |e|G_E^e(p_E - p_e), \quad (2.16)$$

$$\forall E \in \Omega_h : \quad \sum_{e \in \partial E} |e|F_E^e = \int_E f, \quad (2.17)$$

$$\forall e \in \mathcal{E}_{h,\text{int}}, : \quad F_E^e + F_{E'}^e = 0, \quad (2.18)$$

where  $E, E' \in \Omega_h$  are the two elements such that  $e \subset \partial E \cap \partial E'$  for every  $e \in \mathcal{E}_{h,\text{int}}$ . Under suitable assumptions on the regularity of the exact solution  $p$ , the additional unknowns  $p_{\mathcal{E}_h} = (p_e)_{e \in \mathcal{E}_h}$  approximate the face average of the exact solution over each mesh face. We will formalize this concept through the introduction of  $p^J \in H(\mathcal{E}_h)$ , the face interpolation of  $p$ , in Section 3.2, see equation (3.43).

## 2.4 Discretization of the convective term

As discussed in the introduction, two different strategies can be considered for the numerical treatment of the convection term in the HMM discretization of an elliptic problem. In the first strategy, which is reviewed in subsection 2.4.1, we introduce some form of centered or upwind approximation of the convection term in the discretization of the divergence equation provided by the HMM method, cf. [26, 32].



In the second strategy, which is reviewed in subsection 2.4.2, the total flux, which includes both diffusive and convective terms, is approximated, thus leading to a centered-type approximation of the convection terms, cf. [24]. Both approaches have been considered for the Mixed Finite Element method in [29, 30, 39]. It turns out that in the new framework of HMM methods a unified formulation is possible, which is the topic of subsection 2.4.3. We end this section with a discussion on an alternative hybridized form of the numerical convection terms, cf. subsection 2.4.4.

In the rest of this section, we assume that the velocity field  $V$  is a continuous function with continuous derivative, i.e.,  $V \in C^1(\overline{\Omega})^d$ . The cell restriction of its interpolation in  $X_h$  is given by the set of real numbers  $(V_E^e)_{e \in \partial E} \in X_E$  such that

$$\forall e \in \partial E : \quad V_E^e = \frac{1}{|e|} \int_e V \cdot \mathbf{n}_E^e. \quad (2.19)$$

### 2.4.1 FV-based discretizations

Several discretization schemes for the convection term are available in the Finite Volume literature, e.g., the second-order centered scheme, the first-order upwind scheme, the  $\theta$ -scheme, the Scharfetter-Gummel scheme, etc. In these methods, the convection flux of the exact solution field  $p$  is approximated through the numerical convection flux of the discrete scalar field  $p_h \in Q_h$ ; this numerical convection flux is given by the collection of real numbers  $F_c(p_h) = (F_c(p_h))_E^e$  such that

$$\forall E \in \Omega_h, \forall e \in \partial E : \quad \frac{1}{|e|} \int_e V p \cdot \mathbf{n}_E^e \approx (F_c(p_h))_E^e. \quad (2.20)$$

We list below the schemes that we will explicitly consider in the section of numerical experiments. We let  $E'$  be the cell on the other side of  $e$  if  $e \in \mathcal{E}_{h,\text{int}}$  and assume for notation's simplicity that  $p_{E'} = 0$  if  $e \in \mathcal{E}_{h,\text{ext}}$ .

- The *second-order centered scheme* is given by the approximation

$$\frac{1}{|e|} \int_e V p \cdot \mathbf{n}_E^e \approx (F_c(p_h))_E^e = V_E^e \frac{p_E + p_{E'}}{2}.$$

- The *first-order upwind scheme* is given by the approximation

$$\frac{1}{|e|} \int_e V p \cdot \mathbf{n}_E^e \approx (F_c(p_h))_E^e = (V_E^e)^+ p_E - (V_E^e)^- p_{E'}$$

with  $s^\pm = \max(\pm s, 0)$ .

- The  *$\theta$ -scheme* is given by the approximation

$$\begin{aligned} \frac{1}{|e|} \int_e V p \cdot \mathbf{n}_E^e &\approx (F_c(p_h))_E^e = (V_E^e)^+ ((1 - \theta)p_E + \theta p_{E'}) - (V_E^e)^- ((1 - \theta)p_{E'} + \theta p_E) \\ &= (1 - 2\theta) ((V_E^e)^+ p_E - (V_E^e)^- p_{E'}) + \theta V_E^e (p_E + p_{E'}) \end{aligned}$$

with  $\theta \in [0, 1/2]$ ; this choice is clearly intermediate between the centered and the upwind schemes.

- The *Scharfetter-Gummel scheme* [46] is given by the approximation

$$\frac{1}{|e|} \int_e V p \cdot \mathbf{n}_E^e \approx (F_c(p_h))_E^e = \frac{1}{d_e} (A_{\text{sg}}(d_e V_E^e) p_E - A_{\text{sg}}(-d_e V_E^e) p_{E'}), \quad (2.21)$$

with

$$A_{\text{sg}}(s) = \frac{-s}{e^{-s} - 1} - 1. \quad (2.22)$$

Note that the first three approaches above can be found also in the Finite Element literature, see for instance [28, 38]. As pointed out in [26], the Scharfetter-Gummel scheme in [46] was written for an isotropic homogeneous material, i.e.,  $\Lambda = I$ . In the original formulation, diffusion and convection terms were simultaneously treated to define the numerical flux. Removing the diffusive part in the numerical flux formulation allows us to obtain the formulas (2.21)-(2.22). This definition of a pure convective flux through the simple elimination of the diffusive part is somewhat basic in the general case  $\Lambda \neq I$ , especially if some eigenvalues of  $\Lambda$  are small. Although the above definition of  $A_{\text{sg}}$  ensures the  $L^2$ -stability of the scheme, it can give quite bad solutions in convection-dominated cases. This fact can be understood if one comes back to the 2-points Finite Volume scheme for  $-\epsilon \Delta p + \text{div}(Vp) = f$ : choice (2.22) ensures the maximum principle of the scheme only if  $\epsilon \geq 1$ , while the maximum principle is lost numerically if  $\epsilon < 1$ . When applying the Scharfetter-Gummel method to compute the numerical convective flux, a better choice is provided by locally scaling  $A_{\text{sg}}$  in accordance with the smallest eigenvalue of  $\Lambda$ . If  $e$  is the face between  $E$  and  $E'$  and  $\lambda_e$  is the smallest eigenvalue of  $\Lambda_E$  and  $\Lambda_{E'}$ , we use

$$A_{\text{sg},\Lambda,\epsilon}(s) = \min(1, \lambda_e) A_{\text{sg}} \left( \frac{s}{\min(1, \lambda_e)} \right) \quad (2.23)$$

instead of  $A_{\text{sg}}(s)$  in (2.21). In this way, the numerical flux automatically and locally adjusts the upwinding of the convection term depending on its strength with respect to the diffusive term without perturbing the consistency property of  $A_{\text{sg}}$ . Note that  $\lambda_e \rightarrow 0$  implies that  $\lambda_e A_{\text{sg}}(s/\lambda_e) \rightarrow s^+$ . Therefore, if the local diffusion is very small, this implementation of the Scharfetter-Gummel method allows the flux to adjust to upwinding automatically, thus bringing enough numerical diffusion to ensure a better stability.

Once an FV-based discretization of the convective term has been chosen, the divergence of the convection term in (2.1), i.e.,  $\text{div}(Vp)$ , is approximated on  $E$  by

$$(\text{div}(Vp))_E^I \approx \frac{1}{|E|} \sum_{e \in \partial E} |e| (F_c(p_h))_E^e = \text{div}_h(F_c(p_h))|_E$$

and the HMM approximation to the model problem (2.1)-(2.2) then reads:

*find*  $(p_h, F_h) \in Q_h \times X_h$  *such that*

$$\forall G \in X_h : \quad [F_h, G]_{\hat{X}_h} = [\text{div}_h(G), p_h]_{Q_h}, \quad (2.24)$$

$$\forall q \in Q_h : \quad [\text{div}_h(F_h + F_c(p_h)), q]_{Q_h} = [f^I, q]_{Q_h}. \quad (2.25)$$

## 2.4.2 MFD-based discretizations

From the theoretical standpoint, Mimetic Finite Differences have only very recently approached problems different than the pure diffusion one (see for instance [7, 9, 10]). To our knowledge, the only paper considering development and error analysis of convection-diffusion equations directly in the framework of MFD is found in [24]. In this subsection, we briefly review the formulation and the major convergence results of the method considered in that paper, and we show how it can be reformulated as an HMM method.

Let  $H(\text{div}, \Omega)$  be the space of vector fields all of whose components are square integrable functions and that have square integrable divergence. Formally,

$$H(\text{div}, \Omega) = \{ \mathbf{v} \in (L^2(\Omega))^d \text{ such that } \text{div}(\mathbf{v}) \in L^2(\Omega) \}$$

is a Hilbert space when equipped with the scalar product

$$[\mathbf{v}, \mathbf{u}]_{H(\text{div}, \Omega)} = \int_{\Omega} \mathbf{v} \cdot \mathbf{u} + \int_{\Omega} \text{div}(\mathbf{v}) \text{div}(\mathbf{u})$$

and the corresponding norm

$$\|\mathbf{v}\|_{H(\operatorname{div},\Omega)}^2 = \|\mathbf{v}\|_{L^2(\Omega)}^2 + \|\operatorname{div}(\mathbf{v})\|_{L^2(\Omega)}^2.$$

In [24], it is considered a numerical approximation to the *mixed variational formulation* of problem (2.1)-(2.2), which reads as [19]:

find  $(\tilde{F}, p) \in H(\operatorname{div}, \Omega) \times L^2(\Omega)$  such that

$$\forall \mathbf{v} \in H(\operatorname{div}, \Omega) : \quad \left[ \Lambda^{-1} \tilde{F}, \mathbf{v} \right]_{L^2} - [p, \operatorname{div}(\mathbf{v})]_{L^2} - [\Lambda^{-1} V p, \mathbf{v}]_{L^2} = 0 \quad (2.26)$$

$$\forall q \in L^2(\Omega) : \quad \left[ \operatorname{div}(\tilde{F}), q \right]_{L^2} = [f, q]_{L^2}, \quad (2.27)$$

where  $\tilde{F}$  is the total vector flux defined in (2.3).

To discretize the convection term, we transform the corresponding variational term as follows:

$$\forall \mathbf{v} \in H(\operatorname{div}, \Omega) : \quad [\Lambda^{-1} V p, \mathbf{v}]_{L^2} \approx \sum_{E \in \Omega_h} \int_E \Lambda_E^{-1} V p \cdot \mathbf{v} \quad \rightarrow \quad \forall G \in X_h : \quad \sum_{E \in \Omega_h} p_E [V^I, G]_E,$$

where the components of the interpolated velocity field  $V^I \in X_h$ , i.e.,  $(V^I)_E^e$  for all  $E \in \Omega_h$  and  $e \in \partial E$ , are given by (2.19), and the local scalar products are required to satisfy Assumptions (S1)-(S2). The mimetic variational formulation presented in [24] reads as:

find  $(\tilde{F}_h, p_h) \in X_h \times Q_h$  such that

$$\forall G \in X_h : \quad \left[ \tilde{F}_h, G \right]_{\hat{X}_h} - [p_h, \operatorname{div}_h(G)]_{Q_h} - \sum_{E \in \Omega_h} p_E [V^I, G]_E = 0, \quad (2.28)$$

$$\forall q \in Q_h : \quad \left[ \operatorname{div}_h(\tilde{F}_h), q \right]_{Q_h} = [f^I, q]_{Q_h}. \quad (2.29)$$

The convergence analysis of this scheme is carried out in [24] under assumptions on the grid regularity that are substantially equivalent to (HG)-(ME). When the scalar solution  $p$  is in  $H^2(\Omega)$ , the analysis provides the following error estimate

$$\|\tilde{F}_h - \tilde{F}^I\|_{\hat{X}_h} + \|p_h - p^I\|_{Q_h} \leq Ch \|p\|_{H^2(\Omega)} \quad (2.30)$$

where  $\|\cdot\|_{\hat{X}_h}$  and  $\|\cdot\|_{Q_h}$  are the norms induced by the inner products of the spaces  $\hat{X}_h$  and  $Q_h$ , respectively. It is worth mentioning that the approximation of the scalar variable is superconvergent when calculation is performed on a wide set of meshes. Superconvergence was also theoretically proved under some stronger assumptions on the regularity of the domain shape, the source term, and the velocity field.

Despite convergence is proved for  $h \rightarrow 0$ , this scheme is expected to become unstable when the model problem is dominated by convection. This fact usually manifests through spurious effects like numerical undershoots, overshoots, or oscillations that may appear in the approximate solution. To improve stability, we modify the divergence equation by introducing a stabilization term that depends on the solution's jumps at mesh faces. We use the symbols  $E$  and  $E'$  to denote the two distinct cells that share face  $e$  when  $e$  is internal, and assume the orientation of  $e$  such that  $\mathbf{n}_E^e \cdot \mathbf{n}^e = 1$ . Let us now introduce the *jump* of the discrete scalar field  $q_h \in Q_h$ , which is given by:

$$[[q_h]]_e = \begin{cases} q_E - q_{E'} & \text{for } e \in \mathcal{E}_{h,\text{int}}, \\ q_E & \text{for } e \in \mathcal{E}_{h,\text{ext}}. \end{cases} \quad (2.31)$$

Equation (2.29) is substituted by

$$\forall q \in Q_h : \quad \left[ \operatorname{div}_h(\tilde{F}_h) + J_h(p_h), q \right]_{Q_h} = [f^I, q]_{Q_h}, \quad (2.32)$$

where the stabilization term  $J_h(p_h)$  is given by:

$$J_h(p_h)|_E = \frac{\alpha}{2|E|} \sum_{e \in \partial E} |e| |(V^I)_E^e| \llbracket p_h \rrbracket_e, \quad (2.33)$$

and  $\alpha$  is a non-negative parameter that can be tuned to control the amount of numerical dissipation of the scheme.

This approach formally differs from the method introduced in the previous subsection by using FV-based discretizations in that the convection term is numerically treated as part of the mimetic flux equation. However, it is possible to “extract” an explicit form of the numerical convection flux from the scheme given by equations (2.28) and (2.32) to reformulate it as an HMM method. To this purpose, we define the collection of numbers  $F_h = (F_E^e)_{E \in \Omega_h, e \in \partial E}$  by

$$F_E^e = \tilde{F}_E^e - p_E (V^I)_E^e. \quad (2.34)$$

Equation (2.28) shows that  $F_h$  satisfies (2.14), and therefore plays the role of a purely diffusive flux. Moreover, noticing that the stabilization term  $J_h(p_h)$  is locally written as a balance of fluxes, i.e., a discrete divergence, allows us to identify the convective flux as

$$(F_c(p_h))_E^e = p_E (V^I)_E^e + \frac{\alpha}{2} |(V^I)_E^e| (p_E - p_{E'}) \quad (2.35)$$

(we let  $p_{E'} = 0$  if  $e$  is a boundary edge) so that (2.32) is simply given by  $\operatorname{div}_h(F_h + F_c(p_h)) = f^I$ . The stabilized MFD scheme (2.28) and (2.32) can, therefore, be written as:

find  $p_h \in Q_h$  and  $F_h \in \widehat{X}_h$  such that

$$\forall G \in X_h : \quad [F_h, G]_{\widehat{X}_h} = [\operatorname{div}_h(G), p_h]_{Q_h}, \quad (2.36)$$

$$\forall q \in Q_h : \quad [\operatorname{div}_h(F_h + F_c(p_h)), q]_{Q_h} = [f^I, q]_{Q_h}, \quad (2.37)$$

$$\forall e \in \mathcal{E}_{h,\text{int}} : \quad (F_h + (F_c(p_h))_E^e + (F_h + (F_c(p_h))_{E'}^e) = 0. \quad (2.38)$$

Note that the diffusive flux  $F_h$  and the convective flux  $F_c(p_h)$  are not conservative in the sense of (2.6) when considered separately, and, therefore, belong to the linear space  $\widehat{X}_h$ . However, their sum, i.e.,  $F_h + F_c(p_h)$ , is conservative since it belongs to  $X_h$  in view of equation (2.38).

### 2.4.3 Unified setting

A unified formulation exists for the numerical discretization of the convection term. This formulation includes the FV-based discretizations, as was noted in [26], and the MFD-based discretization (2.36)-(2.38). This fact makes it possible to simplify the software implementation and carry out a unified theoretical analysis.

Let us consider two functions  $A, B : \mathbf{R} \rightarrow \mathbf{R}$  and choose the numerical convection flux as the collection of real numbers

$$F_c(p_h) = (F_c(p_h)_E^e)_{E \in \Omega_h, e \in \partial E} \quad (2.39)$$

such that

$$\forall E \in \Omega_h, \forall e \in \partial E : (F_c(p_h))_E^e := \frac{1}{d_e} (A(d_e V_E^e) p_E + B(d_e V_E^e) p_{E'}). \quad (2.40)$$

Since in the MFD discretization of the convection term these flux components are not conservative, the diffusive flux components cannot be conservative either and conservation must be imposed on the total flux. The generic HMM approximation to the model problem (2.1)-(2.2) is thus written as:

find  $p_h \in Q_h$  and  $F_h \in \widehat{X}_h$  such that

$$\forall G \in X_h : \quad [F_h, G]_{\widehat{X}_h} = [\operatorname{div}_h(G), p_h]_{Q_h}, \quad (2.41)$$

$$\forall q \in Q_h : \quad [\operatorname{div}_h(F_h + F_c(p_h)), q]_{Q_h} = [f^I, q]_{Q_h}, \quad (2.42)$$

$$\forall e \in \mathcal{E}_{h,\text{int}} : \quad (F_h + (F_c(p_h))_E^e + (F_h + (F_c(p_h))_{E'}^e = 0. \quad (2.43)$$

The schemes presented in the previous subsections can all be included in this general setting, with the following choices of  $A$  and  $B$ :

- *Centered scheme*:  $A(s) = A_{\text{ce}}(s) := \frac{s}{2}$  and  $B(s) = -A_{\text{ce}}(-s) = \frac{s}{2}$ .
- *Upwind scheme*:  $A(s) = A_{\text{up}}(s) := s^+$  and  $B(s) = -A_{\text{up}}(-s) = -s^-$ .
- $\theta$ -*scheme*:  $A(s) = A_\theta(s) := (1 - 2\theta)A_{\text{up}}(s) + 2\theta A_{\text{ce}}(s)$  and  $B(s) = -A_\theta(-s)$ .
- *Scharfetter-Gummel scheme*:  $A(s) = A_{\text{sg}}(s)$  defined by (2.22) and  $B(s) = -A_{\text{sg}}(-s)$ ; the locally scaled Scharfetter-Gummel scheme is obtained by using  $A_{\text{sg},\Lambda,e}$  defined by (2.23) instead of  $A_{\text{sg}}$ .
- *Stabilized MFD scheme*:  $A(s) = s + \frac{\alpha}{2}|s|$  and  $B(s) = -\frac{\alpha}{2}|s|$ .

The first four choices in (2.40) lead to a conservative definition of the numerical convection flux, whereas the last one does not. However, in all the cases mentioned above, total conservation is ensured by (2.43). We notice that all these choices of  $A$  and  $B$  satisfy the following properties:

(AB1)  $A : \mathbf{R} \rightarrow \mathbf{R}$  and  $B : \mathbf{R} \rightarrow \mathbf{R}$  are Lipschitz-continuous functions and  $A(0) = B(0) = 0$ ;

(AB2)  $A(s) + B(s) = s$  for any real number  $s$ ;

(AB3) one of the two following alternatives holds:

(AB3-s)  $A(s) + B(-s) = 0$  and  $A(s) - B(s) \geq 0$  for any real number  $s$ ;

(AB3-w) the function  $s \rightarrow A(s) + B(-s)$  is odd and there exists  $C > 0$  such that  $A(s) - B(s) \geq -C|s|$  for any real number  $s$ .

We refer to (AB3-s) as the *strong* (AB3) condition, and to (AB3-w) as the *weak* (AB3) condition. Assumption (AB3-s) is satisfied by all the FV-based discretizations listed above whereas the MFD-based discretization satisfies (AB3-w). In fact, condition  $A(s) + B(-s) = 0$  in (AB3-s) is the one ensuring the conservation of the numerical convection flux (2.40). On the other hand, the numerical convection flux extracted from the MFD-based formulation satisfies (AB3-w) and, hence, is not conservative. We will see in Section 3 that Assumptions (AB1)-(AB3) are enough to carry out the theoretical analysis of the scheme in (2.39)-(2.43), with slightly different results depending on which alternative in (AB3) is satisfied.

**Remark 2.6** *It is worth noting that equation (2.42) can be rewritten in a Finite Volume form as the following cell-based flux balance equation:*

$$\forall E \in \Omega_h : \quad \sum_{e \in \partial E} |e| \left( F_E^e + (F_c(p_h))_E^e \right) = \int_E f. \quad (2.44)$$

**Remark 2.7** *We could also choose, in (2.40), different functions  $A = A^e$  and  $B = B^e$  for each edge  $e$ , provided that all these functions satisfy (AB1)-(AB3) and that their Lipschitz constants remains uniformly bounded as the mesh size tends to 0. This setting would allow the scheme to make a finer tuning of the numerical diffusion due to upwinding, thus better adapting the scheme behavior to the location inside the domain or the local geometry of the mesh.*

#### 2.4.4 An alternative hybrid discretization of the convection term

An alternative discretization of the convection term is possible by using the hybridized value  $p_e$  in (2.40) instead of  $p_{E'}$ , an idea introduced in [4]. In such a case, we define the numerical convection flux of the discrete scalar field described by  $(p_h, p_{\mathcal{E}_h}) \in Q_h \times H(\mathcal{E}_h)$  as the collection of real numbers:

$$F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}) = \left( (F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}))_E^e \right)_{E \in \Omega_h, e \in \partial E} \quad (2.45)$$

such that

$$\forall E \in \Omega_h, \forall e \in \partial E : (F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}))_E^e = \frac{1}{d_e} (A(d_e V_E^e) p_E + B(d_e V_E^e) p_e). \quad (2.46)$$

The substantial difference with the preceding choice (2.40) is that no property on  $A$  and  $B$  ensure that the fluxes  $F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h})$  are conservative (and they are not in general). However, this will not bring any additional difficulty in the theoretical study provided that the following weaker form of (AB3) is considered:

(AB3h) one of the following strong or weak alternatives holds:

(AB3h-s)  $A(s) - B(s) \geq 0$  for any real number  $s$ ,

(AB3h-w) there exists  $C > 0$  such that  $A(s) - B(s) \geq -C|s|$  for any real number  $s$ .

The hybrid HMM formulation can then be written as:

find  $(p_h, F_h, p_{\mathcal{E}_h}) \in Q_h \times \widehat{X}_h \times H(\mathcal{E}_h)$  such that

$$\forall E \in \Omega_h, \forall G_E \in X_E : [F_E, G_E]_E = \sum_{e \in \partial E} |e| G_E^e (p_E - p_e), \quad (2.47)$$

$$\forall E \in \Omega_h : \sum_{e \in \partial E} |e| \left( F_E^e + (F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}))_E^e \right) = \int_E f, \quad (2.48)$$

$$\forall e \in \mathcal{E}_{h,\text{int}} : (F_h + (F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}))_E^e) + (F_h + (F_{c,\mathcal{E}_h}(p_h, p_{\mathcal{E}_h}))_{E'}^e) = 0, \quad (2.49)$$

where the local scalar products used in (2.47) satisfy (S1)-(S2), and, thus, may be given in the form (2.11)-(2.13).

**Remark 2.8** *An important advantage of discretizing the convective fluxes by using (2.45)-(2.46) instead of (2.39)-(2.40) is that the unknowns  $p_h$  and  $F_h$  in the resulting numerical formulation (2.47)-(2.49) can be eliminated by static condensation, i.e., through a local Gaussian elimination (this classical technique is not directly applicable to (2.39)-(2.40)). This procedure, which is common for hybrid-Mixed Finite Elements, provides a reduced linear system in the face unknowns  $p_{\mathcal{E}_h}$ . Moreover, when the discretization of the convection term increases significantly the numerical diffusion, as for example in the case of the upwind scheme, the hybrid version of the HMM method is likely to be less diffusive than that provided by (2.39)-(2.40).*

### 3 Theoretical study

In the present section we develop the theoretical analysis for the class of methods that we wish to investigate in this work. In subsection 3.1, we prove the convergence of the numerical approximations to the exact solution and its gradient. The analysis is based on a compactness argument, which is common in the Finite Volume literature, under the weaker assumptions of mesh regularity (G1)-(G2) (see also Definition 2.2). In subsection 3.2, we prove an  $\mathcal{O}(h)$  convergence rate for the numerical approximation of both scalar solution and flux. The analysis is on a stability and consistency argument, which are in the MFD (and FEM) literature, under the stronger mesh regularity Assumptions (HG) and (ME).

Let us introduce the mesh-dependent norms for the spaces  $X_h$  and  $Q_h$ . Let  $\mathcal{D}_h$  be an admissible mesh in accordance with Definition 2.2 that satisfies (G1)-(G2) or, alternatively, (HG)-(ME). The scalar product in  $\widehat{X}_h$  induces the norm:

$$\|G\|_{\widehat{X}_h}^2 = [G, G]_{\widehat{X}_h} \quad \forall G \in \widehat{X}_h, \quad (3.1)$$

and its local counterpart

$$\|G\|_E^2 = [G_E, G_E]_E \quad \forall G_E \in X_E. \quad (3.2)$$

The elements of  $Q_h$  can be identified with the  $\Omega_h$ -piecewise constant functions and the scalar product in  $Q_h$  is, in fact, the  $L^2$ -scalar product for such functions. Therefore, it is quite natural to consider the  $L^2$  norm. However, we will also find it useful to carry out the analysis by using the discrete  $H_0^1$ -like norm

$$\|q_h\|_{1, \mathcal{D}_h} = \left( \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} \left( \frac{|q_E - q_{E'}|}{d_e} \right)^2 \right)^{1/2} \quad \forall q_h \in Q_h, \quad (3.3)$$

where  $E'$  is the cell on the other side of  $e \in \partial E \cap \mathcal{E}_{h,\text{int}}$  and, to ease notation, we take  $q_{E'} = 0$  if  $e \in \partial E \cap \mathcal{E}_{h,\text{ext}}$ . We will also need a discrete  $H^1$  norm on  $Q_h \times H(\mathcal{E}_h)$ :

$$\|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h} = \left( \sum_{E \in \Omega_h} \sum_{e \in \partial E} \frac{|e|}{d_{E,e}} |q_E - q_e|^2 \right)^{1/2} \quad \forall (q_h, q_{\mathcal{E}_h}) \in Q_h \times H(\mathcal{E}_h). \quad (3.4)$$

It is easy to see that this norm is stronger than (3.3). More precisely, if  $\theta \geq \text{regul}(\mathcal{D}_h)$  there exists a constant  $C$  only dependent on  $\theta$  such that, for all  $(q_h, q_{\mathcal{E}_h}) \in Q_h \times H(\mathcal{E}_h)$ , there holds that

$$\|q_h\|_{1, \mathcal{D}_h} \leq C \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}. \quad (3.5)$$

In the following developments, we will number all constants whose value may be zero depending on which alternative is considered in (AB3), i.e., the strong (AB3-s) or the weak (AB3-w) condition. We will also use the symbol  $\lesssim$  to indicate an upper bound that holds up to a positive multiplicative constant independent of  $h$ . However, we will trace explicitly the constants where required by the proofs or that may be zero depending on the choice of Assumption (AB3).

**Lemma 3.1** *Assume that (H1)-(H3) hold. Let  $\mathcal{D}_h$  be an admissible discretization of  $\Omega$  such that  $\theta \geq \text{regul}(\mathcal{D}_h)$ ; let  $F_c(q)$  be the convective flux of  $q \in Q_h$  given by (2.39)-(2.40) for the vector field  $V \in C^1(\overline{\Omega})^d$  with  $A$  and  $B$  satisfying Assumptions (AB1)-(AB3). Then, there exists a non-negative constant  $C_1 \geq 0$  that only depends on  $\theta, V, A, B$  such that*

$$\begin{aligned} \forall (q, q_{\mathcal{E}_h}) \in Q_h \times H(\mathcal{E}_h) : \\ \frac{1}{2} \int_{\Omega} q^2 \text{div}(V) \leq \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F_c(q))_E^e (q_E - q_e) + C_1 h \|(q, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2, \end{aligned} \quad (3.6)$$

and where  $C_1 = 0$  if (AB3-s) holds.

**Proof of Lemma 3.1** By gathering the sum by faces we transform the term involving  $F_c(q)$  in the right-hand side of (3.6) as follows:

$$\begin{aligned} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F_c(q))_E^e (q_E - q_e) &= \sum_{e \in \mathcal{E}_h} |e| ((F_c(q))_E^e (q_E - q_e) + (F_c(q))_{E'}^e (q_{E'} - q_e)) \\ &= \sum_{e \in \mathcal{E}_h} |e| (F_c(q))_E^e (q_E - q_{E'}) \\ &\quad + \sum_{e \in \mathcal{E}_h} |e| ((F_c(q))_E^e + (F_c(q))_{E'}^e) (q_{E'} - q_e). \end{aligned} \quad (3.7)$$

To handle the first term in the right-hand side of (3.7), we note that using (2.40) and writing, thanks to (AB2),

$$A(d_e V_E^e) = \frac{1}{2}(d_e V_E^e + A(d_e V_E^e) - B(d_e V_E^e)) \text{ and } B(d_e V_E^e) = \frac{1}{2}(d_e V_E^e + B(d_e V_E^e) - A(d_e V_E^e)), \quad (3.8)$$

we have

$$(F_c(q))_E^e = \frac{1}{2}V_E^e(q_E + q_{E'}) + \frac{1}{2d_e}(A(d_e V_E^e) - B(d_e V_E^e))(q_E - q_{E'}).$$

Therefore, we infer that

$$\begin{aligned} \sum_{e \in \mathcal{E}_h} |e|(F_c(q))_E^e(q_E - q_{E'}) &= \frac{1}{2} \sum_{e \in \mathcal{E}_h} |e|V_E^e(q_E + q_{E'})(q_E - q_{E'}) \\ &\quad + \frac{1}{2} \sum_{e \in \mathcal{E}_h} \frac{|e|}{d_e}(A(d_e V_E^e) - B(d_e V_E^e))(q_E - q_{E'})^2. \end{aligned} \quad (3.9)$$

Then, let us observe that

$$\sum_{e \in \mathcal{E}_h} |e|(q_E - q_{E'})^2 \lesssim \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e|(q_E - q_e)^2 \lesssim h\|(q, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2. \quad (3.10)$$

By using (AB3), the conservation of  $(V_E^e)_{E \in \Omega_h, e \in \partial E}$ , the fact that  $\sum_{e \in \partial E} |e|V_E^e = \int_E \operatorname{div}(V)$ , and inequality (3.10) we obtain the following estimate:

$$\begin{aligned} \sum_{e \in \mathcal{E}_h} |e|(F_c(q))_E^e(q_E - q_{E'}) &\geq \frac{1}{2} \sum_{e \in \mathcal{E}_h} |e|V_E^e(q_E^2 - q_{E'}^2) - C_2 \sum_{e \in \mathcal{E}_h} |e|(q_E - q_{E'})^2 \\ &\geq \frac{1}{2} \sum_{E \in \Omega_h} q_E^2 \sum_{e \in \partial E} |e|V_E^e - C_3 h\|(q, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \\ &\geq \frac{1}{2} \int_{\Omega} q^2 \operatorname{div}(V) - C_3 h\|(q, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \end{aligned} \quad (3.11)$$

where  $C_2$  and  $C_3$  only depend on  $\theta, V, A, B$ , and  $C_2 = C_3 = 0$  if (AB3-s) holds.

From (2.40) and since  $V_E^e = -V_{E'}^e$ , we have

$$(F_c(q))_E^e + (F_c(q))_{E'}^e = \frac{1}{d_e} ([A(d_e V_E^e) + B(-d_e V_E^e)] q_E + [B(d_e V_E^e) + A(-d_e V_E^e)] q_{E'}).$$

If (AB3-s) holds, this quantity is equal to zero (this is the conservation of the convective flux), and if (AB3-w) holds we have, thanks to (AB1),

$$|(F_c(q))_E^e + (F_c(q))_{E'}^e| = \frac{1}{d_e} |(A(d_e V_E^e) + B(-d_e V_E^e))(q_E - q_{E'})| \leq C_4 \|V\|_{\infty} |q_E - q_{E'}|$$

for some  $C_4$  only dependent on  $A$  and  $B$ . Writing  $|q_E - q_{E'}| \leq |q_E - q_e| + |q_e - q_{E'}|$  and using again inequality (3.10) allows us to estimate the last term of (3.7) as follows:

$$\left| \sum_{e \in \mathcal{E}_h} |e|((F_c(q))_E^e + (F_c(q))_{E'}^e)(q_{E'} - q_e) \right| \leq C_5 \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e|(q_E - q_e)^2 \leq C_5 h\|(q, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \quad (3.12)$$

where  $C_5$  only depends on  $V, A, B$  and there holds that  $C_5 = 0$  if (AB3-s) holds.

The proof terminates by gathering inequalities (3.11) and (3.12) into (3.7). ■



## 3.1 Convergence of the method

### 3.1.1 Preliminary results

Proposition 3.2 below is the key point in the study of scheme (2.39)-(2.43), since it gives the inequality leading to the basic *a priori* estimates of the solution error. To state this proposition, we first notice that, thanks to (2.41), we can introduce the set of face values  $p_{\mathcal{E}_h} \in H(\mathcal{E}_h)$  such that (2.47) holds even if  $F_h$  is not conservative. To this purpose, we simply define  $p_e$  through  $|e|(p_E - p_e) = [F_E, G_E(E, e)]_E$  where  $G_E(E, e)_e = 1$  and  $G_E(E, e)_{e'} = 0$  for  $e \neq e'$ . Then, taking the vector  $G \in X_h$  that vanishes on all mesh faces except  $e$  and is such that  $G_E^e = 1$  and  $G_{E'}^e = -1$  in (2.41) allows us to show that  $p_e$  does not depend on the choice of the cell  $E$  such that  $e \in \partial E$ . This definition also ensures that  $p_e = 0$  whenever  $e \in \mathcal{E}_{h,\text{ext}}$ .

**Proposition 3.2** *Assume that (H1)-(H3) hold. Let  $\mathcal{D}_h$  be an admissible discretization of  $\Omega$  such that  $\theta \geq \text{regul}(\mathcal{D}_h)$ ; let  $F_c(q)$  be the convective flux of  $q \in Q_h$  given by (2.39)-(2.40) for the vector field  $V \in C^1(\bar{\Omega})^d$  with  $A$  and  $B$  satisfying Assumptions (AB1)-(AB3). Then, for all solution  $(p_h, F_h)$  to the HMM scheme (2.41)-(2.43),*

$$\sum_{E \in \Omega_h} [F_E, F_E]_E + \frac{1}{2} \int_{\Omega} \text{div}(V) p_h^2 \leq \int_{\Omega} f p_h + C_1 h \| (p_h, p_{\mathcal{E}_h}) \|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \quad (3.13)$$

where  $C_1$ , which is the same constant of Lemma 3.1, is non-negative, only depends on  $\theta, V, A, B$ , and is zero when (AB3-s) holds.

**Proof of Proposition 3.2** Let us take  $q = p_h$  in (2.42), use the flux conservation (2.43) and property (2.47) of face values to obtain:

$$\begin{aligned} \int_{\Omega} f p_h &= \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| \left( F_E^e + (F_c(p_h))_E^e \right) p_E \\ &= \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| \left( F_E^e + (F_c(p_h))_E^e \right) (p_E - p_e) \\ &= \sum_{E \in \Omega_h} [F_E, F_E]_E + \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F_c(p_h))_E^e (p_E - p_e). \end{aligned} \quad (3.14)$$

The proposition follows by applying Lemma 3.1 with  $q = p_h$  and  $q_{\mathcal{E}_h} = p_{\mathcal{E}_h}$ . ■

**Corollary 3.3** *Under the assumptions of Proposition 3.2, if  $V$  satisfies (H4) then, for all  $(p_h, F_h)$  solution to the scheme (2.41)-(2.43) we have*

$$\| (p_h, p_{\mathcal{E}_h}) \|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \lesssim \| f \|_{L^2(\Omega)} \| p_h \|_{L^2(\Omega)} + C_1 h \| (p_h, p_{\mathcal{E}_h}) \|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \quad (3.15)$$

where  $C_1$ , which is the same constant of Lemma 3.1 and Proposition 3.2, is non-negative, only depends on  $\theta, V, A, B$ , and is zero when (AB3-s) holds.

In particular, for all  $h$  small enough (or any  $h$  if (AB3-s) holds), the scheme (2.41)-(2.43) has a unique solution.

**Proof of Corollary 3.3** We apply Proposition 3.2 and use (H4) and the form (2.11)-(2.13) of the local scalar products  $([\cdot, \cdot]_E)_{E \in \Omega_h}$  to write, thanks to (C),

$$\sum_{E \in \Omega_h} |E| |\mathbf{v}_E(F_E)|^2 + \alpha \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} |T_{E,e}(F_E)|^2 \leq \int_{\Omega} f p_h + C_1 h \| (p_h, p_{\mathcal{E}_h}) \|_{1, \mathcal{D}_h, \mathcal{E}_h}^2. \quad (3.16)$$

From (2.47) and (2.11)-(2.13) we have

$$|e|(p_E - p_e) = |E| \Lambda_E \mathbf{v}_E(F_E) \cdot \mathbf{v}_E(G_E(e)) + T_E(G_E(e))^T \mathbb{B}_E T_E(F_E), \quad (3.17)$$

where  $G_E(e) \in X_E$  is equal to 1 on the face  $e$  and 0 on the other faces. But  $\mathbf{v}_E(G_E(e)) = -\frac{1}{|E|}\Lambda_E^{-1}|e|(\bar{x}_e - x_E)$ , and thus, by the bound on  $\text{regul}(\mathcal{D}_h)$ ,  $|\mathbf{v}_E(G_E(e))| \lesssim \frac{|e|d_{E,e}}{|E|}$  and, for all  $e' \in \partial E$ ,  $|T_{E,e'}(G_E(e))| \lesssim |G_E(e)^{e'}| + \frac{|e|d_{E,e}}{|E|}$ . In particular, by using the Cauchy-Schwarz inequality and (C), since  $\sum_{e' \in \partial E} |e'|d_{E,e'} = d|E|$ ,

$$\begin{aligned} |T_E(G_E(e))^T \mathbb{B}_E T_E(F_E)| &\lesssim \left( \sum_{e' \in \partial E} |e'|d_{E,e'} |T_{E,e'}(G_E(e))|^2 \right)^{1/2} \left( \sum_{e' \in \partial E} |e'|d_{E,e'} |T_{E,e'}(F_E)|^2 \right)^{1/2} \\ &\lesssim \left( |e|d_{E,e} + \frac{|e|^2 d_{E,e}^2}{|E|} \right)^{1/2} \left( \sum_{e' \in \partial E} |e'|d_{E,e'} |T_{E,e'}(F_E)|^2 \right)^{1/2}. \end{aligned}$$

Plugged into (3.17), this estimate and  $|E||\mathbf{v}_E(G_E(e))| \lesssim |e|d_{E,e}$  lead to

$$|p_E - p_e| \lesssim d_{E,e} |\mathbf{v}_E(F_E)| + \left( \frac{d_{E,e}}{|e|} + \frac{d_{E,e}^2}{|E|} \right)^{1/2} \left( \sum_{e' \in \partial E} |e'|d_{E,e'} |T_{E,e'}(F_E)|^2 \right)^{1/2}.$$

We then obtain, from (3.16),

$$\begin{aligned} \sum_{E \in \Omega_h} \sum_{e \in \partial E} \frac{|e|}{d_{E,e}} |p_E - p_e|^2 &\lesssim \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e|d_{E,e} |\mathbf{v}_E(F_E)|^2 \\ &\quad + \sum_{E \in \Omega_h} \sum_{e \in \partial E} \left( 1 + \frac{|e|d_{E,e}}{|E|} \right) \left( \sum_{e' \in \partial E} |e'|d_{E,e'} |T_{E,e'}(F_E)|^2 \right) \\ &\lesssim \int_{\Omega} f p_h + C_1 h \| (p_h, p_{\mathcal{E}_h}) \|_{1, \mathcal{D}_h, \mathcal{E}_h}^2, \end{aligned}$$

and the proof of (3.15) is completed.

Existence and uniqueness of the numerical solution readily follows from (3.15). In fact, when the right-hand side  $f$  vanishes, this inequality implies that the mesh-dependent norm  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}$  is zero, and, thus, that  $(p_h, p_{\mathcal{E}_h})$  are zero at least for a sufficiently small mesh size  $h$ . In such a case, the numerical flux  $F_h$  is also zero by (2.47). ■

**Remark 3.4** (Estimates for the hybrid discretization of the convection) *For the hybrid discretization in (2.45)-(2.49) with  $A$  and  $B$  satisfying (AB1)-(AB2) and (AB3h) there holds a similar result as that given in Proposition 3.2 and Corollary 3.3. However, the proof is simpler. In fact, by using (3.8) we have that*

$$\begin{aligned} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e|(F_{c, \mathcal{E}_h}(q, q_{\mathcal{E}_h}))_E^e (q_E - q_e) &= \frac{1}{2} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| V_E^e (q_E + q_e) (q_E - q_e) \\ &\quad + \frac{1}{2} \sum_{E \in \Omega_h} \sum_{e \in \partial E} \frac{|e|}{d_e} (A(d_e V_E^e) - B(d_e V_E^e)) (q_E - q_e)^2. \end{aligned}$$

*The right-hand side of this equation is similar to the right-hand side of (3.9) with  $q_e$  instead of  $q_{E'}$  and, reasoning as in the proof of Lemma 3.1, can be bounded from below by the the right-hand side of (3.11). The resulting estimate is then used in (3.14) with  $F_{c, \mathcal{E}_h}(p_h, p_{\mathcal{E}_h})$  instead of  $F_c(p_h)$  in order to prove Proposition 3.2.*

We conclude this preliminary subsection by reporting two technical lemmas that we will use in the analysis of the next subsection. The first lemma is a direct consequence of [36, Lemmas 5.2,5.3] and, for this reason, is given without proof.

**Lemma 3.5** (Discrete Sobolev inequalities) *Let  $\mathcal{D}_h$  be an admissible discretization of  $\Omega$  in the sense of Definition 2.2 with  $\theta > 0$  and such that  $\theta \leq \frac{d_{E,e}}{d_{E',e}} \leq \theta^{-1}$  for all  $e \in \mathcal{E}_{h,int}$ . Let  $r = \frac{2d}{d-2}$  if  $d > 2$  and  $r < +\infty$  if  $d = 2$ . Then, there exists a real positive constant  $C$  that only depends on  $\Omega$ ,  $\theta$  and  $r$  such that, for all  $q_h \in Q_h$ ,  $\|q_h\|_{L^r(\Omega)} \leq C\|q_h\|_{1,\mathcal{D}_h}$ .*

Let  $\mathbf{v}_h(F_h)$  denote the piecewise-constant function equal to  $\mathbf{v}_E(F_E)$  on  $E \in \Omega_h$  as defined in (2.12).

**Lemma 3.6** (Discrete Rellich theorem) *Let  $\Lambda : \Omega \rightarrow M_d(\mathbf{R})$  be a diffusion tensor satisfying hypothesis (H2); let  $(\mathcal{D}_h)_{h \rightarrow 0}$  be a family of admissible discretization of  $\Omega$  in the sense of Definition 2.2 with mesh size  $h$  tending to 0 and satisfying the regularity Assumptions (G1)-(G2); let  $(p_h, p_{\mathcal{E}_h}) \in Q_h \times H(\mathcal{E}_h)$  be a numerical scalar field such that  $\|(p_h, p_{\mathcal{E}_h})\|_{1,\mathcal{D}_h,\mathcal{E}_h}$  remains bounded as  $h \rightarrow 0$ ; let  $F_h = (F_E^e)_{E \in \Omega_h, e \in \partial E}$  be a collection of numbers that satisfy equation (2.47) for the assigned  $(p_h, p_{\mathcal{E}_h})$  and with the local scalar products defined accordingly to (2.11)-(2.13).*

*Then, there exists a scalar field  $p \in H_0^1(\Omega)$  such that, up to a subsequence as  $h \rightarrow 0$ ,*

- (i)  $p_h \rightarrow p$  in  $L^r(\Omega)$  for all  $r < \frac{2d}{d-2}$ ;
- (ii)  $\mathbf{v}_h(F_h) \rightarrow \nabla p$  weakly in  $L^2(\Omega)^d$ .

**Proof of Lemma 3.6** Using [36, Lemma 5.6], Lemma 3.5, Vitali's theorem and the fact that the quantity  $\|(p_h, p_{\mathcal{E}_h})\|_{1,\mathcal{D}_h,\mathcal{E}_h}$  is uniformly bounded ensures that  $(p_h)_{h \rightarrow 0}$  is relatively compact in  $L^r(\Omega)$  for all  $r < \frac{2d}{d-2}$ . After defining the discrete gradient  $\tilde{\nabla}(p_h, p_{\mathcal{E}_h}) : \Omega \rightarrow \mathbf{R}^d$  by

$$\forall E \in \Omega_h, \forall x \in E : \tilde{\nabla}(p_h, p_{\mathcal{E}_h})(x) = \frac{1}{|E|} \sum_{e \in \partial E} |e|(p_e - p_E)\mathbf{n}_E^e,$$

we see from the bound on  $\|(p_h, p_{\mathcal{E}_h})\|_{1,\mathcal{D}_h,\mathcal{E}_h}$  that  $\tilde{\nabla}(p_h, p_{\mathcal{E}_h})$  remains bounded in  $L^2(\Omega)^d$ . The technique used to prove [36, Lemma 5.7] ensures that if  $p_h \rightarrow p$  in  $L^2(\Omega)$  (up to a subsequence), then,  $p$  belongs to  $H_0^1(\Omega)$  and  $\tilde{\nabla}(p_h, p_{\mathcal{E}_h})$  is weakly convergent to  $\nabla p$  in  $L^2(\Omega)^d$ . The lemma is therefore true since the argument discussed in [33, Remark 2.7] implies that  $\tilde{\nabla}(p_h, p_{\mathcal{E}_h}) = \mathbf{v}_h(F_h)$  if  $(p_h, p_{\mathcal{E}_h}, F_h)$  are linked through (2.11)-(2.13) and (2.47). ■

### 3.1.2 Convergence without regularity assumption

Let us consider the HMM method on  $(\mathcal{D}_h)_{h \rightarrow 0}$ , a family of meshes that are admissible according to Definition 2.2, with mesh size  $h$  tending to 0 and all of which satisfy the regularity conditions (G1)-(G2). We also assume that all the local scalar products in the scheme formulation are defined by (2.11)-(2.13) through a set of symmetric and positive definite matrices  $(\mathbb{B}_E)_{E \in \Omega_h}$  that verifies the coercivity condition (C). Moreover, the numerical convection flux  $F_c(p_h)$  in (2.42) is built by using (2.39)-(2.40) through some instance of the functions  $A$  and  $B$  that satisfy (AB1)-(AB3). Finally, we recall that  $\mathbf{v}_h(F_h) : \Omega \rightarrow \mathbf{R}^d$  is the piecewise-constant function equal to  $\mathbf{v}_E(F_E)$  on  $E$  for all  $E \in \Omega_h$ . The convergence result of this sub-section is stated in the following theorem.

**Theorem 3.7** *Let  $p \in H_0^1(\Omega)$  be the weak solution to (2.1)-(2.2) under Assumptions (H1)-(H4), and  $(p_h, F_h)$  the numerical solution to problem (2.41)-(2.43) built along the guidelines summarized above. Then, for  $h \rightarrow 0$  there holds that:*

- (i)  $p_h \rightarrow p$  in  $L^r(\Omega)$  for all  $r < \frac{2d}{d-2}$ ;
- (ii)  $\mathbf{v}_h(F_h) \rightarrow \nabla p$  in  $L^2(\Omega)^d$ .

#### Proof of Theorem 3.7

The proof of Theorem 3.7 is based on compactness tools developed for Mixed Finite Volume or Hybrid Finite Volume for the pure diffusion equation [31, 36] and on techniques from the classical Finite Volume

schemes [26, 35] to handle the numerical convection term. We report the full proof for the sake of completeness since none of these methods has ever been formulated in the new HMM framework.

**Step 1:** *compactness of the approximate solutions.*

Using Corollary 3.3 we have  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \lesssim \|f\|_{L^2(\Omega)} \|p_h\|_{L^2(\Omega)}$  (at least for  $h$  small enough if (AB3-s) does not hold). In view of Lemma 3.5 and inequality (3.5), we obtain an upper bound on  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}$ . Then, the result of Lemma 3.6 implies the existence of a function  $p \in H_0^1(\Omega)$  such that, up to a subsequence,  $p_h \rightarrow p$  in  $L^r(\Omega)$  for all  $r < \frac{2d}{d-2}$  and  $\mathbf{v}_h(F_h) \rightarrow \nabla p$  weakly in  $L^2(\Omega)^d$ .

**Step 2:** *the limit function  $p$  is the weak solution to (2.1)-(2.2).* Since the exact solution is unique, this step allows us to prove the convergence to  $p$  of the whole sequence of discrete solutions  $p_h$  for  $h \rightarrow 0$ . We take  $\varphi \in C_c^\infty(\Omega)$ , define  $\varphi_h \in Q_h$  by  $\varphi_h = \varphi(x_E)$  on  $E \in \Omega_h$  and plug  $q = \varphi_h$  in (2.42). Since  $F_h + F_c(p_h)$  is conservative, we obtain:

$$\begin{aligned}
\int_{\Omega} f \varphi_h &= \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| F_E^e (\varphi(x_E) - \varphi(\bar{x}_e)) \\
&+ \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (\varphi(x_E) - \varphi(\bar{x}_e)) \frac{1}{d_e} \left( A(d_e V_E^e) p_E + B(d_e V_E^e) p_{E'} \right) \\
&= \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| F_E^e (x_E - \bar{x}_e) \cdot \nabla \varphi(x_E) + \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| F_E^e R_{E,e}^h(\varphi) \\
&+ \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (\varphi(x_E) - \varphi(\bar{x}_e)) \frac{1}{d_e} \left( A(d_e V_E^e) + B(d_e V_E^e) \right) p_E \\
&+ \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (\varphi(x_E) - \varphi(\bar{x}_e)) \frac{1}{d_e} B(d_e V_E^e) (p_{E'} - p_E) \\
&= \mathsf{T}_1 + \mathsf{T}_2 + \mathsf{T}_3 + \mathsf{T}_4
\end{aligned} \tag{3.18}$$

where the residual term  $R_{E,e}^h(\varphi)$  in  $\mathsf{T}_2$  is such that  $|R_{E,e}^h(\varphi)| \lesssim d_{E,e} h \|\nabla^2 \varphi\|_{\infty}$ .

By (2.12) we have that

$$\mathsf{T}_1 = \sum_{E \in \Omega_h} |E| \Lambda_E \mathbf{v}_E(F_E) \cdot \nabla \varphi(x_E) = \int_{\Omega} \Lambda \mathbf{v}_h(F_h) \cdot (\nabla \varphi)_h$$

where  $(\nabla \varphi)_h = \nabla \varphi(x_E)$  on  $E \in \Omega_h$ . The regularity of  $\varphi$  together with the weak convergence of  $\mathbf{v}_h(F_h)$  implies that

$$\mathsf{T}_1 \rightarrow \int_{\Omega} \Lambda \nabla p \cdot \nabla \varphi \quad \text{as } h \rightarrow 0. \tag{3.19}$$

From (2.13) we have  $|F_E^e| \lesssim |T_{E,e}(F_E)| + |\mathbf{v}_E(F_E)|$  and, since  $\|p_h\|_{L^2(\Omega)}$  and  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}$  are bounded, inequality (3.16) implies that

$$\sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} |F_E^e| \leq \left( \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} \right)^{1/2} \left( \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} |F_E^e|^2 \right)^{1/2} \lesssim 1$$

(recall that  $\sum_{e \in \partial E} |e| d_{E,e} = d|E|$ ). Therefore, we obtain that

$$|\mathsf{T}_2| \lesssim h \|\nabla^2 \varphi\|_{\infty} \rightarrow 0 \quad \text{as } h \rightarrow 0. \tag{3.20}$$

Assumption (AB2) makes it possible to show that

$$\begin{aligned}
\mathsf{T}_3 &= \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} |e| (\varphi(x_E) - \varphi(\bar{x}_e)) V_E^e \\
&= \sum_{E \in \Omega_h} p_E \varphi(x_E) \sum_{e \in \partial E} |e| V_E^e - \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} |e| \varphi(\bar{x}_e) V_E^e \\
&= \int_{\Omega} p_h \varphi_h \operatorname{div}(V) - \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} \int_e \varphi V \cdot \mathbf{n}_E^e + \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} \int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e \\
&= \int_{\Omega} p_h \varphi_h \operatorname{div}(V) - \int_{\Omega} p_h \operatorname{div}(\varphi V) + \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} \int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e.
\end{aligned}$$

The regularity of  $\varphi$  and the convergence of  $p_h$  ensure that, as  $h \rightarrow 0$ , the first two terms in this right-hand side tend to  $\int_{\Omega} p \varphi \operatorname{div}(V)$  and  $\int_{\Omega} p \operatorname{div}(\varphi V)$ . As for the last term, using the fact that  $\int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e$  vanishes for boundary faces ( $\varphi$  has a compact support) and is conservative for interior faces (i.e. changing  $E$  in  $E'$ , the cell on the other side of  $e$ , only changes the sign), we find

$$\begin{aligned}
\left| \sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} \int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e \right| &= \left| \sum_{E \in \Omega_h} \sum_{e \in \partial E} (p_E - p_e) \int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e \right| \\
&\lesssim h \|\nabla \varphi\|_{\infty} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| |p_E - p_e|.
\end{aligned}$$

But Cauchy-Schwarz inequality and the bound on  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}$  gives

$$\sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| |p_E - p_e| \leq (d |\Omega|)^{1/2} \|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h} \lesssim 1, \quad (3.21)$$

and  $\sum_{E \in \Omega_h} p_E \sum_{e \in \partial E} \int_e (\varphi - \varphi(\bar{x}_e)) V \cdot \mathbf{n}_E^e$  thus tends to 0 with  $h$ . We deduce that

$$\mathsf{T}_3 \rightarrow \int_{\Omega} p \varphi \operatorname{div}(V) - \int_{\Omega} p \operatorname{div}(\varphi V) = - \int_{\Omega} V p \cdot \nabla \varphi \quad \text{as } h \rightarrow 0. \quad (3.22)$$

To handle  $\mathsf{T}_4$ , we start noting that Assumption (AB1) implies that  $\frac{1}{d_e} |B(d_e V_E^e)| \lesssim 1$ . Thus, writing  $p_{E'} - p_E = p_{E'} - p_e + p_e - p_E$  and using (3.21), we obtain:

$$\begin{aligned}
|\mathsf{T}_4| &\lesssim h \|\nabla \varphi\|_{\infty} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (|p_{E'} - p_e| + |p_e - p_E|) \\
&\lesssim 2h \|\nabla \varphi\|_{\infty} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| |p_E - p_e| \rightarrow 0 \quad \text{as } h \rightarrow 0.
\end{aligned} \quad (3.23)$$

Eventually, the convergence properties (3.19), (3.20), (3.22) and (3.23) allows us to get the limit of (3.18) for  $h \rightarrow 0$  and show that  $p$  is the weak solution to (2.1)-(2.2).

**Step 3: strong convergence of the gradient.** Estimate (3.13) and Relation (2.11) imply that

$$\int_{\Omega} \Lambda \mathbf{v}_h(F_h) \cdot \mathbf{v}_h(F_h) + \frac{1}{2} \int_{\Omega} \operatorname{div}(V) p_h^2 \leq \int_{\Omega} f p_h + Ch \|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2.$$

Taking the upper limit of this inequality, recalling that  $\|(p_h, p_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}$  stays bounded and noting that  $p_h$  is strongly convergent to  $p$  in  $L^2(\Omega)$  and that  $p$  is the weak solution to (2.1)-(2.2) lead to

$$\limsup_{h \rightarrow 0} \int_{\Omega} \Lambda \mathbf{v}_h(F_h) \cdot \mathbf{v}_h(F_h) + \frac{1}{2} \int_{\Omega} \operatorname{div}(V) p^2 \leq \int_{\Omega} f p = \int_{\Omega} \Lambda \nabla p \cdot \nabla p + \frac{1}{2} \int_{\Omega} \operatorname{div}(V) p^2,$$

from which we deduce that

$$\limsup_{h \rightarrow 0} \int_{\Omega} \Lambda \mathbf{v}_h(F_h) \cdot \mathbf{v}_h(F_h) \leq \int_{\Omega} \Lambda \nabla p \cdot \nabla p. \quad (3.24)$$

Since  $\mathbf{w} \rightarrow (\int_{\Omega} \Lambda \mathbf{w} \cdot \mathbf{w})^{1/2}$  is a norm in  $L^2(\Omega)^d$  equivalent to the usual norm, equation (3.24) proves that the weak convergence  $\mathbf{v}_h(F_h) \rightarrow \nabla p$  in  $L^2(\Omega)^d$  is, in fact, strong. ■

**Remark 3.8** *Adjusting these same arguments makes it possible to prove a similar convergence result for the hybrid HMM formulation (2.47)-(2.49), which is based on the numerical convection flux (2.45)-(2.46) instead of (2.39)-(2.40).*

### 3.1.3 About the regularity assumption on $V$

Oftentimes, the velocity field  $V$  is not given but comes from the resolution of another problem (see e.g. [25]). In this case, it is not obvious that it satisfies the regularity assumption  $V \in C^1(\bar{\Omega})^d$ : we can in general ensure that  $V \in H(\text{div}, \Omega)$ , but not more. How does this impact the preceding convergence study?

We first of all have to be able to define the fluxes  $V_E^e$  of the velocity; this is in general quite straightforward, either using (2.19) and the fact that  $V$  belongs to  $H(\text{div}, \Omega)$ , or even more directly by looking at the discretization of the equation providing  $V$  (this discretization usually also provides the fluxes of the velocity, as in [25]). The minimal requirement on these fluxes is their conservativity

$$\forall e \in \mathcal{E}_{h,\text{int}}, : V_E^e + V_{E'}^e = 0$$

(where  $E, E' \in \Omega_h$  are the two elements such that  $e \subset \partial E \cap \partial E'$  for every  $e \in \mathcal{E}_{h,\text{int}}$ ) and their compatibility with the coercitivity assumption  $\text{div}(V) \geq 0$ :

$$\forall E \in \Omega_h, : \sum_{e \in \partial E} |e| V_E^e \geq 0$$

(usually,  $\sum_{e \in \partial E} |e| V_E^e$  plays the role of an approximation of  $\int_E \text{div}(V)$ ). Under these two requirements and the strong version of Assumption (AB3) (i.e. (AB3s)), it is then easy to see that the *a priori* estimates still hold (see Lemma 3.1, Proposition 3.2 and Corollary 3.3).

As for the convergence (Theorem 3.7), we have to check if  $\mathsf{T}_3$  and  $\mathsf{T}_4$  behave well. For  $\mathsf{T}_3$  we need that

$$\forall E \in \Omega_h, : \sum_{e \in \partial E} |e| V_E^e = \int_E \text{div}(V)$$

(or at least that  $\sum_{e \in \partial E} |e| V_E^e$  approximates  $\int_E \text{div}(V)$  as the size of the mesh tends to 0), which is usually the case from the definition of  $V_E^e$  using (2.19) or an expression of these fluxes coming from the resolution of another elliptic equation, and that, for any smooth function  $\varphi$  with compact support, denoting by  $\Phi_h : \Omega \rightarrow \mathbf{R}$  the function defined by

$$\forall E \in \Omega_h, \forall x \in E : \Phi_h(x) = \sum_{e \in \partial E} \varphi(\bar{x}_e) |e| V_E^e,$$

the function  $\Phi_h$  weakly converges in  $L^2(\Omega)$ , as  $h \rightarrow 0$ ,  $\text{div}(\varphi V)$ . Since, for any  $W \in H(\text{div}, \Omega)$ , defining  $W_E^e = \frac{1}{|e|} \int_e W \cdot \mathbf{n}_E^e$  (in the usual weak sense), we have

$$|W_E^e|^2 \leq Ch^{-d} \|W\|_{L^2(E)}^2 + Ch^{-d+2} \|\text{div}(W)\|_{L^2(E)}^2 \quad (3.25)$$

(this is the usual Agmon scaling of trace estimates), the estimates we provide in the proof of Theorem 3.7 on the last part of  $\mathsf{T}_3$  indicate that  $\Phi_h$  behaves as needed if  $V_E^e$  comes from 2.19 with  $V \in H(\text{div}, \Omega)$ ; if

these velocity fluxes come from the approximation of another elliptic equation, then the expected behavior of  $\Phi_h$  is usually a straightforward consequence of the properties of the scheme used on this other equation (see e.g. [25]).

For  $\mathbb{T}_4$ , we require that

$$\sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| d_{E,e} |V_E^e|^2 \quad (3.26)$$

remains bounded as the size of the mesh tends to 0. When  $V_E^e$  comes from the approximation of an elliptic equation (i.e.  $V_E^e$  is the “ $F_E^e$ ” of this other equation), then this estimate is usually a basic one (for HMM methods, for example, it is a direct consequence of (S1) or (C)); if  $V_E^e$  is constructed from (2.19) with  $V \in H(\text{div}, \Omega)$ , then (3.25) shows that (3.26) also remains bounded independently of the mesh size.

In other words, although the preceding study has been made, for the sake of simplicity, with regular velocity fields, it is easy to adapt to more realistic fields, and the convergence result still hold for these fields.

### 3.2 Error estimates

In the theoretical developments of this section we assume that (HG) and (ME) hold.

Now, we consider the bilinear form

$$\begin{aligned} \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G'_h, q'_h, q'_{\mathcal{E}_h}) &= [G_h, G'_h]_{\widehat{X}_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (G'_h)_E^e (q_E - q_e) \\ &+ [\text{div}_h(G_h + F_c(q_h)), q'_h]_{Q_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (G_h + F_c(q_h))_E^e q'_e \end{aligned} \quad (3.27)$$

for all couple of triplets  $(G_h, q_h, q_{\mathcal{E}_h})$  and  $(G'_h, q'_h, q'_{\mathcal{E}_h})$  in  $\widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$ . Problem (2.41)-(2.43) can be reformulated as

Find  $(F_h, p_h, p_{\mathcal{E}_h}) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$  such that:

$$\mathcal{B}(F_h, p_h, p_{\mathcal{E}_h}; G_h, q_h, q_{\mathcal{E}_h}) = [f^I, q_h]_{Q_h} \quad \forall (G_h, q_h, q_{\mathcal{E}_h}) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h). \quad (3.28)$$

In order to prove the convergence result, we need the following stability lemma.

**Lemma 3.9** *Assume (AB1)-(AB3) with either  $h$  small enough if (AB3-w) holds or any  $h$  if (AB3-s) holds. For any triple  $(G_h, q_h, q_{\mathcal{E}_h}) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$  there exists a triple  $(G'_h, q'_h, q'_{\mathcal{E}_h}) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$  with*

$$\|G'_h\|_{\widehat{X}_h} + \|q'_h\|_{1, \mathcal{D}_h} + \|(q'_h, q'_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h} \leq 1 \quad (3.29)$$

for which there holds that:

$$\mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G'_h, q'_h, q'_{\mathcal{E}_h}) \gtrsim \|G_h\|_{\widehat{X}_h} + \|q_h\|_{1, \mathcal{D}_h} + \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}. \quad (3.30)$$

#### Proof of Lemma 3.9.

A straightforward calculation shows that

$$\mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G_h, q_h, q_{\mathcal{E}_h}) = \|G_h\|_{\widehat{X}_h}^2 + \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F_c(q_h))_E^e (q_E - q_e). \quad (3.31)$$

Since  $\text{div}(V) \geq 0$ , applying Lemma 3.1 yields the inequality

$$\mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G_h, q_h, q_{\mathcal{E}_h}) \geq \|G_h\|_{\widehat{X}_h}^2 - C_1 h \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2. \quad (3.32)$$

The non negative real constant  $C_1$ , which is provided by Lemma 3.1, is zero if Assumption (AB3-s) holds.

Let us consider the *non-conservative* vector field  $\widehat{G}_h \in \widehat{X}_h$  given by

$$\forall E \in \Omega_h, \forall e \in \partial E : (\widehat{G}_h)_E^e = \frac{q_e - q_E}{d_{E,e}}.$$

Since (M2) implies that  $|E| \lesssim |e|d_{E,e}$ , we have that

$$\|\widehat{G}_h\|_{\widehat{X}_h}^2 = \sum_{E \in \Omega_h} \|\widehat{G}_E\|_E^2 \leq \sigma^* \sum_{E \in \Omega_h} \sum_{e \in \partial E} \frac{|E|}{d_{E,e}^2} (q_E - q_e)^2 \leq \widehat{C} \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2, \quad (3.33)$$

where  $\widehat{C} > 0$  is independent of  $h$  and only depends on the constant  $\sigma^*$  of Assumption (S1) and on the mesh regularity constants of (M2). We infer, from the Cauchy-Schwarz inequality and Young's inequality, that

$$\left| [G_h, \widehat{G}_h]_{\widehat{X}_h} \right| \leq \|G_h\|_{\widehat{X}_h} \|\widehat{G}_h\|_{\widehat{X}_h} \leq \frac{\widehat{C}}{2} \|G_h\|_{\widehat{X}_h}^2 + \frac{1}{2} \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2.$$

By using the definitions (3.27) and (3.4) we obtain that

$$\begin{aligned} \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; \widehat{G}_h, 0, 0) &= [G_h, \widehat{G}_h]_{\widehat{X}_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (\widehat{G}_h)_E^e (q_E - q_e) \\ &= [G_h, \widehat{G}_h]_{\widehat{X}_h} + \sum_{E \in \Omega_h} \sum_{e \in \partial E} \frac{|e|}{d_{E,e}} (q_E - q_e)^2 \\ &\geq -\frac{\widehat{C}}{2} \|G_h\|_{\widehat{X}_h}^2 + \frac{1}{2} \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2. \end{aligned} \quad (3.34)$$

In the following development, it is natural to use the  $H^1$ -like norm for the elements of  $Q_h$  given by:

$$\|q_h\|_{1,h}^2 = \sum_{e \in \mathcal{E}_h} |e| h_e^{-1} (\llbracket q_h \rrbracket_e)^2, \quad (3.35)$$

where  $\llbracket q_h \rrbracket_e$  is the jump of  $q_h$  at edge  $e$  defined accordingly to (2.31). Assumptions (HG)-(ME) implies that mesh-dependent norm  $\|\cdot\|_{1,h}$  in (3.35) is *uniformly* equivalent to norm  $\|\cdot\|_{1, \mathcal{D}_h}$  in (3.3), i.e., there exists two positive constants  $\nu_*$  and  $\nu^*$  independent of the mesh size  $h$  such that there holds:

$$\nu_* \|\cdot\|_{1, \mathcal{D}_h} \leq \|\cdot\|_{1,h} \leq \nu^* \|\cdot\|_{1, \mathcal{D}_h} \quad (3.36)$$

for every instance of the admissible mesh family  $(\mathcal{D}_h)_h$ . As in [11], let us consider the *conservative* vector field  $\widetilde{G}_h \in X_h$  given by

$$\forall E \in \Omega_h, \forall e \in \partial E : (\widetilde{G}_h)_E^e = h_e^{-1} (q_{E'} - q_E).$$

Since (M2) implies that  $|E| \lesssim |e|h_e$ , we have that

$$\|\widetilde{G}_h\|_{\widehat{X}_h}^2 = \sum_{E \in \Omega_h} \|\widetilde{G}_E\|_E^2 \leq \sigma^* \sum_{E \in \Omega_h} \sum_{e \in \partial E} |E| h_e^{-2} (q_{E'} - q_E)^2 \leq \widetilde{C} \|q_h\|_{1,h}^2, \quad (3.37)$$

where  $\widetilde{C}$  is independent of  $h$  and only depends on  $\sigma^*$  and the mesh regularity constants of (M2). We then apply the Cauchy-Schwarz inequality and Young's inequality to obtain:

$$\left| [G_h, \widetilde{G}_h]_{\widehat{X}_h} \right| \leq \|G_h\|_{\widehat{X}_h} \|\widetilde{G}_h\|_{\widehat{X}_h} \leq \frac{\widetilde{C}}{2} \|G_h\|_{\widehat{X}_h}^2 + \frac{1}{2\widetilde{C}} \|\widetilde{G}_h\|_{\widehat{X}_h}^2 \leq \frac{\widetilde{C}}{2} \|G_h\|_{\widehat{X}_h}^2 + \frac{1}{2} \|q_h\|_{1,h}^2.$$



By using the definition of  $\tilde{G}_h$  and norm definition (3.4) we obtain that

$$\begin{aligned} \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (\tilde{G}_h)_E^e (q_E - q_e) &= \sum_{e \in \mathcal{E}_h} |e| \left( (\tilde{G}_h)_E^e q_E + (\tilde{G}_h)_{E'}^e q_{E'} \right) + \sum_{e \in \mathcal{E}_h} |e| \left( (\tilde{G}_h)_E^e + (\tilde{G}_h)_{E'}^e \right) q_e \\ &= - \sum_{e \in \mathcal{E}_h} |e| h_e^{-1} (q_{E'} - q_E)^2 = - \|q_h\|_{1,h}^2. \end{aligned}$$

Therefore, we have that

$$\mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; \tilde{G}_h, 0, 0) = \left[ G_h, \tilde{G}_h \right]_{\tilde{X}_h} + \|q_h\|_{1,h}^2 \geq -\frac{\tilde{C}}{2} \|G_h\|_{\tilde{X}_h}^2 + \frac{1}{2} \|q_h\|_{1,h}^2. \quad (3.38)$$

Let  $G'_h = \theta G_h + \hat{G}_h + \tilde{G}_h$  for some value of  $\theta$ ,  $q'_h = q_h$  and  $q'_{\mathcal{E}_h} = q_{\mathcal{E}_h}$ . From (3.32), (3.34) and (3.38) there holds:

$$\begin{aligned} &\mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G'_h, q'_h, q'_{\mathcal{E}_h}) \\ &= \theta \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G_h, q_h, q_{\mathcal{E}_h}) + \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; \hat{G}_h, 0, 0) + \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; \tilde{G}_h, 0, 0) \\ &\geq \left( \theta - \frac{\hat{C}}{2} - \frac{\tilde{C}}{2} \right) \|G_h\|_{\tilde{X}_h}^2 + \frac{1}{2} \|q_h\|_{1,h}^2 + \left( \frac{1}{2} - \theta C_1 h \right) \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2. \end{aligned} \quad (3.39)$$

Now, we take  $\theta = (1 + \hat{C} + \tilde{C})/2$  and we obtain the inequality

$$\|G_h\|_{\tilde{X}_h}^2 + \|q_h\|_{1,h}^2 + \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}^2 \lesssim \mathcal{B}(G_h, q_h, q_{\mathcal{E}_h}; G'_h, q'_h, q'_{\mathcal{E}_h}), \quad (3.40)$$

which holds for  $h$  small enough under Assumption (AB3-w), and for any  $h$  under Assumption (AB3-s) because  $C_1 = 0$  in this case. Using inequalities (3.33) and (3.37) allows us to obtain:

$$\begin{aligned} \|G'_h\|_{\tilde{X}_h} + \|q'_h\|_{1,h} + \|(q'_h, q'_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h} &\leq \theta \|G_h\|_{\tilde{X}_h} + \left(1 + \sqrt{\hat{C}}\right) \|q_h\|_{1,h} \\ &\quad + \left(1 + \sqrt{\tilde{C}}\right) \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}. \end{aligned} \quad (3.41)$$

Lemma's inequalities (3.29)-(3.30) follow from (3.40)-(3.41) by rescaling the three discrete fields  $G'_h$ ,  $q'_h$ , and  $q'_{\mathcal{E}_h}$  by the positive factor  $\max\left(\theta, 1 + \sqrt{\hat{C}}, 1 + \sqrt{\tilde{C}}\right) \left(\|G_h\|_{\tilde{X}_h} + \|q_h\|_{1,h} + \|(q_h, q_{\mathcal{E}_h})\|_{1, \mathcal{D}_h, \mathcal{E}_h}\right)$ . ■

The following technical lemma provides us with an estimate for the interpolation of a vector field which is locally in  $(H^1(E))^d$ .

**Lemma 3.10** *Let  $G \in (H^1(E))^d$  and  $G^I$  the interpolated field (2.8). Then, we have that*

$$\|G^I\|_E \lesssim \|G\|_{L^2(E)} + h_E |G|_{H^1(E)} \quad (3.42)$$

**Proof of Lemma 3.10** Using the stability condition of Assumption (S1), the Agmon inequality from (M3), and the scaling  $|E|/|e| \lesssim h_E$ , which is a consequence of (M2), it readily follows that

$$\begin{aligned} \|(G^I)\|_E^2 &\lesssim |E| \sum_{e \in \partial E} (G_e^e)^2 \lesssim |E| \sum_{e \in \partial E} |e|^{-1} \|G\|_{L^2(e)}^2 \lesssim h_E \left( h_E^{-1} \|G\|_{L^2(E)}^2 + h_E |G|_{H^1(E)}^2 \right) \\ &\lesssim \|G\|_{L^2(E)}^2 + h_E^2 |G|_{H^1(E)}^2, \end{aligned}$$

from which the lemma's statement immediately follows. ■

We can now prove the main result of this sub-section that is stated in the following theorem. This theorem provides a bound on the approximation error that is defined by comparing the numerical solution  $(p_h, F_h, p_{\mathcal{E}_h}) \in Q_h \times \tilde{X}_h \times H(\mathcal{E}_h)$  with the interpolations  $p^I$  and  $F^I$  of the exact solution and flux given by (2.8) and by the interpolated field  $p^J = \{(p^J)^e\}^{e \in \mathcal{E}_h} \in H(\mathcal{E}_h)$  given by

$$\forall e \in \mathcal{E}_h : (p^J)^e = \frac{1}{|e|} \int_e p. \quad (3.43)$$

**Theorem 3.11** *Let  $p$  be the solution of the continuous problem (2.1)-(2.2) under Assumptions (H1)-(H4) with  $\Lambda$  locally Lipschitz continuous on  $\Omega_h$ , c.f. Remark 2.5, and  $F$  given by (2.3). Let  $(F_h, p_h)$  be the solution of problem (2.41)-(2.42) under Assumptions (HG)-(ME) and (AB1)-(AB3) with either  $h$  small enough if (AB3-w) holds or any  $h$  if (AB3-s) holds. Then, there holds that:*

$$\|F_h - F^I\|_{\widehat{X}_h} + \|p_h - p^I\|_{1, \mathcal{D}_h} + \|(p_h - p^I, p_{\mathcal{E}_h} - p^J)\|_{1, \mathcal{D}_h, \mathcal{E}_h} \lesssim h \|p\|_{H^2(\Omega)}. \quad (3.44)$$

**Proof of Theorem 3.11.** Let us consider the triplet of error fields  $(F_h - F^I, p_h - p^I, p_{\mathcal{E}_h} - p^J) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$ . Due to Lemma 3.9 there exist a triplet  $(G_h, q_h, q_{\mathcal{E}_h}) \in \widehat{X}_h \times Q_h \times H(\mathcal{E}_h)$  with

$$\|G_h\|_{\widehat{X}_h} + \|q_h\|_{1, \mathcal{D}_h} + \|q_{\mathcal{E}_h}\|_{1, \mathcal{D}_h, \mathcal{E}_h} \leq 1 \quad (3.45)$$

such that

$$\begin{aligned} & \|F_h - F^I\|_{\widehat{X}_h} + \|p_h - p^I\|_{1, \mathcal{D}_h} + \|(p_h - p^I, p_{\mathcal{E}_h} - p^J)\|_{1, \mathcal{D}_h, \mathcal{E}_h} \\ & \lesssim \mathcal{B}(F_h - F^I, p_h - p^I, p_{\mathcal{E}_h} - p^J; G_h, q_h, q_{\mathcal{E}_h}). \end{aligned} \quad (3.46)$$

By using equations (3.28) a straightforward calculation gives:

$$\mathcal{B}(F_h - F^I, p_h - p^I, p_{\mathcal{E}_h} - p^J; G_h, q_h, q_{\mathcal{E}_h}) = \mathsf{T}_1 + \mathsf{T}_2 + \mathsf{T}_3 \quad (3.47)$$

where

$$\begin{aligned} \mathsf{T}_1 &= [p^I, \operatorname{div}_h(G_h)]_{Q_h} - [F^I, G_h]_{\widehat{X}_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (p^J)^e G_E^e, \\ \mathsf{T}_2 &= [f^I, q_h]_{Q_h} - [\operatorname{div}_h(F^I) + \operatorname{div}_h(F_c(p^I)), q_h]_{Q_h}, \\ \mathsf{T}_3 &= \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F^I + F_c(p^I))_E^e q_e. \end{aligned}$$

For convenience, we will separately bound  $\mathsf{T}_1$  and  $\mathsf{T}_2 + \mathsf{T}_3$ . To this purpose, let us first introduce the discontinuous  $\Omega_h$ -piecewise linear function  $p^1$ , which is such that  $p^1|_E$  is the  $L^2$  orthogonal projection of  $p$  on the linear polynomials defined on  $E \in \Omega_h$ . Let us start by noting that  $\|p - p^1\|_{L^2(E)} \leq \|p - q^1\|_{L^2(E)}$  for any linear polynomial  $q^1$  defined in  $E$ ; hence, taking  $q^1 = \mathcal{L}_1(p)$ , the linear interpolation of  $p$  on  $E$  provided by (M4), allows us to use the estimate for the interpolation error. Moreover, adding and subtracting  $\mathcal{L}_1(p)$ , applying the triangular inequality, using (M4) and a standard inverse inequality yields:

$$|p - p^1|_{H^1(E)} \leq |p - \mathcal{L}_1(p)|_{H^1(E)} + |\mathcal{L}_1(p) - p^1|_{H^1(E)} \lesssim h_E |p|_{H^2(E)} + h_E^{-1} \|\mathcal{L}_1(p) - p^1\|_{L^2(E)} \quad (3.48)$$

The second term in the last inequality of (3.48) is developed by adding and subtracting  $p$ , applying the triangular inequality and the estimate for the interpolation error of (M4):

$$\|\mathcal{L}_1(p) - p^1\|_{L^2(E)} \leq \|\mathcal{L}_1(p) - p\|_{L^2(E)} + \|p - p^1\|_{L^2(E)} \lesssim h_E^2 |p|_{H^2(E)}. \quad (3.49)$$

Substituting (3.49) in (3.48) yields the final inequality:

$$\|p - p^1\|_{L^2(E)} + h_E |p - p^1|_{H^1(E)} \lesssim h_E^2 \|p\|_{H^2(E)}. \quad (3.50)$$

We recall that, for convenience, we may identify the elements of  $Q_h$  with the piecewise constant functions whose restriction to each cell  $E$  is the degree-of-freedom of that cell. Therefore, it is possible to reformulate the first term of  $\mathsf{T}_1$  as an  $L^2$ -scalar product, so that  $[p^I, \operatorname{div}_h(G_h)]_{Q_h} = [p, \operatorname{div}_h(G_h)]_{L^2}$ . Then, we split  $\mathsf{T}_1$  in four sub-terms by recalling that  $F = -\Lambda \nabla p$ , adding and subtracting  $(\Lambda_E \nabla p)^I$  and

$(\Lambda_E \nabla p^1)^I$ , using the local consistency assumption (S2), and noting that  $|e|(p^J)^e = \int_e p$ . We have the following developments:

$$\begin{aligned}
\mathsf{T}_1 &= [p - p^1, \operatorname{div}_h(G_h)]_{L^2} + [p^1, \operatorname{div}_h(G_h)]_{L^2} - [F^I, G_h]_{\hat{X}_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} G_E^e \int_e p \\
&= \mathsf{T}_{1,1} + [p^1, \operatorname{div}_h(G_h)]_{L^2} + [(\Lambda \nabla p)^I, G_h]_{\hat{X}_h} - \sum_{E \in \Omega_h} \sum_{e \in \partial E} G_E^e \int_e p \\
&= \mathsf{T}_{1,1} + [p^1, \operatorname{div}_h(G_h)]_{L^2} + \sum_{E \in \Omega_h} [(\Lambda_E \nabla p)^I, G_h]_E \\
&\quad + \sum_{E \in \Omega_h} [((\Lambda - \Lambda_E) \nabla p)^I, G_h]_E - \sum_{E \in \Omega_h} \sum_{e \in \partial E} G_E^e \int_e p \\
&= \mathsf{T}_{1,1} + \sum_{E \in \Omega_h} \left( [\operatorname{div}_h(G_h), p^1]_{L^2(E)} + [(\Lambda_E \nabla p^1)^I, G_h]_E \right) + \sum_{E \in \Omega_h} [(\Lambda_E \nabla (p - p^1))^I, G_h]_E \\
&\quad + \sum_{E \in \Omega_h} [(\Lambda - \Lambda_E) \nabla p]^I, G_h]_E - \sum_{E \in \Omega_h} \sum_{e \in \partial E} G_E^e \int_e p \\
&= \mathsf{T}_{1,1} + \sum_{E \in \Omega_h} \sum_{e \in \partial E} G_E^e \int_e (p^1 - p) + \sum_{E \in \Omega_h} [(\Lambda_E \nabla (p - p^1))^I, G_h]_E + \sum_{E \in \Omega_h} [(\Lambda - \Lambda_E) \nabla p]^I, G_h]_E \\
&= \mathsf{T}_{1,1} + \mathsf{T}_{1,2} + \mathsf{T}_{1,3} + \mathsf{T}_{1,4} .
\end{aligned}$$

To estimate  $\mathsf{T}_{1,1}$ , let us first note that the definition of  $\operatorname{div}_h$ , the scalings  $h_E^d \leq |E| \lesssim h_E^d$  and  $|e| \lesssim h_E^{d-1}$  from (M2) and Assumption (S1) imply that

$$\|\operatorname{div}_h(G_h)\|_{L^2(E)}^2 = |E| |(\operatorname{div}_h(G_h))_E|^2 \lesssim \frac{1}{|E|} \sum_{e \in \partial E} |e|^2 (G_E^e)^2 \lesssim h_E^{-2} \|G_h\|_E^2 . \quad (3.51)$$

Thus, using the Cauchy-Schwarz inequality for each scalar product in  $X_E$ , error estimate (3.50), inequality (3.51), the Cauchy-Schwarz inequality again, and finally noting that (3.45) implies that  $\|G_h\|_{\hat{X}_h} \leq 1$  yield

$$\begin{aligned}
\mathsf{T}_{1,1} &\lesssim \sum_{E \in \Omega_h} \|p - p^1\|_{L^2(E)} \|\operatorname{div}_h(G_h)\|_{L^2(E)} \lesssim \sum_{E \in \Omega_h} \left( h_E^2 |p|_{H^2(E)} \right) \left( h_E^{-1} \|G_h\|_E \right) \\
&\lesssim h \left( \sum_{E \in \Omega_h} |p|_{H^2(E)}^2 \right)^{1/2} \left( \sum_{E \in \Omega_h} \|G_h\|_E^2 \right)^{1/2} \lesssim h |p|_{H^2(\Omega)} . \quad (3.52)
\end{aligned}$$

The second term is bounded using a scaling argument and inequality (3.45). We obtain that:

$$\mathsf{T}_{1,2} \lesssim h \|p\|_{H^2(\Omega)} \|G_h\|_{\hat{X}_h} \lesssim h \|p\|_{H^2(\Omega)} . \quad (3.53)$$

To get an upper bound for  $\mathsf{T}_{1,3}$ , we use the Cauchy-Schwarz inequality for the local scalar product in  $X_E$ , the result of Lemma 3.10, an upper bound on  $\Lambda_E$  that easily follows from the upper bound of  $\Lambda$  in (H2), again the Cauchy-Schwarz inequality, the estimate of the interpolation error given by (3.50), which follows from (M4), and the fact that  $\|G_h\|_{\hat{X}_h} \leq 1$  due to inequality (3.45). We obtain the following chain

of inequalities:

$$\begin{aligned}
\mathsf{T}_{1,3} &= \sum_{E \in \Omega_h} [(\Lambda_E \nabla(p - p^1))^I, G_h]_E \\
&\leq \sum_{E \in \Omega_h} \|(\Lambda_E \nabla(p - p^1))^I\|_E \|G_h\|_E \\
&\lesssim \sum_{E \in \Omega_h} \left( \|\Lambda_E \nabla(p - p^1)\|_{L^2(E)} + h_E |\Lambda_E \nabla(p - p^1)|_{H^1(E)} \right) \|G_h\|_E \\
&\lesssim \sum_{E \in \Omega_h} (|p - p^1|_{H^1(E)} + h_E |p|_{H^2(E)}) \|G_h\|_E \\
&\lesssim \left( \sum_{E \in \Omega_h} |p - p^1|_{H^1(E)}^2 + h_E^2 |p|_{H^2(E)}^2 \right)^{1/2} \left( \sum_{E \in \Omega_h} \|G_h\|_E^2 \right)^{1/2} \\
&\lesssim \left( \sum_{E \in \Omega_h} h_E^2 |p|_{H^2(E)}^2 \right)^{1/2} \|G_h\|_{\hat{X}_h} \\
&\lesssim h |p|_{H^2(\Omega)}. \tag{3.54}
\end{aligned}$$

Using Cauchy-Schwarz inequality and inequality (3.42) we get

$$\begin{aligned}
\mathsf{T}_{1,4} &\lesssim \sum_{E \in \Omega_h} \|((\Lambda - \Lambda_E) \nabla p)^I\|_E \|G_h\|_E \\
&\lesssim \left( \sum_{E \in \Omega_h} \|((\Lambda - \Lambda_E) \nabla p)^I\|_E^2 \right)^{1/2} \\
&\lesssim \left( \sum_{E \in \Omega_h} \|(\Lambda - \Lambda_E) \nabla p\|_{L^2(E)}^2 + h_E^2 |(\Lambda - \Lambda_E) \nabla p|_{H^1(E)}^2 \right)^{1/2}. \tag{3.55}
\end{aligned}$$

Due to the definition of  $\Lambda_E$  and since the restriction  $\Lambda|_E$  belongs to  $W^{1,\infty}(E)$  for all  $E \in \Omega_h$ , we obtain that

$$\|\Lambda - \Lambda_E\|_{L^\infty(E)} + h_E |\Lambda - \Lambda_E|_{W^{1,\infty}(E)} \lesssim h_E \quad \forall E \in \Omega_h.$$

Combining the above bound with (3.55) easily yields

$$\mathsf{T}_{1,4} \lesssim h \|p\|_{H^2(\Omega)}. \tag{3.56}$$

Combining (3.52), (3.53), (3.54), and (3.56) yields the following upper bound of  $\mathsf{T}_1$

$$\mathsf{T}_1 \lesssim h \|p\|_{H^2(\Omega)}. \tag{3.57}$$

To get an upper bound for  $\mathsf{T}_2 + \mathsf{T}_3$  we note that using the commuting property of the divergence operator (2.9), c.f. also Remark 2.4, the flux definition given in (2.3), and the model's equation (2.1) allows us to write:

$$\operatorname{div}_h(F^I) = (\operatorname{div}(F))^I = f^I - (\operatorname{div}(Vp))^I. \tag{3.58}$$

Equation (3.58) makes it possible to reformulate  $\mathsf{T}_2$  as follows:

$$\mathsf{T}_2 = [(\operatorname{div}(Vp))^I - \operatorname{div}_h(F_c(p^I)), q_h]_{Q_h}. \tag{3.59}$$

As before, we identify the elements of  $Q_h$  with the space of  $\Omega_h$ -piecewise constant functions, and the scalar product in  $Q_h$  with the  $L^2$  scalar product. Then, we split  $\mathsf{T}_2$  into two sub-terms by applying the divergence theorem to each cell's contribution and adding and subtracting the term  $V_E^e p$ :

$$\mathsf{T}_2 = \sum_{E \in \Omega_h} \sum_{e \in \partial E} q_E \int_e (V \cdot \mathbf{n}_E^e - V_E^e) p + \sum_{E \in \Omega_h} \sum_{e \in \partial E} q_E \int_e (V_E^e p - (F_c(p^I))^e) = \mathsf{T}_{2,1} + \mathsf{T}_{2,2}. \tag{3.60}$$

Noting that  $V_E^e + V_{E'}^e = 0$  and  $\mathbf{n}_E^e + \mathbf{n}_{E'}^e = 0$  for any  $e \in \mathcal{E}_{h,\text{int}}$ , and using definition (2.31) for the jump of  $q_h$ , i.e.,  $[[q_h]]_e$ , allows us to reformulate  $\mathsf{T}_{2,1}$  as follows

$$\mathsf{T}_{2,1} = \sum_{e \in \mathcal{E}_h} [[q_h]]_e \int_e (V \cdot \mathbf{n}_E^e - V_E^e) p$$

where  $E = E(e)$  is the unique cell to the boundary of which  $e$  belongs and such that  $\mathbf{n}_E^e \cdot \mathbf{n}^e = 1$ . Due to the definition of  $V_E^e$ , on each edge  $e$  the quantity  $(V \cdot \mathbf{n}_E^e - V_E^e)$  is orthogonal to constants. Therefore we can write

$$\mathsf{T}_{2,1} = \sum_{e \in \mathcal{E}_h} [[q_h]]_e \int_e (V \cdot \mathbf{n}_E^e - V_E^e) (p - \bar{p}_e) \quad (3.61)$$

where  $\bar{p}_e$  is the average of  $p$  on  $e$ . Applying the Hölder inequality to each face's term, the interpolation estimates for the face's velocity, the Cauchy-Schwarz inequality, and the equivalence between norms  $\|\cdot\|_{1,h}$  and  $\|\cdot\|_{1,\mathcal{D}_h}$  give:

$$\begin{aligned} \mathsf{T}_{2,1} &\lesssim \sum_{e \in \mathcal{E}_h} |[q_h]_e| h_e^{\frac{d-1}{2}} \|(V \cdot \mathbf{n}_E^e - V_E^e)(p - \bar{p}_e)\|_{L^2(e)} \\ &\lesssim \sum_{e \in \mathcal{E}_h} h_e^{\frac{d-2}{2}} |[q_h]_e| h_e^{\frac{3}{2}} |V|_{W^{1,\infty}(\Omega)} \|p - \bar{p}_e\|_{L^2(e)} \\ &\lesssim |V|_{W^{1,\infty}(\Omega)} \left( \sum_{e \in \mathcal{E}_h} h_e^{d-2} |[q_h]_e|^2 \right)^{1/2} \left( \sum_{e \in \mathcal{E}_h} h_e^3 \|p - \bar{p}_e\|_{L^2(e)}^2 \right)^{1/2} \\ &\lesssim |V|_{W^{1,\infty}(\Omega)} \|q_h\|_{1,\mathcal{D}_h} \left( \sum_{e \in \mathcal{E}_h} h_e^5 \|\nabla p\|_{L^2(e)}^2 \right)^{1/2}, \end{aligned} \quad (3.62)$$

where in the last line we also used a standard approximation result. Now, the Agmon inequality for  $\nabla p$ , c.f. (M3) with  $\phi = \nabla p$ , and the fact that for  $e \subseteq \partial E \cap \partial E'$  there holds that  $h_e \leq \max(h_E, h_{E'})$  imply that

$$\sum_{e \in \mathcal{E}_h} h_e^5 \|\nabla p\|_{L^2(e)}^2 \lesssim \sum_{E \in \Omega_h} \sum_{e \in \partial E} h_E^5 \|\nabla p\|_{L^2(e)}^2 \lesssim \sum_{E \in \Omega_h} h_E^5 \left( h_E^{-1} |p|_{H^1(E)}^2 + h_E |p|_{H^2(E)}^2 \right).$$

The bound for  $\mathsf{T}_{2,1}$  readily follows recalling (3.45)

$$\mathsf{T}_{2,1} \lesssim h^2 \|p\|_{H^1(\Omega)} + h^3 \|p\|_{H^2(\Omega)} \lesssim h^2 \|p\|_{H^2(\Omega)}, \quad (3.63)$$

where we included the data factor  $|V|_{W^{1,\infty}(\Omega)}$  in the inequality's constant. As a byproduct we observe here that, looking at (3.63) and at the estimate (3.68), it becomes clear that the error coming from the approximation of the datum  $V$  is a higher order term.

Now, let us search for an upper bound for  $\mathsf{T}_{2,2} + \mathsf{T}_3$ . First, note that, since  $V^I$  and  $F^I$  are conservative fields,

$$\sum_{E \in \Omega_h} \sum_{e \in \partial E} q_e \int_e V_E^e p = \sum_{e \in \mathcal{E}_h} q_e (V_E^e + V_{E'}^e) \int_e p = 0 \quad \text{and} \quad \sum_{E \in \Omega_h} \sum_{e \in \partial E} |e| (F^I)_E^e q_e = 0$$

and thus

$$\mathsf{T}_{2,2} + \mathsf{T}_3 = \sum_{E \in \Omega_h} \sum_{e \in \partial E} (q_E - q_e) \int_e \left( V_E^e p - (F_c(p^I))_E^e \right). \quad (3.64)$$

Moreover, Assumption (A2) implies that  $V_E^e = \frac{1}{d_e} (A(d_e V_E^e) + B(d_e V_E^e))$  and therefore, by using definition (2.40) and the triangle inequality, a straightforward calculation gives:

$$\|V_E^e p - (F_c(p^I))_E^e\|_{L^2(e)}^2 \leq \left\| (p - (p^I)_E) \frac{A(d_e V_E^e)}{d_e} \right\|_{L^2(e)}^2 + \left\| (p - (p^I)_{E'}) \frac{B(d_e V_E^e)}{d_e} \right\|_{L^2(e)}^2. \quad (3.65)$$

From (A1) and the definition of  $V_E^e$  it easily follows that  $\max(|A(d_e V_E^e)|, |B(d_e V_E^e)|) \lesssim d_e$ . Then, by using the Agmon inequality of (M3) and the standard first-order interpolation estimate for cell averages, i.e.,  $\|p - (p^I)_E\|_{L^2(E)} \lesssim h_E |p|_{H^1(E)}$ , we have that

$$\begin{aligned} \left\| (p - (p^I)_E) \frac{A(d_e V_E^e)}{d_e} \right\|_{L^2(e)}^2 &\lesssim \|p - (p^I)_E\|_{L^2(e)}^2 \\ &\lesssim h_E^{-1} \|p - (p^I)_E\|_{L^2(E)}^2 + h_E |p - (p^I)_E|_{H^1(E)}^2 \\ &\lesssim h_E |p|_{H^1(E)}^2. \end{aligned} \quad (3.66)$$

A similar inequality can be derived by repeating the same argument for the second term in the right hand side of (3.65) when  $e \in \mathcal{E}_{h,\text{int}}$ , and noting that the second term is zero if  $e$  is a boundary face. Finally, we obtain:

$$\left\| V_E^e p - (F_c(p^I))_E^e \right\|_{L^2(e)}^2 \lesssim h |p|_{H^1(E \cup E')}^2.$$

Therefore, by using a Hölder inequality on the faces and an  $l^2$  Cauchy-Schwarz inequality, from (3.64) we obtain

$$\mathsf{T}_{2,2} + \mathsf{T}_3 \lesssim h^{\frac{1}{2}} \sum_{e \in \mathcal{E}_h} |q_E - q_e| |e|^{\frac{1}{2}} |p|_{H^1(E \cup E')} \lesssim h^{\frac{1}{2}} \|q_h\|_{1, \mathcal{D}_h, \mathcal{E}_h} \left( \sum_{e \in \mathcal{E}_h} h_e |p|_{H^1(E \cup E')}^2 \right)^{1/2}. \quad (3.67)$$

Recalling (3.45) yields

$$\mathsf{T}_{2,2} + \mathsf{T}_3 \lesssim h \|p\|_{H^1(\Omega)}. \quad (3.68)$$

Combining (3.63) and (3.68) we have the bound for  $\mathsf{T}_2 + \mathsf{T}_3$ , and considering also (3.46), (3.47) and (3.57) we conclude the proof. ■

From Theorem 3.11 we get immediately two corollaries that we state without proof.

**Corollary 3.12** *Under the same hypotheses of Theorem 3.11 it holds*

$$\|\tilde{F}_h - \tilde{F}^I\|_{\tilde{X}_h} \lesssim h \|p\|_{H^2(\Omega)}, \quad (3.69)$$

where the total fluxes are defined through  $\tilde{F}^I = -(\Lambda \nabla p + V p)^I$  and  $\tilde{F}_h = F_h + F_c(p_h)$ .

**Corollary 3.13** *Under the same hypotheses of Theorem 3.11 (and applying Lemma 3.5) it holds*

$$\|p^I - p_h\|_{L^r(\Omega)} \lesssim h \|p\|_{H^2(\Omega)}$$

where  $r = \frac{2d}{d-2}$  if  $d > 2$  and  $r < +\infty$  if  $d = 2$ .

**Remark 3.14** *Repeating the same arguments makes it possible to prove a similar error estimate for the hybrid HMM formulation (2.47)-(2.49), which is based on the numerical convection flux (2.45)-(2.46).*

**Remark 3.15** *It must be noted that the proofs in this paper are not uniform with respect to the Peclet number, i.e., the estimates degenerate when the convection becomes dominant. On the other hand, uniform estimates cannot be derived under the general framework considered here, since it comprehends also methods which are not stable in the limit. Nevertheless, the general approach used here can be followed in order to develop uniform error estimates for certain methods. For example, we believe that an uniform error bound can be developed for the upwind scheme starting from a uniform version of the stability results in Proposition 3.2 and Corollary 3.3. A deeper theoretical investigation of the convection dominated case will be the objective of future communications.*

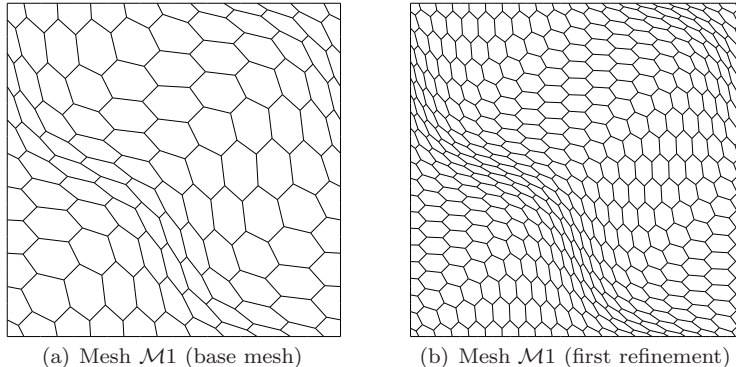


Figure 2: The base mesh and its first refinement of mesh family  $\mathcal{M}1$ ; the mesh construction parameter  $n$  is initially taken equal to 10 and doubled at each refinement step. Details about the mesh characteristics are given in Table 1.

## 4 Numerical experiments

In this section we present a number of examples of problem (2.1)-(2.2), whose solutions are computed over uniform and non-uniform meshes. The performance of these discretization methods is investigated by evaluating the rate of convergence when the meshes are refined and the shock-capturing capability when strong layers develop in the convection-dominated regime.

To this purpose, we consider the sequence of meshes of mesh family  $\mathcal{M}1$  on  $\Omega = ]0, 1[ \times ]0, 1[$ . These meshes are built by remapping the position  $(\xi, \eta)$  of the nodes of an  $n \times n$  uniform grid of quadrilaterals into final positions  $(x, y)$  through

$$x = \xi + (1/10) \sin(2\pi\xi) \sin(2\pi\eta), \quad (4.1)$$

$$y = \eta + (1/10) \sin(2\pi\xi) \sin(2\pi\eta). \quad (4.2)$$

Then, we split each quadrilateral-shaped cell into two triangles, which gives *the primal mesh*, and then we connect the barycenters of adjacent triangular cells by a straight segment. We complete the mesh construction at the domain boundary  $\partial\Omega$  by connecting the barycenters of triangular cells close to  $\partial\Omega$  to the midpoints of boundary edges and these latter to the boundary vertices of the primal mesh. For this mesh family, the base mesh of the refinement process is obtained by setting  $n = 10$ ; refined meshes are generated by doubling this parameter and repeating the construction procedure. The plots of Figure 2 illustrate the base mesh and the first refined mesh of  $\mathcal{M}1$ . Details about the mesh characteristics are reported in Table 1.

The numerical implementation is partially based on P2MESH [15], a C++ public domain library designed to manage data structures of unstructured meshes in the implementation of solvers of partial differential equations. For convenience, we will use the labels listed below to refer to the different instances of the HMM family of schemes considered in our numerical experiments. In each one of these schemes the diffusion term is discretized along the lines described in subsection 2.3 while the numerical treatment of the convection term differs as specified in the item's description:

- HMM-Cnt, two-point centered flux formula;
- HMM-Upw, two-point upwind flux formula;
- HMM-SG, two-point Scharfetter-Gummel formula with local adjustment (2.23);
- HMM-(no stabilization), central mimetic method without any form of stabilization;
- HMM-Jmp, central mimetic method with jump stabilization (2.32);

Table 1: Mesh parameters of the mesh sequence  $\mathcal{M}1$ ;  $r$  is the refinement level (0 refers to the base mesh),  $N_E$  is the number of cells,  $N_e$  is the number of mesh edges,  $N_V$  is the number of mesh vertices.

| $r$ | $N_E$  | $N_e$  | $N_V$  | $h$                   |
|-----|--------|--------|--------|-----------------------|
| 0   | 121    | 400    | 280    | $9.477 \cdot 10^{-2}$ |
| 1   | 441    | 1400   | 960    | $4.843 \cdot 10^{-2}$ |
| 2   | 1681   | 5200   | 3520   | $2.445 \cdot 10^{-2}$ |
| 3   | 6561   | 20000  | 13440  | $1.225 \cdot 10^{-2}$ |
| 4   | 25921  | 78400  | 52480  | $6.130 \cdot 10^{-3}$ |
| 5   | 103041 | 310400 | 207360 | $3.066 \cdot 10^{-3}$ |

## 4.1 Accuracy

In this test case, the forcing term  $f$  in (2.1) and the boundary condition function  $g^D$  in (2.2) are set accordingly to the exact solution

$$p(x, y) = \left( x - e^{\frac{2(x-1)}{\nu}} \right) \left( y^2 - e^{\frac{3(y-1)}{\nu}} \right) \quad (4.3)$$

and  $V = (2, 3)^t$ . We assume that the diffusion tensor  $\Lambda$  is given by the identity matrix scaled by the positive real factor  $\nu$ . By taking  $\nu = 10^{-4}$  the problem is strongly convection-dominated and the solution is characterized by an exponential boundary layer near top and right sides of  $\Omega$ .

We are mainly interested in showing that the shock-capturing capability does not deteriorate too much the convergence behavior where the solution is enough smooth, i.e. away from the boundaries where the layer develops. As pointed out in [16, 27, 42, 43] (to which we also refer the reader interested in the comparison with performance of the mixed-hybrid Finite Element and different kind of Finite Volume schemes on this test case), the errors due to the approximation of the solution gradient in the narrow strip around the boundary where the layer develops are so large that including them in the error measurements would prevent to see any convergence at all. For this reason, we restrict the error measurement to the sub-domain  $[0, 0.95] \times [0, 0.95]$ . Convergence rates are measured by the relative errors

$$\mathcal{E}_{Q_h} = \frac{\|p^I - p_h\|_{Q_h}}{\|p^I\|_{Q_h}} \quad \text{and} \quad \mathcal{E}_{X_h} = \frac{\|\widetilde{F}^I - \widetilde{F}_h\|_{\widehat{X}_h}}{\|\widetilde{F}^I\|_{\widehat{X}_h}}, \quad (4.4)$$

where, in the second error definition, we use the total fluxes  $\widetilde{F}^I$  and  $\widetilde{F}_h$  are defined in Corollary 3.12. Practically speaking, the quantity  $\mathcal{E}_{Q_h}$  is a measure of the approximation error of cell averages and is calculated by using a mesh-dependent  $L^2$ -like norm. On its turn, error  $\mathcal{E}_{X_h}$  compares the edge-based flux  $\widetilde{F}^I$  with the numerical flux  $\widetilde{F}_h$  through the mesh-dependent norm induced in  $X_h$  by the mimetic scalar product.

In Figure 3 we present the log-log plots of the errors  $\mathcal{E}_{Q_h}$  (on the left) and  $\mathcal{E}_{X_h}$  (on the right) versus the characteristic mesh size  $h$ . Herein, we compare the convergence behavior of the various implementation of the HMM schemes considered in this paper. The actual order of accuracy shown by these methods is reflected by the slopes of the experimental error curves, and can be approximately evaluated by comparison with the “theoretical” slopes represented in the bottom-left corner of each plot, c.f., also the caption’s comment. These plots document the *optimal* convergence behavior of all the numerical approximations in the diffusive regime, c.f., the top side plots. When the problem becomes convection-dominated, i.e., for the smallest value of the diffusion coefficient, convergence is still provided for both scalar and flux unknowns by all methods except HMM-(no stabilization).

When we use HMM-SG, HMM-Cnt, and HMM-(no stabilization) in the diffusive regime, a superconvergence effect is visible for the approximation of the scalar and the flux variable. The numerical approximation of the scalar unknown is second-order accurate, while we have  $\mathcal{O}(h^{3/2})$  for the flux variable.



Instead, both HMM-Upw and HMM-Jmp provides a first-order accurate approximation for both  $p$  and  $F$ . It is also worth noting that the error curves of HMM-SG and HMM-Cnt almost coincide. Moreover, the errors from HMM-(no stabilization) are a little bit smaller than those obtained by the centered schemes of Finite Volume type. Instead, in this test case scheme HMM-Upw gives better results than HMM-Jmp.

In the convection-dominated case, i.e., for  $\nu = 10^{-4}$ , the central approximation HMM-(no stabilization) is not at all convergent on the meshes considered by  $\mathcal{M}1$ ; instead, HMM-Cnt is convergent, but the numerical solution (not shown in the paper) is affected by large amplitude oscillations that almost completely destroy the solution's profile. This fact is consistent with the error curves displayed in Figure 3. The numerical approximation of the scalar and flux variable provided by the methods HMM-Upw and HMM-SG is linearly convergent, while the one provided by HMM-Jmp seems to converge at a rate proportional to  $\mathcal{O}(h^{1/2})$ , even if this last effect might be due to an insufficient mesh resolution.

**Remark 4.1** *As noticed at the end of Section 2.4.1, in the convection-dominated regime the  $A$  function given by (2.23) is numerically nearly indistinguishable from the upwind functions  $A_{\text{up}}$ . Figures 3 and 4 confirm this. On the other hand, in the diffusion regime, the modified Scharfetter-Gummel scheme has better convergence properties than the upwind scheme. This behavior is a very interesting characteristic of the choice (2.23) when the convection term is discretized by (2.40): it automatically adjusts to either provide a good order of convergence in the diffusive regime, or enough numerical diffusion to stabilize the calculation in the convection-dominated regime. Note that if one takes  $A$  and  $B$  satisfying (AB1)-(AB2) and (AB3-s) and such that  $A(s) \sim s$  as  $s \rightarrow +\infty$ ,  $A(s)$  has a finite limit as  $s \rightarrow -\infty$  and  $A(s)$  is regular around  $s = 0$ , then a scheme using such functions modified in the same way as (2.23) is expected to show the same kind of behavior (this has been numerically tested on several choices of such functions). Note also that this approach does not hold for the upwind scheme since  $A_{\text{up}}(s)$  is not regular at  $s = 0$ .*

## 4.2 Shock-capturing behavior

Shock-capturing behavior is investigated by solving (2.1)-(2.2) in the convection-dominated regime. The exact solution may be characterized by boundary layers of exponential and parabolic types and is approximated on the sequence of meshes of  $\mathcal{M}1$ . The numerical solution is plotted at mesh vertices. Vertex values are obtained by interpolating the approximate cell averages provided by the scheme.

### 4.2.1 Exponential boundary layers

We experimentally investigate how these methods approximate a solution with an exponential boundary layer, which forms on those sides of the domain boundary where  $V$  points outward. To this purpose, we solve problem (2.1)-(2.2) with the same data of the accuracy benchmark test in the convection-dominated regime, i.e., for  $\nu = 10^{-4}$ .

In Figure 4, we compare the numerical solutions produced by the following implementations: HMM-SG, HMM-Upw, HMM-(no stabilization) and HMM-Jmp.

In plots (a)–(b) the non-oscillatory solution produced by schemes HMM-Upw and HMM-SG is displayed. Instead, from plot (c) it is evident that when calculation is performed using the HMM-(no stabilization) without any stabilization the numerical solution suffers of severe oscillations. These oscillations disappear when we introduce a stabilizing term in the divergence equation, which is based on the solution's jump at mesh edges. However, a great numerical diffusion is introduced by this form of upwinding and the resolution of the boundary layer is poor and generally worst than that obtained through the other HMM implementations.

### 4.2.2 Exponential and parabolic boundary layers

On  $\Omega = ]0, 1[ \times ]0, 1[$ , we numerically solve (2.1)-(2.2) with Dirichlet boundary condition

$$p(x, 0) = (1 - x)^3, \quad p(x, 1) = (1 - x)^2, \quad p(0, y) = 1, \quad p(1, y) = 0,$$

and  $V = (1, 0)^T$  in the convection-dominated regime for  $\nu = 10^{-4}$ . The solution has an exponential boundary layer at the side  $x = 1$  and two parabolic boundary layers at  $y = 0$  and  $y = 1$ . Figure 5 shows the numerical results obtained from calculations using HMM-SG, HMM-Upw, HMM-(no stabilization), and HMM-Jmp. The behavior is similar to the behavior documented in the previous subsection for the case of the single exponential layer.

### 4.3 Strongly anisotropic heterogeneous and convection-dominated case

In this third example we consider the test case proposed in [34], where problem (2.1)-(2.2) is solved for a strongly anisotropic and heterogeneous diffusion tensor and a rotating convection field. A zero-th order term proportional to  $p$  is also present in the model's equations; its discretization is straightforward (see e.g. [24]). The domain  $\Omega = ]0, 1[ \times ]0, 1[$  is split into four subdomains  $\Omega_1 = ]0, 2/3[ \times ]0, 2/3[$ ,  $\Omega_2 = ]0, 2/3[ \times ]2/3, 1[$ ,  $\Omega_3 = ]2/3, 1[ \times ]2/3, 1[$ ,  $\Omega_4 = ]2/3, 1[ \times ]0, 2/3[$ . The diffusion tensor is diagonal in each sub-domain and is characterized by a very small value along one principal direction:

$$\Lambda = \begin{pmatrix} 10^{-6} & 0 \\ 0 & 1 \end{pmatrix} \quad \text{in } \Omega_1 \text{ and } \Omega_3,$$

and

$$\Lambda = \begin{pmatrix} 1 & 0 \\ 0 & 10^{-6} \end{pmatrix} \quad \text{in } \Omega_2 \text{ and } \Omega_4.$$

Note that the directions along which diffusion is small are interchanged for adjacent subdomains. Convection is given by the clockwise rotating solenoidal field  $V(x, y) = 40(x(2y-1)(x-1), -y(2x-1)(y-1))$  and the right-hand side is a gaussian bump positioned at distance 0.35 of the domain center,  $f(x, y) = 10^{-2} \exp(-(r - 0.35)^2/0.005)$  with  $r^2 = (x - 0.5)^2 + (y - 0.5)^2$ . This problem is convection-dominated, thus requiring some sort of upwinding in the numerical treatment of the convection term. Moreover, the exact solution is continuous, but internal layers form near the interfaces that separate the subdomains due to the small diffusion value in the switching directions. The strong solution gradients cannot be resolved by the attainable grid sizes and the numerical approximations are expected to be discontinuous at the internal interfaces.

In the test cases presented in the previous subsections there was no significant difference between the numerical approximations provided by the cell-based version of the upwind scheme HMM-Upw, c.f. (2.41)-(2.42), and its edge-based version, c.f. (2.47)-(2.49). This is no longer the case herein, as illustrated in Figure 6. The calculations, which use the two alternative versions of the HMM-Upw scheme, are performed on a grid obtained by a  $30 \times 30$  periodic reproduction of the pattern reported in Figure 6-(a). Since the *exact* solution of this problem is unknown, a reference solution is calculated, for comparison's sake, on a very fine cartesian grid. The reference solution is displayed in Figure 6-(b). The numerical solution provided by the cell-based upwind scheme is shown in Figure 6-(c) and is clearly affected by spurious oscillations. Instead, this undesirable effect is almost completely absent in the numerical solution provided by the edge-based upwind scheme, which is shown in Figure 6-(d).

It is also worth mentioning the behavior of these two different implementations of the HMM-Upw scheme as far as minimum and maximum principles are concerned. To this purpose, we recall that the numerical solutions obtained by first-order upwind two-point Finite Volume schemes in convection-dominated problems are characterized by numerical properties like positivity, monotonicity, etc. A thorough inspection of our numerical results reveals that both cell-based and edge-based schemes respects the minimum value, which is zero for the reference solution, and provides  $6.6 \times 10^{-4}$  and  $6.9 \times 10^{-4}$ , respectively, for the maximum value against a reference value of approximately  $6.7 \times 10^{-4}$ . Nonetheless, we noticed a minimum value of approximately  $-1.1 \times 10^{-5}$ , which corresponds to a numerical undershoot of around 1.6%, when we applied the cell-based scheme on a different mesh given by splitting every other rectangular cell of a  $120 \times 60$  regular partition of  $\Omega$  in two sub-triangles. On this latter mesh the edge-based upwind scheme was still seen to respect the zero minimum value. We do not show the other solution plots for these latter

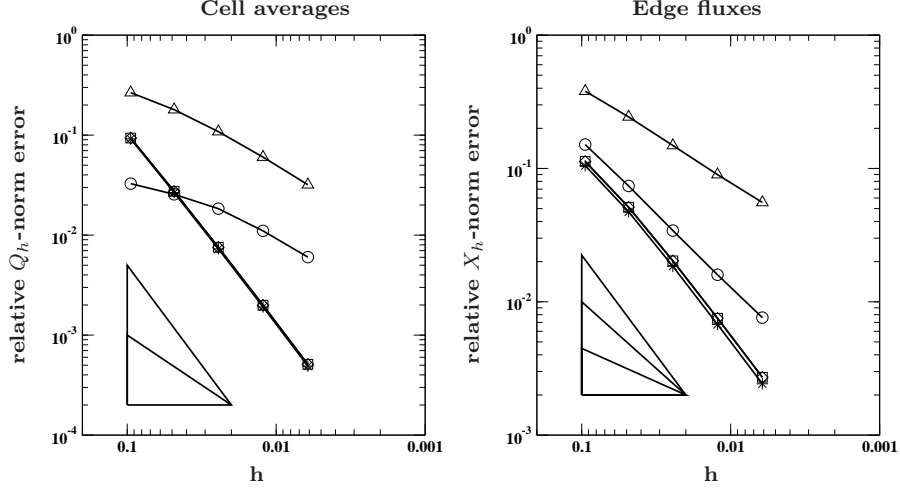
calculations because their behaviors are very similar to those of the solutions shown in Figure 6. From these qualitative comparisons we deduce that the edge-based upwind scheme may be more stable and accurate than the cell-based upwind scheme. We also remark that the edge-based upwind scheme has the advantage of being fully hybridizable, thus leading to a linear system in the edge unknowns through local variable eliminations like, for example, in the static condensation of Mixed Finite Elements. For these reasons, the edge-based upwind scheme may be preferable when dealing with stiff problems on coarse meshes.

**Remark 4.2** *As a final comment we observe that, in all the developed tests, the schemes HMM-(no stabilization) and HMM-Jmp, which satisfy only (AB3-w), do not show particular pathologies for coarse meshes. Therefore, at least on the basis of the presented tests, the  $h$ -small-enough condition appearing in Theorem 3.11 does not seem to pose a true limitation in practice.*

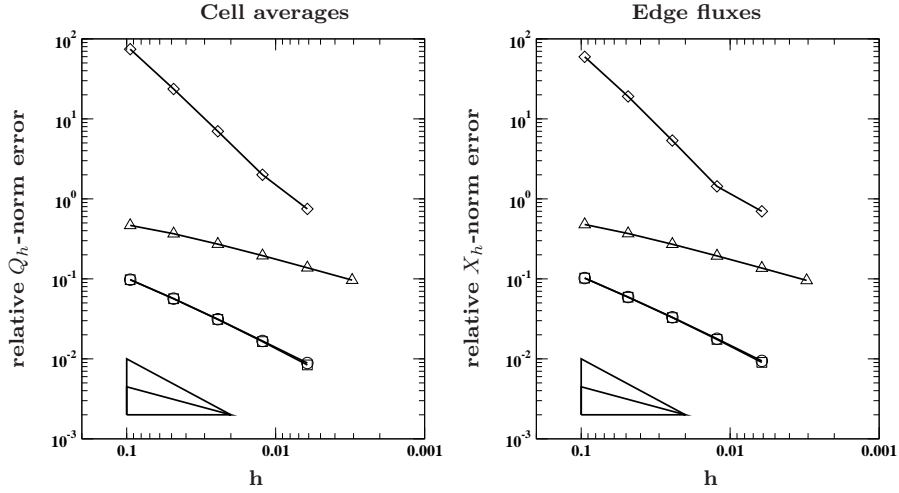
## 5 Conclusions

We presented a new family of methods for the numerical approximation to the solution of the steady convection-diffusion equation. These methods, which are referred to as *Hybrid Mimetic Mixed* methods (HMM), are based on a unified formulation for the Hybrid Finite Volume method, the Mixed Finite Volume method and the Mimetic Finite Difference method, and differ mainly in the approximation of the convection term. In particular, we considered centered, upwind, weighted and locally scaled Scharfetter-Gummel type discretizations, for which we provided a full proof of convergence under very general regularity conditions of the solution field, and derived an error estimate when the scalar solution is in  $H^2(\Omega)$ .

In the last part of the paper, we numerically compared the performance of these schemes on a set of test cases selected from the literature in both diffusion and convection-dominated regimes. As expected, the methods, including a centered-type discretization of the convective term, showed a better behavior in the test cases dominated by diffusion, exhibiting a superconvergence in the approximation of both scalar and vector variables. On the other hand, such schemes showed a strong loss of convergence rate in the convection dominated tests, while on that same tests the methods with upwinding or stabilization exhibited a better behavior. Finally, we showed a test with strong anisotropy and jumps in the coefficients. The results seem to suggest that the hybridized formulation gives more stable results for this kind of problems.



Diffusive regime:  $V = (2, 3)^T$ ,  $\nu = 1$ ; slopes are  $h^2$  and  $h$  on the left plot,  $h^{3/2}$ ,  $h$ ,  $h^{1/2}$  on the right plot.



Convection-Dominated regime:  $V = (2, 3)^T$ ,  $\nu = 10^{-4}$ ; slopes are  $h$  and  $h^{1/2}$  on both plots.

*Symbols:*  $\circ$  HMM-Upw,  $\square$  HMM-SG,  $\diamond$  HMM-Cnt,  $\triangle$  HMM-Jmp,  $\star$  HMM-(no stabilization).

Figure 3: Test case 1: error curves for the numerical approximation of an exact solution that is smooth in the diffusive regime (top) and shows an exponential boundary layer on right and top sides of the computational domain in the convection-dominated regime (bottom). Approximation errors are measured on the reduced domain  $[0, 0.95] \times [0, 0.95]$ , i.e. away from the critical region where the layer may develop. All calculations are performed on the mesh sequence  $\mathcal{M}1$ .

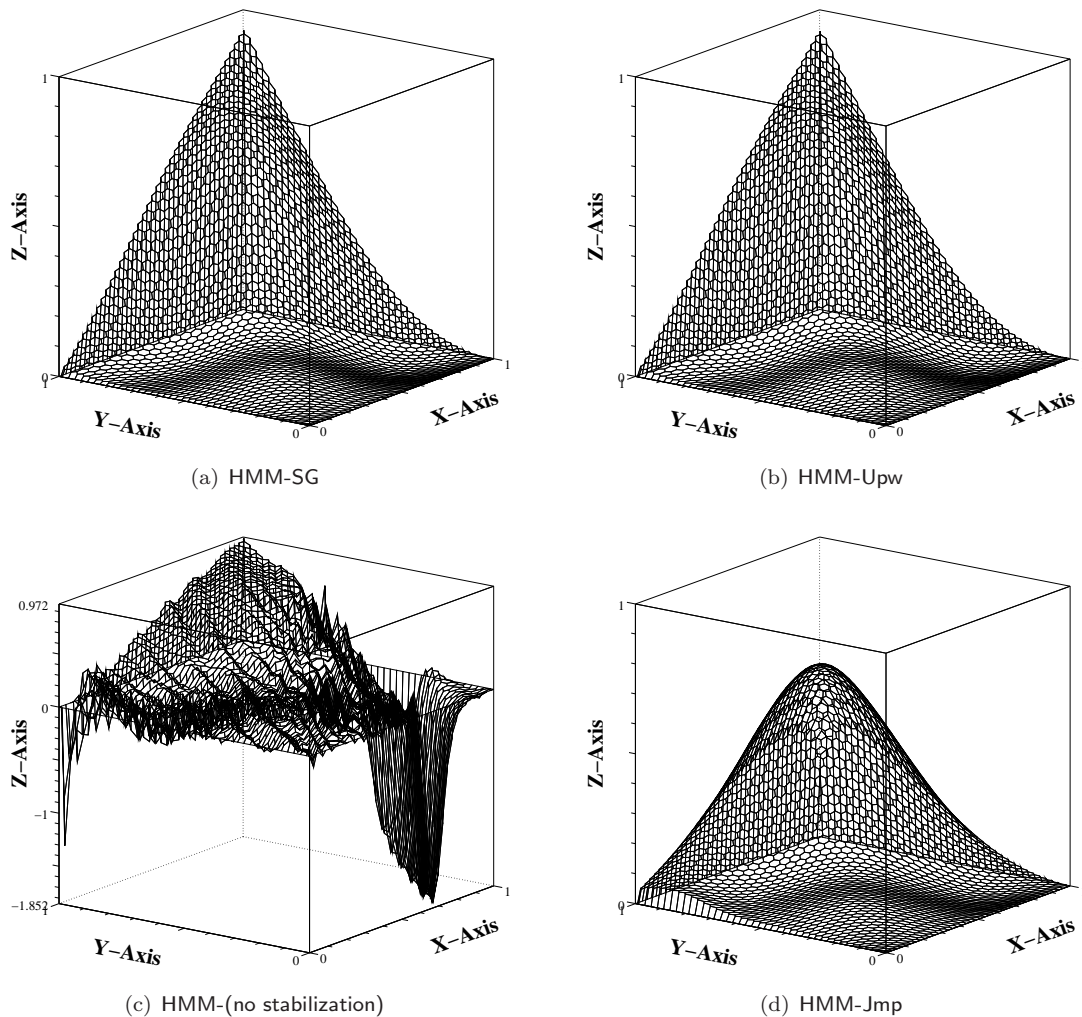


Figure 4: Shock-capturing test case: the exact solution has an exponential boundary layer on right and top sides of the computational domain. Calculations are performed on the second mesh of mesh family  $\mathcal{M}1$  by taking constant velocity field  $V = (2, 3)^T$  and  $\nu = 10^{-4}$ . Numerical solution is displayed at mesh vertices through linear interpolation. Severe oscillations are visible in plot (c) when we use scheme HMM-(no stabilization), i.e, the central mimetic discretization without any stabilization (note the different scale along  $Z$ ). This phenomenon disappears in plot (d) when jump stabilization is turned on using scheme HMM-Jmp.

## References

- [1] I. Aavatsmark, T. Barkve, O. Boe, and T. Mannseth. Discretization on unstructured grids for inhomogeneous, anisotropic media. Part I: Derivation of the methods. *SIAM J. Sci. Comput.*, 19(5):1700–1716, 1998.
- [2] I. Aavatsmark, T. Barkve, O. Boe, and T. Mannseth. Discretization on unstructured grids for inhomogeneous, anisotropic media. Part II: Discussion and numerical results. *SIAM J. Sci. Comput.*, 19(5):1717–1736, 1998.

- [3] T. Arbogast, C. N. Dawson, P. T. Keenan, M. F. Wheeler, and I. Yotov. Enhanced cell-centered finite differences for elliptic equations on general geometry. *SIAM J. Sci. Comput.*, 19(2):404–425, 1998.
- [4] D.N. Arnold and F. Brezzi. Mixed and nonconforming finite element methods. Implementation, post processing and error estimates. *Math. Mod. Numer. Anal.*, 19:7–32, 1985.
- [5] D.N. Arnold, F. Brezzi, B. Cockburn, and L.D. Marini. Unified analysis of discontinuous Galerkin methods for elliptic problems. *Siam J. Numer. Anal.*, 39:1749–1779, 2002.
- [6] J. Baranger, J.-F. Maitre, and F. Oudin. Connection between finite volume and mixed finite element methods. *RAIRO Modl. Math. Anal. Numr.*, 30:445–465, 1996.
- [7] L. Beirão da Veiga. A mimetic finite difference method for linear elasticity. accepted for publication in *Math. Mod. Numer. Anal.*
- [8] L. Beirão da Veiga. A residual based error estimator for the mimetic finite difference method. *Numer. Math.*, 108(3):387–406, 2008.
- [9] L. Beirão da Veiga, V. Gyrya, K. Lipnikov, and G. Manzini. Mimetic finite difference method for the Stokes problem on polygonal meshes. *J. Comput. Phys.*, 228(19):7215–7232, 2009.
- [10] L. Beirão da Veiga, K. Lipnikov, and G. Manzini. The mimetic finite difference method for the steady Stokes problem on polyhedral meshes. Technical Report 6PV09/5/0, IMATI-CNR, 2007.
- [11] L. Beirão da Veiga, K. Lipnikov, and G. Manzini. Convergence analysis of the high-order mimetic finite difference method. *Numer. Math.*, 113(3):325–356, 2009.
- [12] L. Beirão da Veiga and G. Manzini. An a posteriori error estimator for the mimetic finite difference approximation of elliptic problems. *Int. J. for Numer. Meth. in Engrn.*, 76(11):1696–1723, 2008.
- [13] L. Beirão da Veiga and G. Manzini. A higher-order formulation of the mimetic finite difference method. *SIAM J. Sci. Comput.*, 31(1):732–760, 2008.
- [14] M. Berndt, K. Lipnikov, J. D. Moulton, and M. Shashkov. Convergence of mimetic finite difference discretizations of the diffusion equation. *East-West J. Numer. Math.*, 9:253–284, 2001.
- [15] E. Bertolazzi and G. Manzini. Algorithm 817 P2MESH: generic object-oriented interface between 2-D unstructured meshes and FEM/FVM-based PDE solvers. *ACM Trans. Math. Softw.*, 28(1):101–132, 2002.
- [16] E. Bertolazzi and G. Manzini. A cell-centered second-order accurate finite volume method for convection-diffusion problems on unstructured meshes. *Math. Models Methods Appl. Sci.*, 8:1235–1260, 2004.
- [17] S. Brenner and L. Scott. *The mathematical theory of finite element methods*. Springer-Verlag, Berlin/Heidelberg, 1994.
- [18] F. Brezzi, A. Buffa, and K. Lipnikov. Mimetic finite differences for elliptic problems. *Math. Mod. Numer. Anal.*, 43:277–295, 2009.
- [19] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, New York, 1991.
- [20] F. Brezzi, K. Lipnikov, and M. Shashkov. Convergence of the mimetic finite difference method for diffusion problems on polyhedral meshes. *SIAM J. Numer. Anal.*, 43(5):1872–1896, 2005.

- [21] F. Brezzi, K. Lipnikov, M Shashkov, and V. Simoncini. A new discretization methodology for diffusion problems on generalized polyhedral meshes. *Comput. Methods Appl. Mech. Engrg.*, 196:3682–3692, 2007.
- [22] F. Brezzi, K. Lipnikov, and V. Simoncini. A family of mimetic finite difference methods on polygonal and polyhedral meshes. *Math. Models Methods Appl. Sci.*, 15(10):1533–1551, 2005.
- [23] A. Cangiani and G. Manzini. Flux reconstruction and pressure post-processing in mimetic finite difference methods. *Comput. Methods Appl. Mech. Engrg.*, 197/9-12:933–945, 2008.
- [24] A. Cangiani, G. Manzini, and A. Russo. Convergence analysis of a mimetic finite difference method for general second-order elliptic problems. *SIAM J. Numer. Anal.*, 47(4):2612–2637, 2009.
- [25] C. Chainais-Hillairet and J. Droniou. Convergence analysis of a mixed finite volume scheme for an elliptic-parabolic system modeling miscible fluid flows in porous media. *SIAM J. Numer. Anal.*, 45(5):2228–2258, 2007.
- [26] C. Chainais-Hillairet and J. Droniou. Finite volume schemes for non-coercive elliptic problems with Neumann boundary conditions, 2009. to appear in *IMAJNA*.
- [27] Yves Coudière and Gianmarco Manzini. The discrete duality finite volume method for convection-diffusion problems. *SIAM Journal on Numerical Analysis*, 47(6):4163–4192, 2010.
- [28] C. Dawson and V. Aizinger. Upwind-mixed methods for transport equations. *Comput. Geosci.*, 3:93–110, 1999.
- [29] J. Jr. Douglas and J. E. Roberts. Mixed finite element methods for second order elliptic problems. *Mat. Apl. Comput.*, 1(1):91–103, 1982.
- [30] J. Jr. Douglas and J. E. Roberts. Global estimates for mixed methods for second order elliptic equations. *Math. Comp.*, 44:39–52, 1985.
- [31] J. Droniou and R. Eymard. A mixed finite volume scheme for anisotropic diffusion problems on any grid. *Numer. Math.*, 105(1):35–71, 2006.
- [32] J. Droniou and R. Eymard. Study of the mixed finite volume method for Stokes and Navier-Stokes equations. *Numer. Meth. P. D. E.*, 25(1):137–171, 2009.
- [33] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. A unified approach to mimetic finite difference, hybrid finite volume and mixed finite volume methods. *Math. Models Methods Appl. Sci. (M3AS)*, 20(2):265–295, 2010. DOI No: 10.1142/S0218202510004222.
- [34] A. Ern, A.F. Stephansen, and P. Zunino. A discontinuous Galerkin method with weighted averages for advection-diffusion equations with locally small and anisotropic diffusivity. *IMA J. Numer. Anal.*, 29(2):235–256, 2009.
- [35] R. Eymard, T. Gallouët, and R. Herbin. The finite volume method. In P. Ciarlet and J.L. Lions, editors, *Techniques of Scientific Computing, Part III, Handbook for Numerical Analysis*, pages 715–1022. North Holland, 2000.
- [36] R. Eymard, T. Gallouët, and R. Herbin. Discretization of heterogeneous and anisotropic diffusion problems on general non-conforming meshes, SUSHI: a scheme using stabilisation and hybrid interfaces, 2009. To appear in *IMAJNA*.
- [37] J. Hyman, J. Morel, M. Shashkov, and S. Steinberg. Mimetic finite difference methods for the diffusion equation. *Comput. Geosci.*, 6:333–352, 2002.

- [38] J. Jaffre. Decentrage et elements finis mixtes pour les equations de diffusion-convection. *Calcolo*, 21:171–197, 1984.
- [39] J. Jaffre and J. E. Roberts. Upstream weighting and mixed finite elements in the simulation of miscible displacements. *RAIRO Modl. Math. Anal. Numr.*, 19(3):443–460, 1985.
- [40] Y. Kuznetsov, K. Lipnikov, and M. Shashkov. The mimetic finite difference method on polygonal meshes for diffusion-type problems. *Comput.Geosci.*, 8:301–324, 2005.
- [41] K. Lipnikov, M. Shashkov, and I. Yotov. Local flux mimetic finite difference methods. *Numer. Math.*, 112:115–152, 2009.
- [42] Gianmarco Manzini and Alessandro Russo. A finite volume method for advection-diffusion problems in convection-dominated regimes. *Computer Methods in Applied Mechanics and Engineering*, 197(13-16):1242 – 1261, 2008.
- [43] G. Rapin and G. Lube. A stabilized scheme for the lagrange multiplier method for advection-diffusion equations. *Math. Models Methods Appl. Sci.*, 14:1035–1060, 2004.
- [44] B. Riviere. *Discontinuous Galerkin methods for solving Elliptic and Parabolic Equations: Theory and Implementation*. SIAM, 2008.
- [45] T.F. Russell and M.F. Wheeler. Finite element and finite difference methods for continuous flows in porous media. In R.E. Ewing, editor, *The Mathematics of Reservoir Simulation*, pages 35–106, Philadelphia, 1983. SIAM.
- [46] D. L. Scharfetter and H. K. Gummel. Large signal analysis of a silicon read diode. *IEEE Trans. on Elec. Dev.*, 16:64–77, 1969.
- [47] M. Vohralik. Equivalence between lowest-order mixed finite element and multi-point finite volume methods on simplicial meshes. *M2AN Math. Model. Numer. Anal.*, 40(2):367–391, 2006.
- [48] M. F. Wheeler and I. Yotov. A multipoint flux mixed finite element method. *SIAM J. Numer. Anal.*, 44(5):2082–2106, 2006.
- [49] A. Younes, P. Ackerer, and G. Chavent. From mixed finite elements to finite volumes for elliptic pdes in two and three dimensions. *Internat. J. Numer. Methods Engrg.*, 59(3):365–388, 2004.



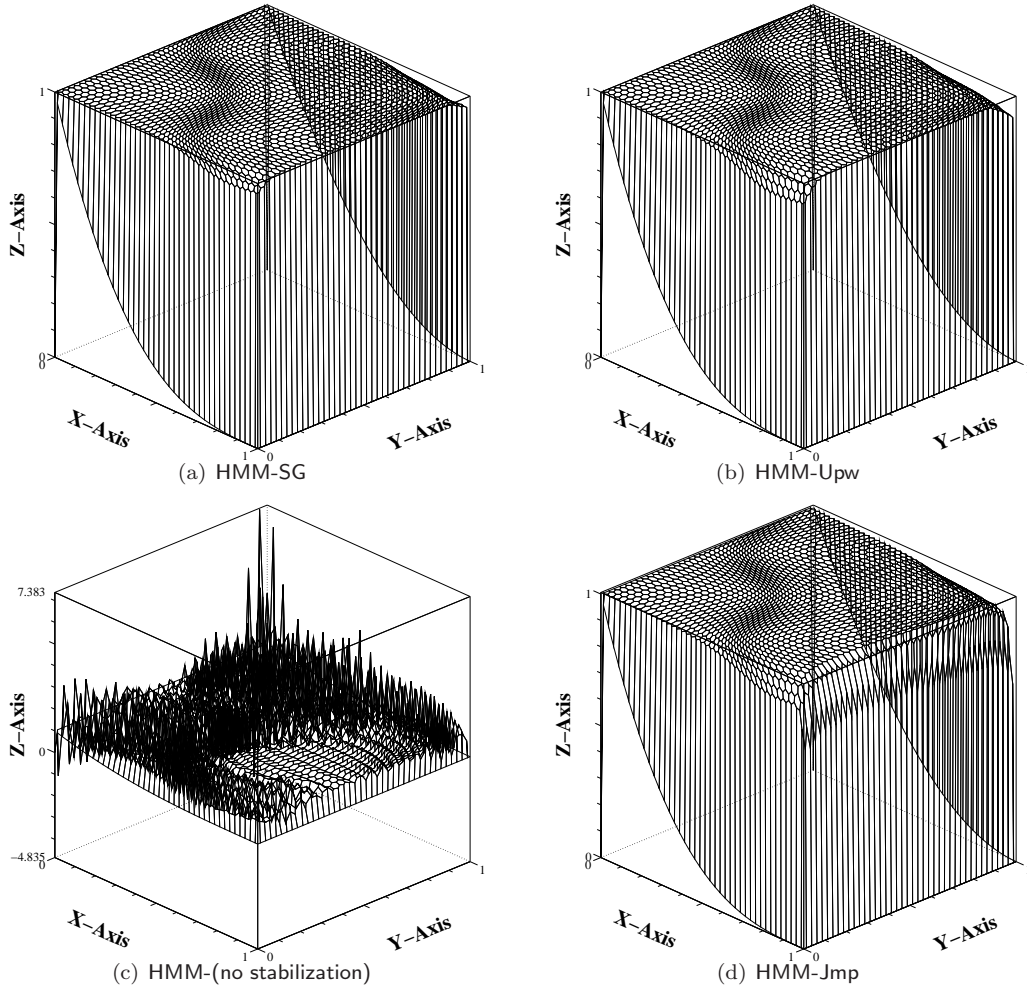


Figure 5: Shock-capturing test case: the exact solution has an exponential boundary layer on the right side and two parabolic layers on top and bottom side of the computational domain. Calculations are performed on the second mesh of mesh family  $\mathcal{M}1$  by taking the constant velocity field  $V = (1, 0)^T$  and  $\nu = 10^{-4}$ . Numerical solution is displayed at mesh vertices through linear interpolation. Severe oscillations are visible in plot (c) when we use scheme HMM-(no stabilization), i.e, the central mimetic discretization without any stabilization (note the different scale along  $Z$ ). This phenomenon disappears in plot (d) when jump stabilization is turned on using scheme HMM-Jmp.

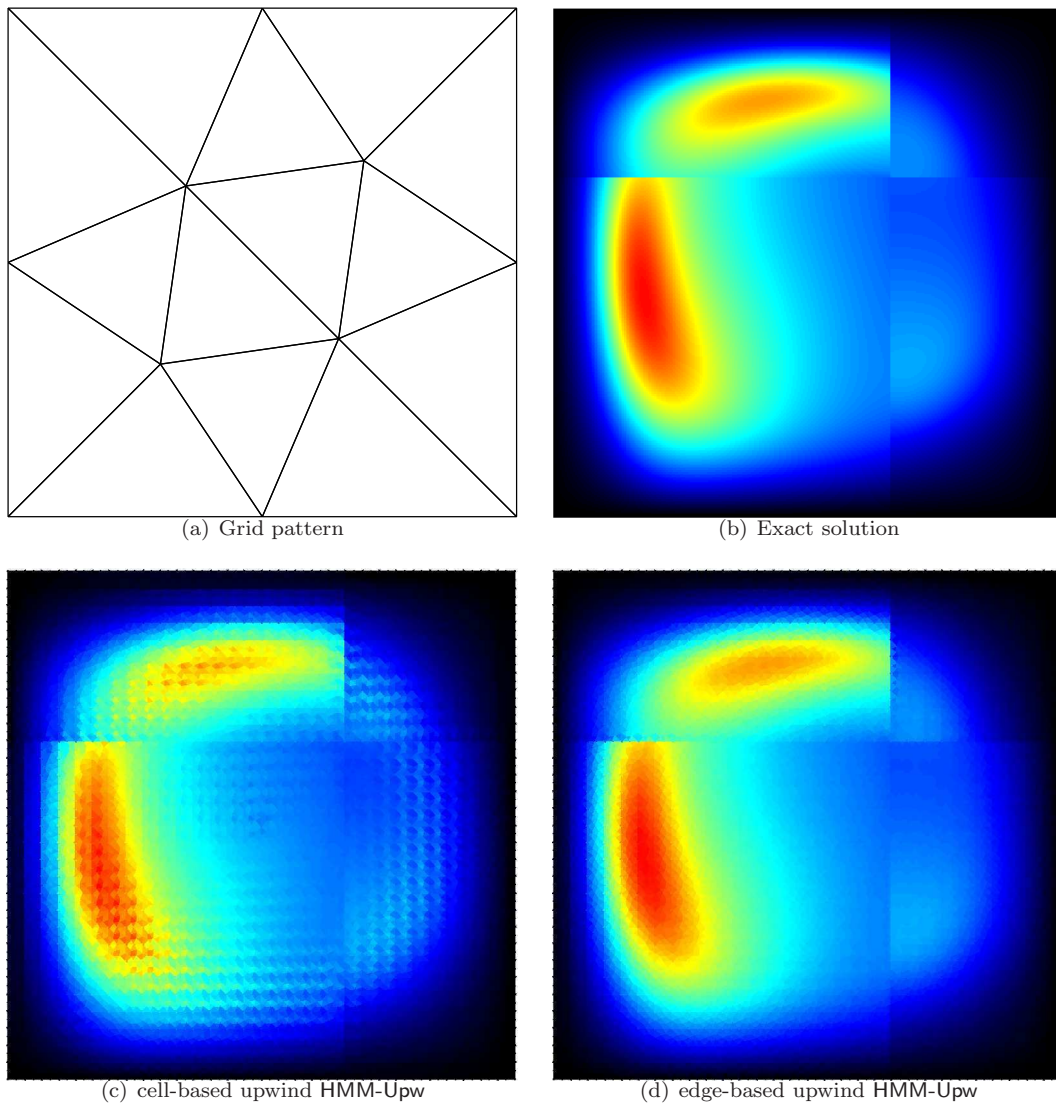


Figure 6: Strongly anisotropic heterogeneous and convection-dominated test case: on a coarse mesh, the cell-based upwinding of the convection provokes spurious oscillations, which are completely absent in the edge-based upwinding discretization.