

# Discretisation of heterogeneous and anisotropic diffusion problems on general nonconforming meshes SUSHI: a scheme using stabilisation and hybrid interfaces <sup>1</sup>

R. Eymard<sup>2</sup>, T. Gallouët<sup>3</sup> and R. Herbin<sup>4</sup>

**Abstract:** A symmetric discretisation scheme for heterogeneous anisotropic diffusion problems on general meshes is developed and studied. The unknowns of this scheme are the values at the centre of the control volumes and at some internal interfaces which may for instance be chosen at the diffusion tensor discontinuities. The scheme is therefore completely cell-centred if no edge unknown is kept. It is shown to be accurate on several numerical examples. Convergence of the approximate solution to the continuous solution is proved for general (possibly discontinuous) tensors, general (possibly nonconforming) meshes, and with no regularity assumption on the solution. An error estimate is then deduced under suitable regularity assumptions on the solution.

**Keywords :** Heterogeneous anisotropic diffusion, nonconforming grids, finite volume schemes

## 1 Introduction

Anisotropic heterogeneous diffusion problems arise in a wide range of scientific fields such as hydrogeology, oil reservoir simulation, plasma physics, semiconductor modelling, biology, etc.. When implementing numerical methods for this kind of problem, one needs to find an approximation of  $u$ , weak solution to the following equation:

$$-\operatorname{div}(\Lambda(\mathbf{x})\nabla u) = f \text{ in } \Omega, \quad (1)$$

with boundary condition

$$u = 0 \text{ on } \partial\Omega, \quad (2)$$

where we denote by  $\partial\Omega = \overline{\Omega} \setminus \Omega$  the boundary of the domain  $\Omega$ , under the following assumptions:

$$\Omega \text{ is an open bounded connected polyhedral subset of } \mathbb{R}^d, \quad d \in \mathbb{N} \setminus \{0\}, \quad (3)$$

$$\Lambda \text{ is a measurable function from } \Omega \text{ to } \mathcal{M}_d(\mathbb{R}), \quad (4)$$

where we denote by  $\mathcal{M}_d(\mathbb{R})$  the set of  $d \times d$  matrices, such that for a.e.  $\mathbf{x} \in \Omega$ ,  $\Lambda(\mathbf{x})$  is symmetric, and such that the set of its eigenvalues is included in  $[\underline{\lambda}, \overline{\lambda}]$ , with  $\underline{\lambda}$  and  $\overline{\lambda} \in \mathbb{R}$  satisfying  $0 < \underline{\lambda} \leq \overline{\lambda}$ , and

$$f \in L^2(\Omega). \quad (5)$$

Under these hypotheses, the weak solution of (1)–(2) is the unique function  $u$  satisfying:

$$\begin{cases} u \in H_0^1(\Omega), \\ \int_{\Omega} \Lambda(\mathbf{x})\nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x})d\mathbf{x} = \int_{\Omega} f(\mathbf{x})v(\mathbf{x})d\mathbf{x} \quad \forall v \in H_0^1(\Omega). \end{cases} \quad (6)$$

Usual discretisation schemes for Problem (6) include finite difference, finite element or finite volume methods. Finite volume methods are actually very popular in oilreservoir engineering, a probable reason being that complex coupled physical phenomena may be discretised on the same grids. The well-known five-point scheme on rectangles (see e.g. [29]) and four-point scheme on triangles [23] are

<sup>1</sup>This work was supported by Groupement MOMAS, CNRS/PACEN

<sup>2</sup>Université Paris-Est, France, Robert.Eymard@univ-mlv.fr

<sup>3</sup>Université de Provence, France, Thierry.Gallouet@cmi.univ-mrs.fr

<sup>4</sup>Université de Provence, France, Raphaelle.Herbin@cmi.univ-mrs.fr

not easily adapted to heterogeneous anisotropic diffusion operators [24]. A scheme with an enlarged stencil, which handles anisotropy on meshes satisfying an orthogonality property, was proposed and analysed in [18]. Another problem that has to be faced in several fields of applications (such as hydrogeology and oil reservoir engineering) is the fact that the discretisation meshes are imposed by engineering and computing considerations; therefore, we have to deal with distorted and possibly nonconforming meshes.

A huge literature exists in the engineering setting, so we shall not try to be exhaustive. Let us nevertheless mention the finite volume schemes using the well-known multipoint flux approximation [1, 2, 3]. These schemes involve the reconstruction of the gradient in order to evaluate the fluxes, which is also the case in [13, 28]. Among other approaches let us cite [22], which uses a parametrisation technique. However, even though these schemes perform well in a number of cases, their convergence analysis often seems to remain out of reach, except under additional geometrical conditions [13].

More recently, finite volume schemes using interface values have been studied. In [19] we presented a “hybrid finite volume” (HVF) scheme for any space dimension, which involves edge unknowns in addition to the usual cell unknowns, and in [15], a “mixed finite volume” scheme (MFV) was proposed, which involves the fluxes and the values as unknowns. This is also the case for the mimetic finite difference (MFD) schemes [9, 10], which were introduced previously; in spite of their name, mimetic schemes are very much in the finite volume spirit, since they rely on both a flux balance equation and on the local conservativity of the numerical fluxes, that are probably the two “pillars” of the finite volume philosophy; but then, finite volume schemes are also often called finite difference schemes in the engineering literature because of the finite difference approximation of the fluxes. In fact, a recent benchmark [25] provided sufficient information to suspect that the methods HFV, MFV and MFD indeed coincide at the algebraic level and establishing this is the aim of ongoing work [16]. Let us mention that the Raviart-Thomas mixed finite element method, which also involves edge unknowns, was generalised to handle distorted hexahedral meshes [27]. These schemes require the fluxes or edge unknowns as additional values (or as sole values after hybridisation), and they may be more expensive than cell-centred schemes, especially in the 3D case.

In the two-dimensional case, we also mention [6], which discusses a scheme based on vertex reconstructions, and the family of double mesh schemes [26, 14, 7]. The generalisation of this type of scheme to 3D is the subject of ongoing work.

The scheme that we present here is designed on very general polygonal, possibly non-convex and nonconforming meshes, with the following two priorities in mind:

- For cost reasons and data structure issues, we wish to obtain a symmetric scheme which is as close as possible to a cell-centred scheme, that is to a scheme involving one unknown per control volume (or grid cell).
- For accuracy reasons, we require the local conservativity of the numerical flux to hold at the interfaces between highly heterogeneous media.

In [21], we introduced a cell-centred scheme for the approximation of the Laplace operator on nonconforming grids in the framework of the incompressible Navier-Stokes equations [21] and which may be viewed as a low order nonconforming Galerkin approximation. The scheme (called “SUCCES” in [4]) was also implemented for anisotropic and heterogeneous problems on general meshes, and was shown to be highly competitive for oil reservoir simulation in comparison with other well-known schemes such as the multiple point flux approximation schemes. It is cheaper than the above mentioned hybrid type schemes (HVF, MFV and MFD) because it is based on cell unknowns only. However, it is not as accurate as the hybrid schemes for strongly heterogeneous problems, very likely because of the weaker approximation of the normal fluxes at the heterogeneous interfaces. In the present work, we construct a discretisation scheme (SUSHI) for any kind of polyhedral mesh, which incorporates the best properties of the cell-centred (SUCCES) and hybrid (HFV) schemes: unknowns on the edges are

only introduced when needed, for instance when there is strong medium heterogeneity at these edges. If the set of edge unknowns is empty, then SUSHI reduces to the above mentioned cell-centred scheme; if unknowns are associated to all internal edges, then SUSHI is the hybrid scheme HFV.

The outline of this paper is as follows. In Section 2, we present the guidelines which led us in the construction of convergent schemes on general nonconforming meshes. The practical properties of the resulting schemes are shown through numerical examples in Section 3. Then the mathematical analysis of convergence and error estimation are performed in Section 4. This analysis is based on some discrete functional analytic tools, such as discrete Sobolev inequalities, which are provided in Section 5. Conclusions and perspectives are discussed in Section 6.

## 2 Fundamentals for a class of nonconforming schemes

Let us first present the desired properties which have led us to the design of the schemes under study:

- (P1) The schemes must apply on any type of grid: conforming or nonconforming, 2D and 3D (or more, see for instance the frameworks of kinetic formulations or financial mathematics), consisting of control volumes which are only assumed to be polyhedral (the boundary of each control volume is a finite union of subsets of hyperplanes).
- (P2) The matrices of the linear systems generated are expected to be sparse, symmetric and positive definite.
- (P3) We wish to be able to prove the convergence of the family of discrete solutions to the solution of the continuous problem as the mesh size tends to 0, and of the family of associate gradients to the gradient of the solution, with no regularity assumption on the solution of the continuous problem, and to derive error estimates when the analytic solution is regular enough.

In order to describe the schemes we now introduce some notations for the space discretisation.

**Definition 2.1 (Space discretisation)** *Let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$ , with  $d \in \mathbb{N} \setminus \{0\}$ , and  $\partial\Omega = \overline{\Omega} \setminus \Omega$  its boundary. A discretisation of  $\Omega$ , denoted by  $\mathcal{D}$ , is defined as the triplet  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P})$ , where:*

1.  $\mathcal{M}$  is a finite family of nonempty connected open disjoint subsets of  $\Omega$  (the “control volumes”) such that  $\overline{\Omega} = \cup_{K \in \mathcal{M}} \overline{K}$ . For any  $K \in \mathcal{M}$ , let  $\partial K = \overline{K} \setminus K$  be the boundary of  $K$ ; let  $|K| > 0$  denote the measure of  $K$  and let  $h_K$  denote the diameter of  $K$ .
2.  $\mathcal{E}$  is a finite family of disjoint subsets of  $\overline{\Omega}$  (the “edges” of the mesh), such that, for all  $\sigma \in \mathcal{E}$ ,  $\sigma$  is a nonempty open subset of a hyperplane of  $\mathbb{R}^d$ , whose  $(d-1)$ -dimensional measure  $|\sigma|$  is strictly positive. We also assume that, for all  $K \in \mathcal{M}$ , there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that  $\partial K = \cup_{\sigma \in \mathcal{E}_K} \sigma$ . For any  $\sigma \in \mathcal{E}$ , we denote by  $\mathcal{M}_\sigma = \{K \in \mathcal{M}, \sigma \in \mathcal{E}_K\}$ . We then assume that, for all  $\sigma \in \mathcal{E}$ , either  $\mathcal{M}_\sigma$  has exactly one element and then  $\sigma \subset \partial\Omega$  (the set of these interfaces, called boundary interfaces, is denoted by  $\mathcal{E}_{\text{ext}}$ ) or  $\mathcal{M}_\sigma$  has exactly two elements (the set of these interfaces, called interior interfaces, is denoted by  $\mathcal{E}_{\text{int}}$ ). For all  $\sigma \in \mathcal{E}$ , we denote by  $\mathbf{x}_\sigma$  the barycentre of  $\sigma$ . For all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}_K$ , we denote by  $\mathbf{n}_{K,\sigma}$  the unit vector normal to  $\sigma$  outward to  $K$ .
3.  $\mathcal{P}$  is a family of points of  $\Omega$  indexed by  $\mathcal{M}$ , denoted by  $\mathcal{P} = (\mathbf{x}_K)_{K \in \mathcal{M}}$ , such that for all  $K \in \mathcal{M}$ ,  $\mathbf{x}_K \in K$  and  $K$  is assumed to be  $\mathbf{x}_K$ -star-shaped, which means that for all  $\mathbf{x} \in K$ , the inclusion  $[\mathbf{x}_K, \mathbf{x}] \subset K$  holds. Denoting by  $d_{K,\sigma}$  the Euclidean distance between  $\mathbf{x}_K$  and the hyperplane including  $\sigma$ , one assumes that  $d_{K,\sigma} > 0$ . We then denote by  $D_{K,\sigma}$  the cone with vertex  $\mathbf{x}_K$  and basis  $\sigma$ .

**Remark 2.1** *The above definition applies to a large variety of meshes. Note that no hypothesis is made on the convexity of the control volumes; in fact, generalised hexahedra, i.e. with faces which may be composed of several planar sub-faces may be used. Often encountered in subsurface flow simulations, such hexahedra may have up to 12 faces (resp. 24 faces) if each non planar face is composed of two triangles (resp. four triangles), but only 6 neighbouring control volumes.*

## 2.1 From a “hybrid” finite volume scheme...

The idea of the “hybrid” schemes (among them one may include the mixed finite elements, the mixed finite volume or the mimetic finite difference schemes) is to find an approximation to the solution of (1)–(2) by setting up a system of discrete equations for a family of values  $((u_K)_{K \in \mathcal{M}}, (u_\sigma)_{\sigma \in \mathcal{E}})$  in the control volumes and on the interfaces. The number of unknowns is therefore  $\text{card}(\mathcal{M}) + \text{card}(\mathcal{E})$ . Following the idea of the finite volume framework, Equation (1) is integrated over each control volume  $K \in \mathcal{M}$ , which formally gives (assuming sufficient regularity on  $u$  and  $\Lambda$ ) the following balance equation on the control volume  $K$ :

$$\sum_{\sigma \in \mathcal{E}_K} \left( - \int_{\sigma} \Lambda(\mathbf{x}) \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right) = \int_K f(\mathbf{x}) d\mathbf{x}.$$

The flux  $-\int_{\sigma} \Lambda(\mathbf{x}) \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x})$  is approximated by a function  $F_{K,\sigma}(u)$  of the values  $((u_K)_{K \in \mathcal{M}}, (u_\sigma)_{\sigma \in \mathcal{E}})$  at the “centres” and at the interfaces of the control volumes (in all practical cases,  $F_{K,\sigma}(u)$  only depends on  $u_K$  and all  $(u_{\sigma'})_{\sigma' \in \mathcal{E}_K}$ ). A discrete equation corresponding to (1) is then:

$$\sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(u) = \int_K f(\mathbf{x}) d\mathbf{x} \quad \forall K \in \mathcal{M}. \quad (7)$$

The values  $u_\sigma$  on the interfaces are then introduced so as to allow for a consistent approximation of the normal fluxes in the case of an anisotropic operator and a general, possibly nonconforming mesh. We thus have  $\text{card}(\mathcal{E})$  supplementary unknowns, and need  $\text{card}(\mathcal{E})$  equations to ensure that the problem is well posed. For the boundary faces or edges, these equations are obtained by writing the discrete counterpart of the boundary condition (2):

$$u_\sigma = 0 \quad \forall \sigma \in \mathcal{E}_{\text{ext}}. \quad (8)$$

Following the finite volume ideas, we may write the continuity of the discrete flux for all interior edges, that is to say:

$$F_{K,\sigma}(u) + F_{L,\sigma}(u) = 0, \text{ for } \sigma \in \mathcal{E}_{\text{int}} \text{ such that } \mathcal{M}_\sigma = \{K, L\}. \quad (9)$$

We now have  $\text{card}(\mathcal{M}) + \text{card}(\mathcal{E}_{\text{int}})$  unknowns and equations.

**Remark 2.2** *In the case  $\Lambda(\mathbf{x}) = \lambda(\mathbf{x})\text{Id}$ , on meshes satisfying an orthogonality condition as mentioned in the introduction of this paper (this condition states the orthogonality between the line joining the centres of two neighbouring control volumes with their common interface, see [17, Definition 9.1 p. 762]), a consistent numerical flux is obtained using the two-point formula  $F_{K,\sigma}(u) = \lambda_K |\sigma| (u_K - u_\sigma) / d_{K,\sigma}$ , where  $\lambda_K$  is the average value for  $\lambda$  in  $K$ . Then, writing (9) for all  $\sigma \in \mathcal{E}_{\text{int}}$  such that  $\mathcal{M}_\sigma = \{K, L\}$ , we obtain  $u_\sigma$  as a linear combination of  $u_K$  and  $u_L$ . Plugging this expression into (7), we get a scheme with  $\text{card}(\mathcal{M})$  equations and  $\text{card}(\mathcal{M})$  unknowns (see [17, Section 11.1 pp. 815–820] for more details). In the case of a rectangular (resp. triangular) mesh, this is the well-known five points (resp. four points) scheme with harmonic averages of the diffusion.*

With a proper choice of the expression  $F_{K,\sigma}(u)$ , which we shall introduce below, this scheme, first introduced in [19], is quite efficient for the simulation of fluid flow in heterogeneous media (where harmonic averages for  $\Lambda$  are preferred to arithmetic averages [5]) and may be shown to converge. This

scheme does have one drawback: since the number of unknowns is the sum of the number of control volumes and of interior interfaces, the resulting scheme is quite expensive (although it is sometimes possible to algebraically eliminate the values at the control volumes, as in the mixed hybrid finite element method, see [8, pp. 178-181]).

**Remark 2.3** *Note that in the case of regular conforming simplices (triangles in 2D, tetrahedra in 3D), there is an algebraic possibility to express the unknowns  $(u_\sigma)_{\sigma \in \mathcal{E}}$  as local affine combinations of the values  $(u_K)_{K \in \mathcal{M}}$  and therefore to eliminate them [31]. The idea is to remark that the linear system constituted by the equations (7) for all  $K \in \mathcal{M}_S$ , where  $\mathcal{M}_S$  is the set of all simplices sharing the same interior vertex  $S$ , and (9) for all the interior edges such that  $\mathcal{M}_\sigma \subset \mathcal{M}_S$ , presents as many equations as unknowns  $u_\sigma$ , for  $\sigma \in \cup_{K \in \mathcal{M}_S} \mathcal{E}_K$ . Indeed, the number of edges in  $\cup_{K \in \mathcal{M}_S} \mathcal{E}_K$  such that  $\mathcal{M}_\sigma \not\subset \mathcal{M}_S$  is equal to the number of control volumes in  $\mathcal{M}_S$ . Unfortunately, there is at this time no general result on the invertibility or the symmetry of the matrix of this system, and this method does not apply to other types of meshes than simplicial meshes.*

In order to reduce the computational cost of the scheme, we developed in [21] an idea which is in fact close to the finite element philosophy since we express the finite volume scheme in a weak form; to this end, let us first define the sets  $X_{\mathcal{D}}$  and  $X_{\mathcal{D},0}$  where the discrete unknowns lie, that is to say:

$$X_{\mathcal{D}} = \{v = ((v_K)_{K \in \mathcal{M}}, (v_\sigma)_{\sigma \in \mathcal{E}}), v_K \in \mathbb{R}, v_\sigma \in \mathbb{R}\}, \quad (10)$$

$$X_{\mathcal{D},0} = \{v \in X_{\mathcal{D}} \text{ such that } v_\sigma = 0 \quad \forall \sigma \in \mathcal{E}_{\text{ext}}\}. \quad (11)$$

Multiplying, for any  $v \in X_{\mathcal{D},0}$ , Equation (7) by the value  $v_K$  of  $v$  on the control volume  $K$  and summing over  $K \in \mathcal{M}$  leads to:

$$\sum_{K \in \mathcal{M}} v_K \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(u) = \sum_{K \in \mathcal{M}} v_K \int_K f(\mathbf{x}) d\mathbf{x}.$$

Using (9), we get the following discrete weak formulation:

$$\left\{ \begin{array}{l} \text{Find } u \in X_{\mathcal{D},0} \text{ such that:} \\ \langle u, v \rangle_F = \sum_{K \in \mathcal{M}} v_K \int_K f(\mathbf{x}) d\mathbf{x}, \quad \text{for all } v \in X_{\mathcal{D},0}, \end{array} \right. \quad (12)$$

with

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(u)(v_K - v_\sigma). \quad (13)$$

Note that choosing  $v \in X_{\mathcal{D},0}$  such that  $v_K = 1$ ,  $v_L = 0$  for any  $L \in \mathcal{M}, L \neq K$  and  $v_\sigma = 0$  for any  $\sigma \in \mathcal{E}$  yields (7). Similarly, choosing  $v \in X_{\mathcal{D},0}$  such that  $v_K = 0$  for any  $K \in \mathcal{M}$ , and  $v_\sigma = 1$  and  $v_\tau = 0$  for any  $\tau \in \mathcal{E}, \tau \neq \sigma$  leads to (9). Therefore the hybrid finite volume scheme (7)–(9) is equivalent to the discrete weak formulation (12).

## 2.2 ... to a nonconforming finite element scheme...

We may then choose to use the weak discrete form (13) as an approximation of the bilinear form  $a(\cdot, \cdot)$ , but with a space of dimension smaller than that of  $X_{\mathcal{D},0}$ . This can be achieved by expressing the value of  $u$  on any interior interface  $\sigma \in \mathcal{E}_{\text{int}}$  as a consistent barycentric combination of the values  $u_K$ :

$$u_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K u_K, \quad (14)$$

where  $(\beta_\sigma^K)_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_{\text{int}}}}$  is a family of real numbers, with  $\beta_\sigma^K \neq 0$  only for some control volumes  $K$  close to  $\sigma$ , and such that

$$\sum_{K \in \mathcal{M}} \beta_\sigma^K = 1 \text{ and } \mathbf{x}_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K \mathbf{x}_K \quad \forall \sigma \in \mathcal{E}_{\text{int}}. \quad (15)$$

This ensures that if  $\varphi$  is a regular function, then  $\varphi_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K \varphi(\mathbf{x}_K)$  is a consistent approximation of  $\varphi(\mathbf{x}_\sigma)$  for  $\sigma \in \mathcal{E}_{\text{int}}$ . We recall that the values  $u_\sigma, \sigma \in \mathcal{E}_{\text{ext}}$  are set to 0 in order to respect the boundary conditions (2). Hence the new scheme reads:

$$\left\{ \begin{array}{l} \text{Find } u \in X_{\mathcal{D},0} \text{ such that } u_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K u_K \quad \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ and} \\ \langle u, v \rangle_F = \sum_{K \in \mathcal{M}} v_K \int_K f(\mathbf{x}) d\mathbf{x}, \text{ for all } v \in X_{\mathcal{D},0} \text{ with } v_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K v_K \quad \forall \sigma \in \mathcal{E}_{\text{int}}. \end{array} \right. \quad (16)$$

This method has been shown in [21] to be efficient in the case of a problem where  $\Lambda = \text{Id}$  (for the approximation of the viscous terms in the Navier-Stokes problem). With an appropriate choice for the expression of the numerical flux, it also yields conservativity in a certain sense (more on this below), but no longer to the classical (in the finite volume framework) equation (9): indeed, since the degrees of freedom on the edges are no longer present, one may not use  $v_\sigma = 1$  to recover (9). Note also that taking  $v_K = 1$  does not yield (7). This scheme has been implemented for the discretisation of the diffusive term in the incompressible Navier Stokes equations on general two- or three-dimensional grids, and gives excellent results [11, 12]. Unfortunately, because of poor approximation of the local flux at strongly heterogeneous interfaces, this approach is not sufficient to provide accurate results for some types of flows in heterogeneous media, as we shall show in Section 3. This is especially true when using coarse meshes, as is often the case in industrial problems.

### 2.3 ... to an optimal compromise?

Therefore we now propose a scheme which has the advantage of both techniques: we shall use equation (13) and keep the unknowns  $u_\sigma$  on the edges which require them, for instance those where the matrix  $\Lambda$  is discontinuous: hence (9) will hold for all edges associated to these unknowns; for all other interfaces, we shall impose the values of  $u$  using (14), and therefore eliminate these unknowns. Let us decompose the set  $\mathcal{E}_{\text{int}}$  of interfaces into two nonintersecting subsets, that is:  $\mathcal{E}_{\text{int}} = \mathcal{B} \cup \mathcal{H}, \mathcal{H} = \mathcal{E}_{\text{int}} \setminus \mathcal{B}$ . The interface unknowns associated with  $\mathcal{B}$  will be computed by using the barycentric formula (14).

**Remark 2.4** *Note that, although the accuracy of the scheme is increased in practice when the points where the matrix  $\Lambda$  is discontinuous are located within the set  $\bigcup_{\sigma \in \mathcal{H}} \sigma$ , such a property is not needed in the mathematical study of the scheme.*

Let us introduce the space  $X_{\mathcal{D},\mathcal{B}} \subset X_{\mathcal{D},0}$  defined by:

$$X_{\mathcal{D},\mathcal{B}} = \{v \in X_{\mathcal{D}} \text{ such that } v_\sigma = 0 \text{ for all } \sigma \in \mathcal{E}_{\text{ext}} \text{ and } v_\sigma \text{ satisfying (14) for all } \sigma \in \mathcal{B}\}. \quad (17)$$

The composite scheme which we consider in this work reads:

$$\left\{ \begin{array}{l} \text{Find } u \in X_{\mathcal{D},\mathcal{B}} \text{ such that:} \\ \langle u, v \rangle_F = \sum_{K \in \mathcal{M}} v_K \int_K f(\mathbf{x}) d\mathbf{x}, \text{ for all } v \in X_{\mathcal{D},\mathcal{B}}. \end{array} \right. \quad (18)$$

We therefore obtain a symmetric scheme with  $\text{card}(\mathcal{M}) + \text{card}(\mathcal{H})$  equations and unknowns. It is thus less expensive while it remains accurate (for the choice of numerical flux given below) even in the case of strong heterogeneity (see section 3).

Note that with the present scheme, (9) holds for all  $\sigma \in \mathcal{H}$ , but not generally for any  $\sigma \in \mathcal{B}$ . However, fluxes between pairs of control volumes can nevertheless be identified. Indeed, we may write

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \left( \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{H}} F_{K,\sigma}(u)(v_K - v_\sigma) + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{B}} \sum_{L \in \mathcal{M}} F_{K,\sigma}(u) \beta_\sigma^L (v_K - v_L) \right),$$

and therefore:

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{H}} F_{K,\sigma}(u)(v_K - v_\sigma) + \frac{1}{2} \sum_{(K,L) \in \mathcal{N}_D} F_{K,L}(u)(v_K - v_L),$$

where

$$\mathcal{N}_D = \{(K, L) \in \mathcal{M}^2, \exists \sigma \in \mathcal{E}_K \cap \mathcal{B}, \beta_\sigma^L \neq 0 \text{ or } \exists \sigma \in \mathcal{E}_L \cap \mathcal{B}, \beta_\sigma^K \neq 0\},$$

and

$$F_{K,L}(u) = \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{B}} F_{K,\sigma}(u) \beta_\sigma^L - \sum_{\sigma \in \mathcal{E}_L \cap \mathcal{B}} F_{L,\sigma}(u) \beta_\sigma^K.$$

Note that, if  $(K, L) \in \mathcal{N}_D$ , then  $(L, K) \in \mathcal{N}_D$  and  $F_{K,L}(u) = -F_{L,K}(u)$ ; furthermore,  $F_{K,L}(u) \neq 0$  implies  $(K, L) \in \mathcal{N}_D$ , and the scheme's stencil is determined by the set  $\{L \in \mathcal{M} \text{ such that } (K, L) \in \mathcal{N}_D\}$ . Then, taking  $v_K = 1$  and all other degrees of freedom of  $v \in X_{D,\mathcal{B}}$  equal to 0, (18) yields

$$\sum_{\sigma \in \mathcal{E}_K \cap \mathcal{H}} F_{K,\sigma}(u) + \sum_{\substack{L \in \mathcal{M} \\ (K,L) \in \mathcal{N}_D}} F_{K,L}(u) = \int_K f(\mathbf{x}) d\mathbf{x},$$

which shows the “finite volume philosophy” of the scheme.

**Remark 2.5 (Other boundary conditions)** *In the case of Neumann or Robin boundary conditions, the discrete space  $X_{D,\mathcal{B}}$  is modified to include the unknowns associated to the corresponding edges, and the resulting discrete weak formulation is then straightforward.*

**Remark 2.6 (Extension of the scheme)** *For consistency reasons, it is preferable that the coefficients  $\beta_\sigma^K$  associated with  $\sigma \in \mathcal{B}$  be nonzero for points  $x_K$  that lie in the same “regularity zone” of the solution as  $x_\sigma$  (that is with a zone with no diffusion tensor discontinuity). This is not always easy: indeed, in the tilted barrier example described in Section 3.3 below, the barrier contains only one layer of grid cells, so that, for an internal interface of this layer, it is difficult to use points  $x_L$  that are located in the same diffusion regularity zone with respect to  $x_K$ . There is, however, no additional difficulty to replace (14) in the definition of (17) by*

$$u_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K u_K + \sum_{\sigma' \in \mathcal{H}} \beta_{\sigma'}^{\sigma'} u_{\sigma'} \quad \forall \sigma \in \mathcal{B}, \quad (19)$$

$$\sum_{K \in \mathcal{M}} \beta_\sigma^K + \sum_{\sigma' \in \mathcal{H}} \beta_{\sigma'}^{\sigma'} = 1 \text{ and } \mathbf{x}_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K \mathbf{x}_K + \sum_{\sigma' \in \mathcal{H}} \beta_{\sigma'}^{\sigma'} \mathbf{x}_{\sigma'} \quad \forall \sigma \in \mathcal{B}. \quad (20)$$

*This trick solves the consistency issue without switching the edge to the hybrid set  $\mathcal{H}$ , while all the mathematical properties shown below still hold.*

## 2.4 Construction of the fluxes using a discrete gradient

For the definition of the schemes to be complete, there now remains to explain how we find a convenient expression for  $F_{K,\sigma}(u)$  with respect to the discrete unknowns. An idea that has been used in several of the schemes referred to in the Introduction is to look for a consistent expression of the flux by using adequate linear combinations of the unknowns; however, referring to the beginning of Section 2, such a reconstruction does not in general lead to the desired properties (P2) (symmetric definite positive matrices) and (P3) (convergence). Our idea here is different: it is based on the identification of the numerical fluxes  $F_{K,\sigma}(u)$  through the mesh-dependent bilinear form  $\langle \cdot, \cdot \rangle_F$  defined in (13), using the expression of a discrete gradient. Indeed let us assume that, for all  $u \in X_{\mathcal{D}}$ , we have constructed a discrete gradient  $\nabla_{\mathcal{D}}u$ , we then seek a family  $(F_{K,\sigma}(u))_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_K}}$  such that

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(u)(v_K - v_\sigma) = \int_{\Omega} \nabla_{\mathcal{D}}u(\mathbf{x}) \cdot \Lambda(\mathbf{x}) \nabla_{\mathcal{D}}v(\mathbf{x}) d\mathbf{x} \quad \forall u, v \in X_{\mathcal{D}}. \quad (21)$$

**Remark 2.7 (On the construction of the discrete fluxes)** *Note that it is always possible to deduce an expression for  $F_{K,\sigma}(u)$  satisfying (21), under the sufficient condition that, for all  $K \in \mathcal{M}$  and a.e.  $\mathbf{x} \in K$ ,  $\nabla_{\mathcal{D}}u(\mathbf{x})$  is expressed as a linear combination of  $(u_\sigma - u_K)_{\sigma \in \mathcal{E}_K}$ , the coefficients of which are measurable bounded functions of  $\mathbf{x}$ . This property is ensured in the construction of  $\nabla_{\mathcal{D}}u(\mathbf{x})$  given below.*

We prove in Section 4 below that the desired properties (P2) and (P3) hold if the discrete gradient satisfies the following properties:

1. (Weak compactness) For a sequence of space discretisations of  $\Omega$  with mesh size tending to 0, if the sequence of associated grid functions is bounded in some sense, then their discrete gradient converges at least weakly in  $L^2(\Omega)^d$  to the gradient of an element of  $H_0^1(\Omega)$ ;
2. (Consistency) If  $\varphi$  is a regular function from  $\bar{\Omega}$  to  $\mathbb{R}$ , the discrete gradient of the piece-wise function defined by taking the value  $\varphi(\mathbf{x}_K)$  on each control volume  $K$  and  $\varphi(\mathbf{x}_\sigma)$  on each edge  $\sigma$  is a consistent approximation of the gradient of  $\varphi$ .

Let us first define:

$$\nabla_K u = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (u_\sigma - u_K) \mathbf{n}_{K,\sigma} \quad \forall K \in \mathcal{M}, \forall u \in X_{\mathcal{D}}, \quad (22)$$

where  $\mathbf{n}_{K,\sigma}$  is the outward to  $K$  normal unit vector,  $|K|$  and  $|\sigma|$  are the usual measures (volumes, areas, or lengths) of  $K$  and  $\sigma$ . The consistency of formula (22) stems from the following geometrical relation:

$$\sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (\mathbf{x}_\sigma - \mathbf{x}_K)^t = |K| \text{Id} \quad \forall K \in \mathcal{M}, \quad (23)$$

where  $(\mathbf{x}_\sigma - \mathbf{x}_K)^t$  is the transpose of  $\mathbf{x}_\sigma - \mathbf{x}_K \in \mathbb{R}^d$ , and Id is the  $d \times d$  identity matrix. Indeed, for any linear function defined on  $\Omega$  by  $\psi(\mathbf{x}) = \mathbf{G} \cdot \mathbf{x}$  with  $\mathbf{G} \in \mathbb{R}^d$ , assuming that  $u_\sigma = \psi(\mathbf{x}_\sigma)$  and  $u_K = \psi(\mathbf{x}_K)$ , we get  $u_\sigma - u_K = (\mathbf{x}_\sigma - \mathbf{x}_K)^t \mathbf{G} = (\mathbf{x}_\sigma - \mathbf{x}_K)^t \nabla \psi$ , hence (22) leads to  $\nabla_K u = \nabla \psi$ .

Since the coefficient of  $u_K$  in (22) is in fact equal to zero, a re-construction of the discrete gradient  $\nabla_{\mathcal{D}}u$  solely based on (22) cannot lead to a definite discrete bilinear form in the general case. Hence, we now introduce a stabilised gradient:

$$\nabla_{K,\sigma} u = \nabla_K u + R_{K,\sigma} u \mathbf{n}_{K,\sigma}, \quad (24)$$

with

$$R_{K,\sigma} u = \frac{\sqrt{d}}{d_{K,\sigma}} (u_\sigma - u_K - \nabla_K u \cdot (\mathbf{x}_\sigma - \mathbf{x}_K)), \quad (25)$$

(recall that  $d$  is the space dimension and  $d_{K,\sigma}$  is the Euclidean distance between  $\mathbf{x}_K$  and  $\sigma$ ). We may then define  $\nabla_{\mathcal{D}}u$  as the piece-wise constant function equal to  $\nabla_{K,\sigma}u$  a.e. in the cone  $D_{K,\sigma}$  with vertex  $\mathbf{x}_K$  and basis  $\sigma$ :

$$\nabla_{\mathcal{D}}u(\mathbf{x}) = \nabla_{K,\sigma}u \text{ for a.e. } \mathbf{x} \in D_{K,\sigma}. \quad (26)$$

Note that, from the definition (25), thanks to (23) and to the definition (22), we get that

$$\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|d_{K,\sigma}}{d} R_{K,\sigma}u \mathbf{n}_{K,\sigma} = 0 \quad \forall K \in \mathcal{M}. \quad (27)$$

We prove in Lemmata 4.2 and 4.3 below that the discrete gradient defined by (22)-(26) indeed satisfies the above stated weak compactness and consistency properties. In order to identify the numerical fluxes  $F_{K,\sigma}(u)$  through Relation (21), we put the discrete gradient in the form

$$\nabla_{K,\sigma}u = \sum_{\sigma' \in \mathcal{E}_K} (u_{\sigma'} - u_K) \mathbf{y}^{\sigma\sigma'},$$

with

$$\mathbf{y}^{\sigma\sigma'} = \begin{cases} \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} + \frac{\sqrt{d}}{d_{K,\sigma}} \left( 1 - \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \right) \mathbf{n}_{K,\sigma} & \text{if } \sigma = \sigma' \\ \frac{|\sigma'|}{|K|} \mathbf{n}_{K,\sigma'} - \frac{\sqrt{d}}{d_{K,\sigma}|K|} |\sigma'| \mathbf{n}_{K,\sigma'} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \mathbf{n}_{K,\sigma} & \text{otherwise.} \end{cases} \quad (28)$$

Thus,

$$\int_{\Omega} \nabla_{\mathcal{D}}u(\mathbf{x}) \cdot \Lambda(\mathbf{x}) \nabla_{\mathcal{D}}v(\mathbf{x}) d\mathbf{x} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \sum_{\sigma' \in \mathcal{E}_K} A_K^{\sigma\sigma'} (u_\sigma - u_K) (v_{\sigma'} - v_K) \quad \forall u, v \in X_{\mathcal{D}}, \quad (29)$$

with,

$$A_K^{\sigma\sigma'} = \sum_{\sigma'' \in \mathcal{E}_K} \mathbf{y}^{\sigma''\sigma} \cdot \Lambda_{K,\sigma''} \mathbf{y}^{\sigma''\sigma'} \text{ and } \Lambda_{K,\sigma''} = \int_{D_{K,\sigma''}} \Lambda(\mathbf{x}) d\mathbf{x}. \quad (30)$$

Then we get that the local matrices  $(A_K^{\sigma\sigma'})_{\sigma\sigma' \in \mathcal{E}_K}$  are symmetric and positive, and the identification of the numerical fluxes using (21) leads to the expression:

$$F_{K,\sigma}(u) = \sum_{\sigma' \in \mathcal{E}_K} A_K^{\sigma\sigma'} (u_K - u_{\sigma'}). \quad (31)$$

**Remark 2.8 (Link with the MFD method)** *The above technique yields an explicit construction of a particular MFD method. Indeed, if one chooses  $x_K$  as the centre of mass of  $K$ , the matrix  $A_K$  defined by (30) is an adequate choice for the matrix  $\mathbb{W}_E$  which is a parameter in the general formulation of the family of MFD methods as proposed in [10]. The advantages of the specific matrix  $A_K$  are that:*

- on particular meshes, taking a natural choice for  $x_K$  (for instance the circumcenter for a triangular mesh in the case of a 2D isotropic problem), it degenerates to a diagonal matrix (see Lemma 2.1 below);
- it is linked to an explicit formulation of a consistent gradient, which is used to define the discrete bilinear form (29).

Note however that the SUSHI scheme defined by (18) is not the MFD method of [9, 10]; the main reason is that according to the choice  $\mathcal{B}$ , the SUSHI scheme may be either a completely cell-centred scheme, or a partly or fully hybrid scheme, while the MFD method is a pure hybrid scheme. Note

also that in *SUSHI*, one may take any point in cell  $K$  for  $x_K$ , while the MFD schemes [9, 10] are constructed with the centre of mass (however, this choice might be generalised).

Note that the procedure which we describe in Section 2.2 to write a cell-centred scheme could be applied to any mimetic scheme (or low order mixed finite element scheme) to yield a centred scheme. However, further investigations are needed to determine under what conditions the present convergence analysis extends to mimetic schemes, and conversely, whether the mimetic analysis applies to the *SUSHI* scheme (ongoing work, [16]).

The fluxes defined by (22)-(31) satisfy certain properties which are detailed in Lemma 4.4, and which allow us to prove the convergence of the scheme, as is shown in Theorem 4.1. Note that it seems difficult to deduce such properties from fluxes obtained by using natural expansions of regular functions. Note also that both Lemma 4.4 and Theorem 4.1 hold for general heterogeneous, anisotropic and possibly discontinuous fields  $\Lambda$ , for which the solution  $u$  of (6) is not in general more regular than  $u \in H_0^1(\Omega)$ . In the case where  $\Lambda$  and  $u$  are regular enough, the local flux consistency satisfied by (31) is used in order to obtain an error estimate, see Theorem 4.2. The coefficient  $\sqrt{d}$  may be replaced by any positive real number without any change in the proof of convergence; in fact, for certain problems it can be interesting to use another coefficient, as described in [20] for the so called “*SUSHI-P*” scheme (P for parametric, meaning that the user may choose the stabilisation coefficient as well as the set of edges  $\mathcal{B}$ ). The choice  $\sqrt{d}$  is however natural in the sense that with this value, if  $\mathcal{B} = \emptyset$ , the scheme boils down in two dimensions to the well-known harmonic averaging five points scheme on rectangles and a four-point scheme on triangles; more generally, in any space dimension, even if  $\mathcal{B} \neq \emptyset$  and taking the most natural value for  $u_\sigma$  if  $\sigma \in \mathcal{B}$ , the resulting flux is a two-point flux on meshes that satisfy the “superadmissibility condition” (32), not necessarily with a harmonic averaging of  $\Lambda$  in the case  $\sigma \in \mathcal{B}$ ; this is proven in the next lemma. Note that this superadmissibility condition is also satisfied by rectangular parallelepipeds in three dimensions but unfortunately not by tetrahedra.

**Lemma 2.1 (Superadmissible mesh and two-point flux)** *Let  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P})$  be a discretisation of  $\Omega$  in the sense of Definition 2.1, satisfying the following superadmissibility condition:*

$$\mathbf{n}_{K,\sigma} = \frac{\mathbf{x}_\sigma - \mathbf{x}_K}{d_{K,\sigma}} \quad \forall K \in \mathcal{M}, \forall \sigma \in \mathcal{E}_K. \quad (32)$$

Let us furthermore assume that  $\Lambda(\mathbf{x}) = \lambda(\mathbf{x})\text{Id}$ , where  $\lambda$  is a piece-wise constant function from  $\Omega$  to  $\mathbb{R}$ , which is equal to a constant  $\lambda_K$  in each  $K \in \mathcal{M}$ ; then, the inner product defined by (21)-(25) reads:

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \lambda_K \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_K - u_\sigma)(v_K - v_\sigma).$$

Moreover, choosing thanks to (32),  $\mathbf{x}_\sigma = (d_{K,\sigma}\mathbf{x}_L + d_{L,\sigma}\mathbf{x}_K)/(d_{K,\sigma} + d_{L,\sigma})$  for  $\sigma \in \mathcal{E}_{\text{int}}$  with  $\mathcal{M}_\sigma = \{K, L\}$  in (15), the scheme (18) is the following two-point flux scheme:

$$\sum_{K \in \mathcal{M}} F_{K,\sigma} = \int_K f(\mathbf{x}) \, d\mathbf{x}, \quad (33)$$

$$F_{K,\sigma} = \frac{\lambda_K \lambda_L (d_{K,\sigma} + d_{L,\sigma})}{\lambda_K d_{L,\sigma} + \lambda_L d_{K,\sigma}} \frac{|\sigma|}{d_{K,\sigma} + d_{L,\sigma}} (u_K - u_L) \text{ if } \sigma \in \mathcal{E}_{\text{int}} \cap \mathcal{H}, \mathcal{M}_\sigma = \{K, L\}, \quad (34)$$

$$F_{K,\sigma} = \frac{d_{K,\sigma} \lambda_K + d_{L,\sigma} \lambda_L}{d_{K,\sigma} + d_{L,\sigma}} \frac{|\sigma|}{d_{K,\sigma} + d_{L,\sigma}} (u_K - u_L) \text{ if } \sigma \in \mathcal{E}_{\text{int}} \cap \mathcal{B}, \mathcal{M}_\sigma = \{K, L\}, \quad (35)$$

$$F_{K,\sigma} = \lambda_K \frac{|\sigma|}{d_{K,\sigma}} u_K \text{ if } \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K, \quad (36)$$

PROOF. Let us compute  $\langle u, v \rangle_F$  under the assumptions of Lemma 2.1. From (21) and thanks to (27) we get:

$$\begin{aligned} \langle u, v \rangle_{\mathcal{D}} &= \sum_{K \in \mathcal{M}} \lambda_K \int_K \nabla_{\mathcal{D}} u(\mathbf{x}) \cdot \nabla_{\mathcal{D}} v(\mathbf{x}) d\mathbf{x} \\ &= \sum_{K \in \mathcal{M}} \lambda_K \left( |K| \nabla_K u \cdot \nabla_K v + \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} R_{K,\sigma} u R_{K,\sigma} v \right). \end{aligned}$$

Now from the definition (22) and thanks to the assumption (32), the discrete gradient given by (22) may be written as follows:

$$\nabla_K v = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (v_\sigma - v_K) (\mathbf{x}_\sigma - \mathbf{x}_K) \quad \forall K \in \mathcal{M}, \forall v \in X_{\mathcal{D}},$$

From (23), we get

$$\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} \frac{\sqrt{d}}{d_{K,\sigma}} (\mathbf{x}_\sigma - \mathbf{x}_K) \frac{\sqrt{d}}{d_{K,\sigma}} (\mathbf{x}_\sigma - \mathbf{x}_K)^t = |K| \text{Id}.$$

Therefore, we get that

$$\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} R_{K,\sigma} u R_{K,\sigma} v = \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_\sigma - u_K) (v_\sigma - v_K) - |K| \nabla_K u \cdot \nabla_K v,$$

which in turn yields that

$$\langle u, v \rangle_{\mathcal{D}} = \sum_{K \in \mathcal{M}} \lambda_K \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_\sigma - u_K) (v_\sigma - v_K).$$

Hence the matrix  $A_K$  only contains the terms  $\frac{|\sigma|}{d_{K,\sigma}}$  on the diagonal, and the flux  $F_{K,\sigma}(u)$  is given by

$$F_{K,\sigma}(u) = \lambda_K \frac{|\sigma|}{d_{K,\sigma}} (u_K - u_\sigma).$$

Then the scheme (18) can be written as a classical cell-centred finite volume scheme, with two-point fluxes  $F_{K,L}(u) = -F_{L,K}(u)$  for any  $\sigma \in \mathcal{E}_{\text{int}}$  with  $\mathcal{M}_\sigma = \{K, L\}$ . Indeed, in the case  $\sigma \notin \mathcal{B}$ , the above expression of  $F_{K,\sigma}(u)$  allows us to get the following expression of  $u_\sigma$  from (9):

$$u_\sigma = \frac{\frac{\lambda_K}{d_{K,\sigma}} u_K + \frac{\lambda_L}{d_{L,\sigma}} u_L}{\frac{\lambda_K}{d_{K,\sigma}} + \frac{\lambda_L}{d_{L,\sigma}}}.$$

This yields the harmonic averaging two-point flux

$$F_{K,L}(u) = |\sigma| \frac{\frac{\lambda_K}{d_{K,\sigma}} \frac{\lambda_L}{d_{L,\sigma}}}{\frac{\lambda_K}{d_{K,\sigma}} + \frac{\lambda_L}{d_{L,\sigma}}} (u_K - u_L).$$

In the case  $\sigma \in \mathcal{B}$ , the two-point barycentric formula  $u_\sigma = (d_{K,\sigma} u_L + d_{L,\sigma} u_K) / (d_{K,\sigma} + d_{L,\sigma})$  together with (18) leads to the resulting two-point flux

$$F_{K,L}(u) = \frac{d_{K,\sigma} \lambda_K + d_{L,\sigma} \lambda_L}{d_{K,\sigma} + d_{L,\sigma}} \frac{|\sigma|}{d_{K,\sigma} + d_{L,\sigma}} (u_K - u_L).$$

□

### 3 Numerical results

We present some numerical results obtained with various choices of  $\mathcal{B}$  in the scheme (18), (13) with the flux (31), which we synthesise here for the sake of clarity:

$$\left\{ \begin{array}{l} \text{Find } u \in X_{\mathcal{D},\mathcal{B}} \text{ (that is } (u_K)_{K \in \mathcal{M}}, (u_\sigma)_{\sigma \in \mathcal{H}}), \text{ such that:} \\ \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(u)(v_K - v_\sigma) = \sum_{K \in \mathcal{M}} v_K \int_K f(\mathbf{x})d\mathbf{x}, \text{ for all } v \in X_{\mathcal{D},\mathcal{B}}, \\ \text{with } F_{K,\sigma}(u) = \sum_{\sigma' \in \mathcal{E}_K} A_K^{\sigma\sigma'}(u_{\sigma'} - u_K) \quad \forall K \in \mathcal{M}, \forall \sigma \in \mathcal{E}_K. \end{array} \right. \quad (37)$$

where the matrices  $A_K^{\sigma\sigma'}$  are defined by (30)-(28). In the following, we shall use the choices  $\mathcal{B} = \emptyset$  (HFV),  $\mathcal{B} = \mathcal{E}_{\text{int}}$  or  $\mathcal{B}$  the set of edges which are located on the diffusion tensor discontinuity interfaces; this latter choice is reported as SUSHI-NP (for non parametric) in [20], in contrast with SUSHI-P (for parametric) where the choice of the set  $\mathcal{B}$  may be different, along with the value of the stabilisation coefficient in (25).

#### 3.1 Implementation

Let us first describe an implementation aspects of the scheme. The unknowns, *i.e.* the values  $u_K$ , for  $K \in \mathcal{M}$  and the values  $u_\sigma$ ,  $\sigma \in \mathcal{E}_{\text{int}} \cap \mathcal{H}$ , are ordered as  $(u_i)_{i=1,\dots,N}$ . The  $N \times N$  matrix and the  $N \times 1$  right-hand-side of the linear system resulting from (18) are computed thanks to a loop over the control volumes  $K \in \mathcal{M}$  and to an inner loop on each edge  $\sigma \in \mathcal{E}_K$ . Let us detail the matrix computation loop.

1. All stored matrix coefficients are initially set to 0.
2. The expression  $F_{K,\sigma}(u)$  is written in the form  $F_{K,\sigma}(u) = \sum_{i=1,\dots,N} a_{K,\sigma}^{(i)} u_i$ , where the nonzero coefficients  $(a_{K,\sigma}^{(i)})_{i=1,\dots,N}$  are only locally computed (they are not stored for all  $K$  and  $\sigma$ ). These coefficients are obtained after the elimination of all  $(u_\sigma)_{\sigma \in \mathcal{E}_K \cap \mathcal{B}}$  in (31):

$$F_{K,\sigma}(u) = \sum_{\sigma' \in \mathcal{E}_K \cap \mathcal{H}} A_K^{\sigma\sigma'}(u_K - u_{\sigma'}) + \sum_{\sigma' \in \mathcal{E}_K \cap \mathcal{B}} A_K^{\sigma\sigma'} \sum_{L \in \mathcal{M}} \beta_{\sigma'}^L (u_K - u_L).$$

3. The line of the matrix corresponding to the unknown  $u_K$  is incremented at the column  $j$  with the coefficient  $a_{K,\sigma}^{(j)}$ .
4. If  $\sigma \in \mathcal{B}$  with  $v_\sigma = \sum_{L \in \mathcal{M}} \beta_\sigma^L v_L$  for any  $v \in X_{\mathcal{D},\mathcal{B}}$ , the line of the matrix corresponding to each  $L \in \mathcal{M}$  such that  $\beta_\sigma^L \neq 0$  is incremented at the column  $j$  with the coefficient  $-\beta_\sigma^L a_{K,\sigma}^{(j)}$ .
5. If  $\sigma \in \mathcal{E}_{\text{int}} \cap \mathcal{H}$ , the line of the matrix corresponding to the edge  $\sigma$  is incremented at the column  $j$  with the coefficient  $-a_{K,\sigma}^{(j)}$ .

This procedure is identical in the cases  $\mathcal{B} = \emptyset$  (HFV),  $\mathcal{B} \neq \emptyset$  and  $\mathcal{B} = \mathcal{E}_{\text{int}}$ . However, in the case where  $\mathcal{B} = \emptyset$  (HFV), one may eliminate the unknowns  $u_K$  with respect to the unknowns  $u_\sigma$ , as in the hybrid implementation of the mixed finite element method.

#### 3.2 Order of convergence

We consider here the numerical resolution of Equation (1) supplemented by the homogeneous Dirichlet boundary condition (2); the right-hand side is chosen so as to obtain an exact solution to the problem

and easily compute the error between the exact and approximate solutions. We consider Problem (1)-(2) with a constant matrix  $\Lambda$ :

$$\Lambda = \begin{pmatrix} 1.5 & .5 \\ .5 & 1.5 \end{pmatrix}, \quad (38)$$

and choose  $f : \Omega \rightarrow \mathbb{R}$  such that the exact solution to Problem (1)-(2) is  $\bar{u}$  defined by  $\bar{u}(x, y) = 16x(1-x)y(1-y)$  for any  $(x, y) \in \bar{\Omega}$ . Note that in this case, the composite scheme is in fact the cell-centred scheme, there are no edge-unknowns.

Let us first consider conforming meshes, such as the triangular meshes which are depicted on Figure 1, and uniform square meshes.

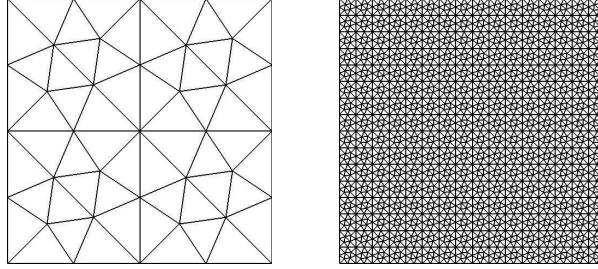


Figure 1: Regular conforming coarse and fine triangular grids

For both  $\mathcal{B} = \emptyset$  (pure hybrid scheme: HFV) and  $\mathcal{B} = \mathcal{E}_{\text{int}}$  (cell-centred scheme), the order of convergence is close to 2 for the unknown  $u$  and 1 for its gradient. Of course, the hybrid scheme is almost three times more costly in terms of number of unknowns than the cell-centred scheme for a given precision. However, the number of nonzero terms in the matrix is, again for a given precision on the approximate solution, larger for the cell-centred scheme than for the hybrid scheme. Hence the number of unknowns is probably not a sufficient criterion for assessing the cost of the scheme.

Results were also obtained in the case of uniform square or rectangular meshes. They show a better rate of convergence of the gradient (order 2 in the case of  $\mathcal{H} = \mathcal{E}_{\text{int}}$  and 1.5 in the case  $\mathcal{B} = \mathcal{E}_{\text{int}}$ ), even though the rate of convergence of the approximate solution remains unchanged and close to 2.

We then use a rectangular nonconforming mesh, obtained by cutting vertically the domain into two parts and using a rectangular grid of  $3n \times 2n$  (resp.  $5n \times 2n$ ) on the first (resp. second side), where  $n$  is the number of the mesh,  $n = 1, \dots, 7$ . Again, the order of convergence which we obtain is 2 for  $u$  and around 1.8 for the gradient. We give in Table 1 below the errors obtained in the discrete  $L^2$  norm for  $u$  and  $\nabla u$  for a nonconforming mesh and (in terms of number of unknowns) and for the rectangular  $4 \times 6$  and  $4 \times 10$  conforming rectangular meshes, for both the hybrid and cell-centred schemes. We show in Figure 2 the solutions for the corresponding grids (which look much the same for the two schemes).

Further detailed results on several problems and conforming, nonconforming and distorted meshes may be found in [20].

### 3.3 The case of a highly heterogeneous tilted barrier

We now turn to the heterogeneous case. The domain  $\Omega = ]0, 1[ \times ]0, 1[$  is composed of 3 sub-domains, which are depicted in Figure 3:  $\Omega_1 = \{(x, y) \in \Omega; \varphi_1(x, y) < 0\}$ , with  $\varphi_1(x, y) = y - \delta(x - .5) - .475$ ,  $\Omega_2 = \{(x, y) \in \Omega; \varphi_1(x, y) > 0, \varphi_2(x, y) < 0\}$ , with  $\varphi_2(x, y) = \varphi_1(x, y) - 0.05$ ,  $\Omega_3 = \{(x, y) \in \Omega; \varphi_2(x, y) > 0\}$ , and  $\delta = 0.2$  is the slope of the drain (see Figure 3). Dirichlet boundary conditions are imposed by setting the boundary values to those of the analytical solution given by  $u(x, y) = -\varphi_1(x, y)$  on  $\Omega_1 \cup \Omega_3$  and  $u(x, y) = -\varphi_1(x, y)/10^{-2}$  on  $\Omega_2$ .

n	NU		NM		$\epsilon(u)$		$\epsilon(\nabla u)$	
	$\mathcal{B} = \emptyset$	$\mathcal{B} = \mathcal{E}_{\text{int}}$	$\mathcal{B} = \emptyset$	$\mathcal{B} = \mathcal{E}_{\text{int}}$	$\mathcal{B} = \emptyset$	$\mathcal{B} = \mathcal{E}_{\text{int}}$	$\mathcal{B} = \emptyset$	$\mathcal{B} = \mathcal{E}_{\text{int}}$
C1	130	48	874	488	1.28E-01	1.20E-01	1.64E-02	3.57E-02
NC	182	64	1334	724	1.03E-01	9.43E-02	1.66E-02	3.69E-02
C2	222	80	1542	864	7.61E-02	7.09E-02	9.18E-03	2.44E-02

Table 1: Error for the nonconforming rectangular mesh, pure hybrid scheme ( $\mathcal{B} = \emptyset$ ) and centred ( $\mathcal{B} = \mathcal{E}_{\text{int}}$ ) schemes. For both schemes NU is the number of unknowns in the resulting linear system, NM is the number of nonzero terms in the matrix,  $\epsilon(u)$  is the discrete  $L^2$  norm of the error of the solution and  $\epsilon(\nabla u)$  is the discrete  $L^2$  norm of the error in the gradient. C1 and C2 are the two conforming meshes represented on the left and the right in Figure 2, and NC is the nonconforming one represented in the middle.

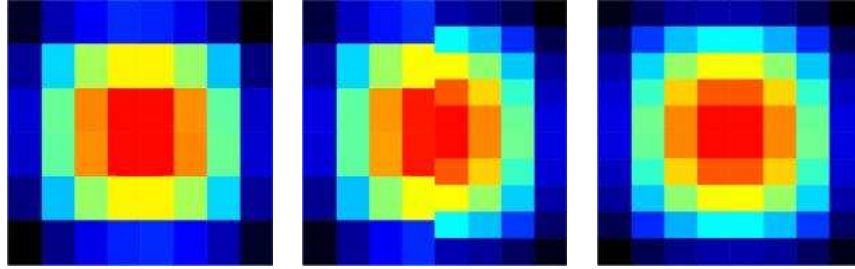


Figure 2: The approximate solution for conforming and nonconforming meshes. Left: conforming  $8 \times 6$  mesh, centre: nonconforming  $4 \times 6, 4 \times 10$  mesh, right: conforming  $10 \times 10$ .

The permeability tensor  $\Lambda$  is heterogeneous and isotropic, given by  $\Lambda(\mathbf{x}) = \lambda(\mathbf{x})\text{Id}$ , with  $\lambda(\mathbf{x}) = 1$  for a.e.  $x \in \Omega_1 \cup \Omega_3$  and  $\lambda(\mathbf{x}) = 10^{-2}$  for a.e.  $x \in \Omega_2$ . Note that the isolines of the exact solution are parallel to the boundaries of the sub-domain, and that the tangential component of the gradient is 0. We use the meshes depicted in Figure 3. Mesh 3 (containing  $10 \times 25$  control volumes) is obtained from Mesh 1 by the addition of two layers of very thin control volumes around each of the two lines of discontinuity of  $\Lambda$ : because of the very low thickness of these layers, equal to  $1/10000$ , the picture representing Mesh 3 is not different from that of Mesh 1.

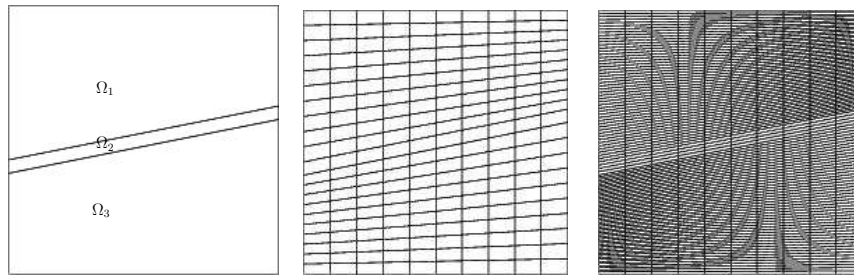


Figure 3: Domain and meshes used for the tilted barrier test: mesh 1 ( $10 \times 21$  centre), mesh 2 ( $10 \times 100$  right)

We get the following results for the approximations of the four fluxes at the boundary.

		nb. unknowns	matrix size	$x = 0$	$x = 1$	$y = 0$	$y = 1$
analytical				-0.2	0.2	1.	-1.
$\mathcal{B} = \mathcal{E}_{\text{int}}$	mesh 1	210	2424	-1.17	1.17	3.51	-3.51
	mesh 2	1000	11904	-0.237	0.237	1.104	-1.104
	mesh 3	250	2904	-0.208	0.208	1.02	-1.02
SUSHI-NP	mesh 1	239	2583	-0.2	0.2	1.	-1.
	mesh 2	1020	12036	-0.2	0.2	1.	-1.
HFV	mesh 1	599	4311	-0.2	0.2	1.	-1.
	mesh 2	2890	21138	-0.2	0.2	1.	-1.

Note that the values of the numerical solution given by the pure hybrid (HFV) and composite (SUSHI-NP) schemes are equal to those of the analytical solution (this holds under the only condition that the interfaces located on the lines  $\varphi_i(x, y) = 0$ ,  $i = 1, 2$ , are not included in  $\mathcal{B}$ , and that, for all  $\sigma \in \mathcal{B}$ , all  $K \in \mathcal{M}$  with  $\beta_\sigma^K \neq 0$  are included in the same sub-domain  $\Omega_i$ ). Note that Mesh 3, which leads to acceptable results for the computation of the fluxes, is not well suited for such a coupled problem, because of too small control volume measures. Hence SUSHI on Mesh 1 appears to be the most suitable method for this problem.

A satisfying natural choice (SUSHI-NP in the above results) is thus to match  $\mathcal{H}$  with the discontinuities of  $\Lambda$ . It is sometimes interesting to choose another set  $\mathcal{B}$ . This is for instance the case for the numerical locking problem for which the choice  $\mathcal{B} = \emptyset$  is best even though the diffusion tensor is homogeneous [20].

It is also sometimes interesting to replace the stabilisation coefficient  $\sqrt{d}$  in (25) by some other coefficient  $\alpha > 0$ . This is the case for instance for very distorted meshes or singular problems, in order to maintain the positivity of the unknown. The coefficient  $\alpha$  is taken to be greater than  $\sqrt{d}$ . The approximate solution remains positive, but the  $L^2$  norm of the error is generally larger. We refer to [20] for such experiments.

## 4 Convergence of the scheme

Let us first introduce some notations related to the mesh. Let  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P})$  be a discretisation of  $\Omega$  in the sense of Definition 2.1. The size of the discretisation  $\mathcal{D}$  is defined by:

$$h_{\mathcal{D}} = \sup\{h_K, K \in \mathcal{M}\},$$

and the regularity of the mesh by:

$$\theta_{\mathcal{D}} = \max \left( \max_{\sigma \in \mathcal{E}_{\text{int}}, K, L \in \mathcal{M}_\sigma} \frac{d_{K,\sigma}}{d_{L,\sigma}}, \max_{K \in \mathcal{M}, \sigma \in \mathcal{E}_K} \frac{h_K}{d_{K,\sigma}} \right). \quad (39)$$

For a given set  $\mathcal{B} \subset \mathcal{E}_{\text{int}}$  and for a given family  $(\beta_\sigma^K)_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_{\text{int}}}}$  satisfying property (15), we introduce a measure of the resulting regularity by

$$\theta_{\mathcal{D}, \mathcal{B}} = \max \left( \theta_{\mathcal{D}}, \max_{K \in \mathcal{M}, \sigma \in \mathcal{E}_K \cap \mathcal{B}} \frac{\sum_{L \in \mathcal{M}} |\beta_\sigma^L| |\mathbf{x}_L - \mathbf{x}_\sigma|^2}{h_K^2} \right). \quad (40)$$

**Remark 4.1** *Note that, for any mesh, it is easy to choose the family  $(\beta_\sigma^K)_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_{\text{int}}}}$  so that  $\theta_{\mathcal{D}, \mathcal{B}}$  remains small. It suffices to express  $\mathbf{x}_\sigma$  as the barycentre of  $d + 1$  points  $\mathbf{x}_L$  (which is always possible), for  $L$  sufficiently close to  $K$ , so that  $\mathbf{x}_L - \mathbf{x}_\sigma$  is close to  $h_K$  when  $\beta_\sigma^K \neq 0$ . Note also that in fact, it would be sufficient to have  $h_K^\eta$  with  $\eta > 1$  instead of  $h_K^2$  in (40) thus allowing the use of farther points.*

Remark that, thanks to the assumption that  $K$  is  $\mathbf{x}_K$ -star-shaped, the following property holds:

$$\sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} = d |K| \quad \forall K \in \mathcal{M}. \quad (41)$$

The space  $X_{\mathcal{D}}$  defined in (10) is equipped with the following semi-norm:

$$\forall v \in X_{\mathcal{D}}, |v|_X^2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (v_\sigma - v_K)^2, \quad (42)$$

which is a norm on the spaces  $X_{\mathcal{D},0}$  and  $X_{\mathcal{D},\mathcal{B}}$  respectively defined by (11) and (17).

Let  $H_{\mathcal{M}}(\Omega) \subset L^2(\Omega)$  be the set of piece-wise constant functions on the control volumes of the mesh  $\mathcal{M}$ . We then denote, for all  $v \in H_{\mathcal{M}}(\Omega)$  and for all  $\sigma \in \mathcal{E}_{\text{int}}$  with  $\mathcal{M}_\sigma = \{K, L\}$ ,  $D_\sigma v = |v_K - v_L|$  and  $d_\sigma = d_{K,\sigma} + d_{L,\sigma}$ , and for all  $\sigma \in \mathcal{E}_{\text{ext}}$  with  $\mathcal{M}_\sigma = \{K\}$ , we denote  $D_\sigma v = |v_K|$  and  $d_\sigma = d_{K,\sigma}$ . We then define the following norm:

$$\forall v \in H_{\mathcal{M}}(\Omega), \|v\|_{1,2,\mathcal{M}} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{D_\sigma v}{d_\sigma} \right)^2 = \sum_{\sigma \in \mathcal{E}} |\sigma| \frac{(D_\sigma v)^2}{d_\sigma}. \quad (43)$$

(Note that this norm is also defined by (74) in Lemma 5.2, setting  $p = 2$ ).

For all  $v \in X_{\mathcal{D}}$ , we denote by  $\Pi_{\mathcal{M}} v \in H_{\mathcal{M}}(\Omega)$  the piece-wise function from  $\Omega$  to  $\mathbb{R}$  defined by  $\Pi_{\mathcal{M}} v(\mathbf{x}) = v_K$  for a.e.  $\mathbf{x} \in K$ , for all  $K \in \mathcal{M}$ . Using the Cauchy-Schwarz inequality, we have for all  $\sigma \in \mathcal{E}_{\text{int}}$  with  $\mathcal{M}_\sigma = \{K, L\}$ ,

$$\frac{(v_K - v_L)^2}{d_\sigma} \leq \frac{(v_K - v_\sigma)^2}{d_{K,\sigma}} + \frac{(v_\sigma - v_L)^2}{d_{L,\sigma}} \quad \forall v \in X_{\mathcal{D}},$$

which leads to the relation

$$\|\Pi_{\mathcal{M}} v\|_{1,2,\mathcal{M}}^2 \leq |v|_X^2 \quad \forall v \in X_{\mathcal{D},0}. \quad (44)$$

For all  $\varphi \in C(\Omega, \mathbb{R})$ , we denote by  $P_{\mathcal{D}} \varphi$  the element of  $X_{\mathcal{D}}$  defined by  $((\varphi(\mathbf{x}_K))_{K \in \mathcal{M}}, (\varphi(\mathbf{x}_\sigma))_{\sigma \in \mathcal{E}})$ , by  $P_{\mathcal{D},\mathcal{B}} \varphi$  the element  $v \in X_{\mathcal{D},\mathcal{B}}$  such that  $v_K = \varphi(\mathbf{x}_K)$  for all  $K \in \mathcal{M}$ ,  $v_\sigma = 0$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ ,  $v_\sigma = \sum_{K \in \mathcal{M}} \beta_\sigma^K \varphi(\mathbf{x}_K)$  for all  $\sigma \in \mathcal{B}$  and  $v_\sigma = \varphi(\mathbf{x}_\sigma)$  for all  $\sigma \in \mathcal{H}$ .

We denote by  $P_{\mathcal{M}} \varphi \in H_{\mathcal{M}}(\Omega)$  the function such that  $P_{\mathcal{M}} \varphi(\mathbf{x}) = \varphi(\mathbf{x}_K)$  for a.e.  $\mathbf{x} \in K$ , for all  $K \in \mathcal{M}$  (we then have  $P_{\mathcal{M}} \varphi = \Pi_{\mathcal{M}} P_{\mathcal{D}} \varphi = \Pi_{\mathcal{M}} P_{\mathcal{D},\mathcal{B}} \varphi$ ).

The following lemma provides an equivalence property between the  $L^2$ -norm of the discrete gradient, defined by (22)-(26) and the norm  $|\cdot|_X$ .

**Lemma 4.1** *Let  $\mathcal{D}$  be a discretisation of  $\Omega$  in the sense of Definition 2.1, and let  $\theta \geq \theta_{\mathcal{D}}$  be given (where  $\theta_{\mathcal{D}}$  is defined by (39)). Then there exists  $C_1 > 0$  and  $C_2 > 0$  only depending on  $\theta$  and  $d$  such that:*

$$C_1 |u|_X \leq \|\nabla_{\mathcal{D}} u\|_{L^2(\Omega)} \leq C_2 |u|_X \quad \forall u \in X_{\mathcal{D}}, \quad (45)$$

where  $\nabla_{\mathcal{D}}$  is defined by (22)-(26).

PROOF. By definition,

$$\|\nabla_{\mathcal{D}} u\|_{L^2(\Omega)^d}^2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} |\nabla_{K,\sigma} u|^2.$$

Therefore, using property (27),

$$\|\nabla_{\mathcal{D}} u\|_{L^2(\Omega)^d}^2 = \sum_{K \in \mathcal{M}} \left( |K| |\nabla_K u|^2 + \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} (R_{K,\sigma} u)^2 \right). \quad (46)$$

Let us now notice that the following inequality holds:

$$(a - b)^2 \geq \frac{\lambda}{1 + \lambda} a^2 - \lambda b^2 \quad \forall a, b \in \mathbb{R}, \forall \lambda > -1. \quad (47)$$

We apply this inequality to  $(R_{K,\sigma}u)^2$  for some  $\lambda > 0$  and obtain

$$(R_{K,\sigma}u)^2 \geq \frac{\lambda d}{1 + \lambda} \left( \frac{u_\sigma - u_K}{d_{K,\sigma}} \right)^2 - \lambda d |\nabla_K u|^2 \left( \frac{|\mathbf{x}_\sigma - \mathbf{x}_K|}{d_{K,\sigma}} \right)^2. \quad (48)$$

This leads to

$$\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} (R_{K,\sigma}u)^2 \geq \frac{\lambda}{1 + \lambda} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{u_\sigma - u_K}{d_{K,\sigma}} \right)^2 - \lambda |K| d |\nabla_K u|^2 \theta^2.$$

Choosing  $\lambda = \frac{1}{d\theta^2}$ , we get that

$$\|\nabla_{\mathcal{D}} u\|_{(L^2(\Omega))^d}^2 \geq \frac{\lambda}{1 + \lambda} |u|_X^2,$$

which shows the left inequality of (45).

Let us now prove the right inequality. On one hand, using the definition (22) of  $\nabla_K u$  and (41), the Cauchy–Schwarz inequality leads to

$$|\nabla_K u|^2 \leq \frac{1}{|K|^2} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_\sigma - u_K)^2 \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} = \frac{d}{|K|} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_\sigma - u_K)^2. \quad (49)$$

On the other hand, by definition (25), and thanks to the definition of the regularity of the mesh (39), we have

$$(R_{K,\sigma}u)^2 \leq 2d \left( \left( \frac{u_\sigma - u_K}{d_{K,\sigma}} \right)^2 + |\nabla_K u|^2 \frac{|\mathbf{x}_\sigma - \mathbf{x}_K|}{d_{K,\sigma}} \right) \leq 2d \left( \left( \frac{u_\sigma - u_K}{d_{K,\sigma}} \right)^2 + \theta^2 |\nabla_K u|^2 \right). \quad (50)$$

From (46), (49) and (50), we conclude that the right inequality of (45) holds.  $\square$

We may now state a weak compactness result for the discrete gradient.

**Lemma 4.2 (Weak discrete  $H^1$  compactness)** *Let  $\mathcal{F}$  be a family of discretisations in the sense of Definition 2.1 such that there exists  $\theta > 0$  with  $\theta \geq \theta_{\mathcal{D}}$  for all  $\mathcal{D} \in \mathcal{F}$ . Let  $(u_{\mathcal{D}})_{\mathcal{D} \in \mathcal{F}}$  be a family of functions, such that:*

- $u_{\mathcal{D}} \in X_{\mathcal{D},0}$  for all  $\mathcal{D} \in \mathcal{F}$ ,
- there exists  $C > 0$  with  $|u_{\mathcal{D}}|_X \leq C$  for all  $\mathcal{D} \in \mathcal{F}$ ,
- there exists  $u \in L^2(\Omega)$  with  $\lim_{h_{\mathcal{D}} \rightarrow 0} \|\Pi_{\mathcal{M}} u_{\mathcal{D}} - u\|_{L^2(\Omega)} = 0$ .

*Then,  $u \in H_0^1(\Omega)$  and  $\nabla_{\mathcal{D}} u_{\mathcal{D}}$  weakly converge in  $L^2(\Omega)^d$  to  $\nabla u$  as  $h_{\mathcal{D}} \rightarrow 0$ , where the operator  $\nabla_{\mathcal{D}}$  is defined by (22)-(26).*

**PROOF.** Let us prolong  $\Pi_{\mathcal{M}} u_{\mathcal{D}}$  and  $\nabla_{\mathcal{D}} u_{\mathcal{D}}$  by 0 outside of  $\Omega$ . Thanks to Lemma 4.1, up to a subsequence, there exists some function  $\mathbf{G} \in L^2(\mathbb{R}^d)^d$  such that  $\nabla_{\mathcal{D}} u_{\mathcal{D}}$  weakly converges in  $L^2(\mathbb{R}^d)^d$  to  $\mathbf{G}$  as  $h_{\mathcal{D}} \rightarrow 0$ . Let us show that  $\mathbf{G} = \nabla u$ . Let  $\psi \in C_c^\infty(\mathbb{R}^d)^d$  be given. Let us consider the term  $T_1^{\mathcal{D}}$  defined by

$$T_1^{\mathcal{D}} = \int_{\mathbb{R}^d} \nabla_{\mathcal{D}} u_{\mathcal{D}}(\mathbf{x}) \cdot \psi(\mathbf{x}) d\mathbf{x}.$$

We get that  $T_1^{\mathcal{D}} = T_2^{\mathcal{D}} + T_3^{\mathcal{D}}$ , with

$$T_2^{\mathcal{D}} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (u_\sigma - u_K) \mathbf{n}_{K,\sigma} \cdot \boldsymbol{\psi}_K, \text{ with } \boldsymbol{\psi}_K = \frac{1}{|K|} \int_K \boldsymbol{\psi}(\mathbf{x}) d\mathbf{x},$$

and

$$T_3^{\mathcal{D}} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} R_{K,\sigma} u \mathbf{n}_{K,\sigma} \cdot \int_{D_{K,\sigma}} \boldsymbol{\psi}(\mathbf{x}) d\mathbf{x}.$$

We compare  $T_2^{\mathcal{D}}$  with  $T_4^{\mathcal{D}}$  defined by

$$T_4^{\mathcal{D}} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (u_\sigma - u_K) \mathbf{n}_{K,\sigma} \cdot \boldsymbol{\psi}_\sigma,$$

with

$$\boldsymbol{\psi}_\sigma = \frac{1}{|\sigma|} \int_\sigma \boldsymbol{\psi}(\mathbf{x}) d\gamma(\mathbf{x}).$$

We get that

$$(T_2^{\mathcal{D}} - T_4^{\mathcal{D}})^2 \leq \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (u_\sigma - u_K)^2 \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} |\boldsymbol{\psi}_K - \boldsymbol{\psi}_\sigma|^2,$$

which leads to  $\lim_{h_{\mathcal{D}} \rightarrow 0} (T_2^{\mathcal{D}} - T_4^{\mathcal{D}}) = 0$ .

Since

$$T_4^{\mathcal{D}} = - \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| u_K \mathbf{n}_{K,\sigma} \cdot \boldsymbol{\psi}_\sigma = - \int_{\mathbb{R}^d} \Pi_{\mathcal{M}} u_{\mathcal{D}}(\mathbf{x}) \operatorname{div} \boldsymbol{\psi}(\mathbf{x}) d\mathbf{x},$$

we get that  $\lim_{h_{\mathcal{D}} \rightarrow 0} T_4^{\mathcal{D}} = - \int_{\mathbb{R}^d} u(\mathbf{x}) \operatorname{div} \boldsymbol{\psi}(\mathbf{x}) d\mathbf{x}$ . Let us now turn to the study of  $T_3^{\mathcal{D}}$ . Noting again that (27) holds, we have:

$$T_3^{\mathcal{D}} = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} R_{K,\sigma} u \mathbf{n}_{K,\sigma} \cdot \int_{D_{K,\sigma}} (\boldsymbol{\psi}(\mathbf{x}) - \boldsymbol{\psi}_K) d\mathbf{x}.$$

Since  $\boldsymbol{\psi}$  is a regular function, there exists  $C_{\boldsymbol{\psi}}$  only depending on  $\boldsymbol{\psi}$  such that  $|\int_{D_{K,\sigma}} (\boldsymbol{\psi}(\mathbf{x}) - \boldsymbol{\psi}_K) d\mathbf{x}| \leq C_{\boldsymbol{\psi}} h_{\mathcal{D}} \frac{|\sigma| d_{K,\sigma}}{d}$ . From (50) and the Cauchy-Schwarz inequality, we thus get:

$$\lim_{h_{\mathcal{D}} \rightarrow 0} T_3^{\mathcal{D}} = 0.$$

This proves that the function  $\mathbf{G} \in L^2(\mathbb{R}^d)^d$  is a.e. equal to  $\nabla u$  in  $\mathbb{R}^d$ . Since  $u = 0$  outside of  $\Omega$ , we get that  $u \in H_0^1(\Omega)$ , and the uniqueness of the limit implies that the whole family  $\nabla_{\mathcal{D}} u_{\mathcal{D}}$  weakly converges in  $L^2(\mathbb{R}^d)^d$  to  $\nabla u$  as  $h_{\mathcal{D}} \rightarrow 0$ .

□

Note that the proof that  $u \in H_0^1(\Omega)$  also results from (44), which allows us to apply Lemma 5.7 of the Appendix in the particular case  $p = 2$ . Let us also remark that several discrete gradients could be chosen, which satisfy the weak compactness property (see for instance the proof of Lemma 5.7). However, we emphasise that the choice of the specific gradient (22) also stems from coercivity and consistency issues. Let us now state the discrete gradient consistency property.

**Lemma 4.3 (Discrete gradient consistency)** *Let  $\mathcal{D}$  be a discretisation of  $\Omega$  in the sense of Definition 2.1, and let  $\theta \geq \theta_{\mathcal{D}}$  be given. Then, for any function  $\varphi \in C^2(\overline{\Omega})$ , there exists  $C_3$  only depending on  $d$ ,  $\theta$  and  $\varphi$  such that:*

$$\|\nabla_{\mathcal{D}} P_{\mathcal{D}} \varphi - \nabla \varphi\|_{(L^\infty(\Omega))^d} \leq C_3 h_{\mathcal{D}}, \quad (51)$$

where  $\nabla_{\mathcal{D}}$  is defined by (22)-(26).

PROOF. From definitions (26) and (24) we get

$$|\nabla_{K,\sigma} P_{\mathcal{D}}\varphi - \nabla\varphi(\mathbf{x}_K)| \leq |\nabla_K P_{\mathcal{D}}\varphi - \nabla\varphi(\mathbf{x}_K)| + |R_{K,\sigma} P_{\mathcal{D}}\varphi|.$$

From (22), we have, for any  $K \in \mathcal{M}$ ,

$$\begin{aligned} \nabla_K P_{\mathcal{D}}\varphi &= \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (\varphi(\mathbf{x}_\sigma) - \varphi(\mathbf{x}_K)) \mathbf{n}_{K,\sigma} \\ &= \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (\nabla\varphi(\mathbf{x}_K) \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) + h_K^2 \rho_{K,\sigma}) \mathbf{n}_{K,\sigma}, \end{aligned}$$

where  $|\rho_{K,\sigma}| \leq C_\varphi$  with  $C_\varphi$  only depending on  $\varphi$ . Thanks to (23) and to the regularity of the mesh, we get

$$|\nabla_K P_{\mathcal{D}}\varphi - \nabla\varphi(\mathbf{x}_K)| \leq \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| h_K^2 |\rho_{K,\sigma}| \leq h_K d C_\varphi \theta.$$

From this last inequality, using Definition 25, we get

$$\begin{aligned} |R_{K,\sigma} P_{\mathcal{D}}\varphi| &= \frac{\sqrt{d}}{d_{K,\sigma}} |\varphi(\mathbf{x}_\sigma) - \varphi(\mathbf{x}_K) - \nabla_K P_{\mathcal{D}}\varphi \cdot (\mathbf{x}_\sigma - \mathbf{x}_K)| \\ &\leq \frac{\sqrt{d}}{d_{K,\sigma}} (h_K^2 \rho_{K,\sigma} + h_K^2 d C_\varphi \theta) \\ &\leq \sqrt{d} \theta (h_K C_\varphi + h_K d C_\varphi \theta), \end{aligned}$$

which concludes the proof.  $\square$

We now give the abstract properties of the discrete fluxes, which are necessary to prove the convergence of the general scheme (18), (13), and then prove that the fluxes that we constructed in Section 2.4 indeed satisfy these properties.

**Definition 4.1 (Continuous, coercive, consistent and symmetric families of fluxes)**

Let  $\mathcal{F}$  be a family of discretisations in the sense of definition 2.1. For  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P}) \in \mathcal{F}$ ,  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}$ , we denote by  $F_{K,\sigma}^{\mathcal{D}}$  a linear mapping from  $X_{\mathcal{D}}$  to  $\mathbb{R}$ , and we denote by  $\Phi = ((F_{K,\sigma}^{\mathcal{D}})_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}}})_{\mathcal{D} \in \mathcal{F}}$ .

We consider the bilinear form defined by

$$\langle u, v \rangle_F = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}^{\mathcal{D}}(u) (v_K - v_\sigma) \quad \forall (u, v) \in X_{\mathcal{D}}^2. \quad (52)$$

The family of numerical fluxes  $\Phi$  is said to be continuous if there exists  $M > 0$  such that

$$\langle u, v \rangle_F \leq M |u|_X |v|_X \quad \forall (u, v) \in X_{\mathcal{D}}^2, \quad \forall \mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P}) \in \mathcal{F}. \quad (53)$$

The family of numerical fluxes  $\Phi$  is said to be coercive if there exists  $\alpha > 0$  such that

$$\alpha |u|_X^2 \leq \langle u, u \rangle_F \quad \forall u \in X_{\mathcal{D}} \quad \forall \mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P}) \in \mathcal{F}. \quad (54)$$

The family of numerical fluxes  $\Phi$  is said to be consistent (with Problem (1)–(2)) if for any family  $(u_{\mathcal{D}})_{\mathcal{D} \in \mathcal{F}}$  satisfying:

- $u_{\mathcal{D}} \in X_{\mathcal{D},0}$  for all  $\mathcal{D} \in \mathcal{F}$ ,
- there exists  $C > 0$  with  $|u_{\mathcal{D}}|_X \leq C$  for all  $\mathcal{D} \in \mathcal{F}$ ,
- there exists  $u \in L^2(\Omega)$  with  $\lim_{h_{\mathcal{D}} \rightarrow 0} \|\Pi_{\mathcal{M}} u_{\mathcal{D}} - u\|_{L^2(\Omega)} = 0$  (recall that, from Lemma 5.7, we get that  $u \in H_0^1(\Omega)$ ),

then

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle u_{\mathcal{D}}, P_{\mathcal{D}}\varphi \rangle_F = \int_{\Omega} \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) \cdot \nabla u(\mathbf{x}) d\mathbf{x} \quad \forall \varphi \in C_c^{\infty}(\Omega). \quad (55)$$

Finally the family of numerical fluxes  $\Phi$  is said to be symmetric if

$$\langle u, v \rangle_F = \langle v, u \rangle_F \quad \forall (u, v) \in X_{\mathcal{D}}^2, \quad \forall \mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P}) \in \mathcal{F}.$$

We now show that the family of fluxes defined by (28)-(31) satisfies the definition of a consistent, coercive and symmetric family of fluxes. Recall that the SUSHI scheme (37) is studied numerically in Section 3 with this choice for the family of fluxes.

**Lemma 4.4 (Flux properties)** *Let  $\mathcal{F}$  be a family of discretisations in the sense of Definition 2.1. We assume that there exists  $\theta > 0$  with*

$$\theta_{\mathcal{D}} \leq \theta \quad \forall \mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P}) \in \mathcal{F}, \quad (56)$$

where  $\theta_{\mathcal{D}}$  is defined by (39). Let  $\Phi = ((F_{K,\sigma}^{\mathcal{D}})_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_K}})_{\mathcal{D} \in \mathcal{F}}$  be the family of fluxes defined by (28)-(31). Then, the family  $\Phi$  is a continuous, coercive, consistent and symmetric family of numerical fluxes in the sense of Definition 4.1.

PROOF. Since the family of fluxes is defined by (28)-(31), it satisfies (21), and therefore we have:

$$\langle u, v \rangle_F = \int_{\Omega} \nabla_{\mathcal{D}} u(\mathbf{x}) \cdot \Lambda(\mathbf{x}) \nabla_{\mathcal{D}} v(\mathbf{x}) d\mathbf{x} \quad \forall u, v \in X_{\mathcal{D}}.$$

Hence the property  $\langle u, v \rangle_F = \langle v, u \rangle_F$  holds. The continuity and coercivity of the family  $\Phi$  result from Lemma 4.1 and the properties of  $\Lambda$ , which give:  $\langle u, v \rangle_F \leq \bar{\lambda} \|\nabla_{\mathcal{D}} u\|_{L^2(\Omega)} \|\nabla_{\mathcal{D}} v\|_{L^2(\Omega)}$  and  $\langle u, u \rangle_F \geq \underline{\lambda} \|\nabla_{\mathcal{D}} u\|_{L^2(\Omega)}^2$  for any  $u, v \in X_{\mathcal{D}}$ . The consistency results from the weak and strong convergence properties in Lemmas 4.2 and 4.3, which give  $\nabla_{\mathcal{D}} u_{\mathcal{D}} \rightarrow \nabla u$  weakly in  $L^2(\Omega)$  and  $\nabla_{\mathcal{D}} P_{\mathcal{D}}\varphi \rightarrow \nabla \varphi$  in  $L^2(\Omega)$  as the mesh size tends to 0.  $\square$

**Theorem 4.1 (Convergence)** *Let  $\mathcal{F}$  be a family of discretisations in the sense of Definition 2.1, for any  $\mathcal{D} \in \mathcal{F}$ , let  $\mathcal{B} \subset \mathcal{E}_{\text{int}}$  and  $(\beta_{\sigma}^K)_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}_{\text{int}}}}$  satisfying (15). Assume that there exists  $\theta > 0$  such that  $\theta_{\mathcal{D},\mathcal{B}} \leq \theta$ , for all  $\mathcal{D} \in \mathcal{F}$ , where  $\theta_{\mathcal{D},\mathcal{B}}$  is defined by (40). Let  $\Phi = ((F_{K,\sigma}^{\mathcal{D}})_{\substack{K \in \mathcal{M} \\ \sigma \in \mathcal{E}}})_{\mathcal{D} \in \mathcal{F}}$  be a continuous, coercive and symmetric and consistent family of numerical fluxes in the sense of Definition 4.1. Let  $(u_{\mathcal{D}})_{\mathcal{D} \in \mathcal{F}}$  be the family of functions satisfying (18) for all  $\mathcal{D} \in \mathcal{F}$ . Then  $\Pi_{\mathcal{M}} u_{\mathcal{D}}$  converges in  $L^2(\Omega)$  to the unique solution  $u$  of (6) as  $h_{\mathcal{D}} \rightarrow 0$ . Moreover  $\nabla_{\mathcal{D}} u_{\mathcal{D}}$  converges to  $\nabla u$  in  $L^2(\Omega)^d$  as  $h_{\mathcal{D}} \rightarrow 0$ .*

PROOF. Letting  $v = u_{\mathcal{D}}$  in (18) and applying the Cauchy-Schwarz inequality yields

$$\langle u_{\mathcal{D}}, u_{\mathcal{D}} \rangle_F = \int_{\Omega} f(\mathbf{x}) \Pi_{\mathcal{M}} u_{\mathcal{D}}(\mathbf{x}) d\mathbf{x} \leq \|f\|_{L^2(\Omega)} \|\Pi_{\mathcal{M}} u_{\mathcal{D}}\|_{L^2(\Omega)}.$$

We apply the Sobolev inequality (77) with  $p = 2$ , which gives in this case

$$\|\Pi_{\mathcal{M}} u_{\mathcal{D}}\|_{L^2(\Omega)} \leq C_4 \|\Pi_{\mathcal{M}} u_{\mathcal{D}}\|_{1,2,\mathcal{M}}.$$

Using (44) and the consistency of the family  $\Phi$  of fluxes, we then have

$$\alpha \|\Pi_{\mathcal{M}} u_{\mathcal{D}}\|_X^2 \leq C_4 \|f\|_{L^2(\Omega)} |u_{\mathcal{D}}|_X.$$

This leads to the inequality

$$\|u_{\mathcal{D}}\|_{1,2,\mathcal{M}} \leq |u_{\mathcal{D}}|_X \leq \frac{C_4}{\alpha} \|f\|_{L^2(\Omega)}. \quad (57)$$

Thanks to Lemma 5.7, we get the existence of  $u \in H_0^1(\Omega)$ , and of a subfamily extracted from  $\mathcal{F}$ , such that  $\|\Pi_{\mathcal{M}}u_{\mathcal{D}} - u\|_{L^2(\Omega)}$  tends to 0 as  $h_{\mathcal{D}} \rightarrow 0$ . For a given  $\varphi \in C_c^\infty(\Omega)$ , let us take  $v = P_{\mathcal{D},\mathcal{B}}\varphi$  in (18) (recall that  $P_{\mathcal{D},\mathcal{B}}\varphi \in X_{\mathcal{D},\mathcal{B}}$ ). We get

$$\langle u_{\mathcal{D}}, P_{\mathcal{D},\mathcal{B}}\varphi \rangle_F = \int_{\Omega} f(\mathbf{x}) P_{\mathcal{M}}\varphi(\mathbf{x}) d\mathbf{x}.$$

Let us remark that, thanks to the continuity of the family  $\Phi$  of fluxes, we have

$$\langle u_{\mathcal{D}}, P_{\mathcal{D},\mathcal{B}}\varphi - P_{\mathcal{D}}\varphi \rangle_F \leq M \frac{C_{13}}{\alpha} \|f\|_{L^2(\Omega)} |P_{\mathcal{D},\mathcal{B}}\varphi - P_{\mathcal{D}}\varphi|_X.$$

Thanks to (15) and (40), we get the existence of  $C_\varphi$  only depending on  $\varphi$  (through its second order partial derivatives) such that, for all  $K \in \mathcal{M}$  and all  $\sigma \in \mathcal{B} \cap \mathcal{E}_K$ ,

$$\left| \sum_{L \in \mathcal{M}} \beta_\sigma^L \varphi(\mathbf{x}_L) - \varphi(\mathbf{x}_\sigma) \right| \leq \sum_{L \in \mathcal{M}} |\beta_\sigma^L| |\mathbf{x}_L - \mathbf{x}_\sigma|^2 C_\varphi \leq \theta_{\mathcal{D},\mathcal{B}} C_\varphi h_K^2. \quad (58)$$

We can then deduce

$$\lim_{h_{\mathcal{D}} \rightarrow 0} |P_{\mathcal{D},\mathcal{B}}\varphi - P_{\mathcal{D}}\varphi|_X = 0. \quad (59)$$

Thanks to the  $\mathcal{F}$ -extracted subfamily properties, we may apply the consistency hypothesis on the family  $\Phi$  of fluxes, which gives

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle u_{\mathcal{D}}, P_{\mathcal{D}}\varphi \rangle_F = \int_{\Omega} \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) \cdot \nabla u(\mathbf{x}) d\mathbf{x}.$$

Gathering the two results above leads to

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle u_{\mathcal{D}}, P_{\mathcal{D},\mathcal{B}}\varphi \rangle_F = \int_{\Omega} \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) \cdot \nabla u(\mathbf{x}) d\mathbf{x},$$

which concludes the proof of the following equality

$$\int_{\Omega} \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) \cdot \nabla u(\mathbf{x}) d\mathbf{x} = \int_{\Omega} f(\mathbf{x}) \varphi(\mathbf{x}) d\mathbf{x}.$$

Therefore,  $u$  is the unique solution of (6), and we get that the whole family  $(u_{\mathcal{D}})_{\mathcal{D} \in \mathcal{F}}$  converges to  $u$  as  $h_{\mathcal{D}} \rightarrow 0$ .

Let us now prove the second part of the theorem.

Let  $\varphi \in C_c^\infty(\Omega)$  be given (this function is devoted to approximate  $u$  in  $H_0^1(\Omega)$ ). Thanks to the Cauchy-Schwarz inequality, we have

$$\int_{\Omega} |\nabla_{\mathcal{D}} u_{\mathcal{D}}(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x} \leq 3 (T_5^{\mathcal{D}} + T_6^{\mathcal{D}} + T_7),$$

with  $T_5^{\mathcal{D}} = \int_{\Omega} |\nabla_{\mathcal{D}} u_{\mathcal{D}}(\mathbf{x}) - \nabla_{\mathcal{D}} P_{\mathcal{D}}\varphi(\mathbf{x})|^2 d\mathbf{x}$ ,  $T_6^{\mathcal{D}} = \int_{\Omega} |\nabla_{\mathcal{D}} P_{\mathcal{D}}\varphi(\mathbf{x}) - \nabla \varphi(\mathbf{x})|^2 d\mathbf{x}$ , and  $T_7 = \int_{\Omega} |\nabla \varphi(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x}$ . Thanks to Lemma 4.3, we have  $\lim_{h_{\mathcal{D}} \rightarrow 0} T_6^{\mathcal{D}} = 0$ .

Thanks to Lemma 4.1 and to the coercivity of the family of fluxes, there exists  $C_5$  such that

$$\|\nabla_{\mathcal{D}} v\|_{L^2(\Omega)^d}^2 \leq C_2^2 |v|_X^2 \leq C_5 \langle v, v \rangle_F \quad \forall v \in X_{\mathcal{D}},$$

with  $C_5 = \frac{C_2^2}{\alpha}$ . Taking  $v = u_{\mathcal{D}} - P_{\mathcal{D}}\varphi$ , we have

$$T_5^{\mathcal{D}} \leq C_5 (\langle u_{\mathcal{D}}, u_{\mathcal{D}} \rangle_F - 2 \langle u_{\mathcal{D}}, P_{\mathcal{D}}\varphi \rangle_F + \langle P_{\mathcal{D}}\varphi, P_{\mathcal{D}}\varphi \rangle_F).$$

By Theorem 4.1 and thanks to and consistency of the family of fluxes, we get

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle u_{\mathcal{D}}, P_{\mathcal{D}}\varphi \rangle_F = \int_{\Omega} \nabla u(\mathbf{x}) \cdot \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) d\mathbf{x}.$$

The sequence  $\|P_{\mathcal{D}}\varphi\|_X$  is bounded; using the regularity of  $\varphi$ , the regularity hypotheses of the family of discretisations, together with the consistency of the family of fluxes implies that

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle P_{\mathcal{D}}\varphi, P_{\mathcal{D}}\varphi \rangle_F = \int_{\Omega} \nabla \varphi(\mathbf{x}) \cdot \Lambda(\mathbf{x}) \nabla \varphi(\mathbf{x}) d\mathbf{x}.$$

Remarking that passing to the limit  $h_{\mathcal{D}} \rightarrow 0$  in (18) with  $v = u_{\mathcal{D}}$  provides that  $\langle u_{\mathcal{D}}, u_{\mathcal{D}} \rangle_F$  converges to  $\int_{\Omega} \nabla u \cdot \Lambda \nabla u d\mathbf{x}$ , we get that

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \langle u_{\mathcal{D}} - P_{\mathcal{D}}\varphi, u_{\mathcal{D}} - P_{\mathcal{D}}\varphi \rangle_F = \int_{\Omega} \nabla(u - \varphi) \cdot \Lambda \nabla(u - \varphi) d\mathbf{x} \leq \bar{\lambda} \int_{\Omega} |\nabla u - \nabla \varphi|^2 d\mathbf{x},$$

which yields

$$\limsup_{h_{\mathcal{D}} \rightarrow 0} T_5^{\mathcal{D}} \leq C_5 \bar{\lambda} \int_{\Omega} |\nabla u - \nabla \varphi|^2 d\mathbf{x}.$$

From the above results, we obtain that there exists  $C_6$ , independent of  $\mathcal{D}$ , such that

$$\int_{\Omega} |\nabla_{\mathcal{D}} u_{\mathcal{D}}(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x} \leq C_6 \int_{\Omega} |\nabla \varphi(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x} + T_8^{\mathcal{D}},$$

with (noting that  $\varphi$  is fixed)  $\lim_{h_{\mathcal{D}} \rightarrow 0} T_8^{\mathcal{D}} = 0$ . Let  $\varepsilon > 0$ ; we may choose  $\varphi$  such that  $\int_{\Omega} |\nabla \varphi(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x} \leq \varepsilon$ , and we may then choose  $h_{\mathcal{D}}$  small enough so that  $T_8^{\mathcal{D}} \leq \varepsilon$ . This completes the proof that

$$\lim_{h_{\mathcal{D}} \rightarrow 0} \int_{\Omega} |\nabla_{\mathcal{D}} u_{\mathcal{D}}(\mathbf{x}) - \nabla u(\mathbf{x})|^2 d\mathbf{x} = 0 \quad (60)$$

in the case of a general continuous, coercive, consistent and symmetric family of fluxes.  $\square$

Let us write an error estimate in the particular case  $\Lambda = \text{Id}$ , assuming a regular exact solution to (6).

**Theorem 4.2 (Error estimate, isotropic case)** *We consider the particular case  $\Lambda = \text{Id}$ , and we assume that the solution  $u \in H_0^1(\Omega)$  of (6) is in  $C^2(\bar{\Omega})$ . Let  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P})$  be a discretisation in the sense of Definition 2.1, let  $\mathcal{B} \subset \mathcal{E}_{\text{int}}$  be given, let  $\mathcal{B} = (\beta_{\sigma}^K)_{\sigma \in \mathcal{B}, K \in \mathcal{M}} \subset \mathbb{R}$  such that (15) holds, and let  $\theta \geq \theta_{\mathcal{D}, \mathcal{B}}$  be given (see (40)). Let  $(F_{K, \sigma})_{K \in \mathcal{M}, \sigma \in \mathcal{E}}$  be a family of linear mappings from  $X_{\mathcal{D}}$  to  $\mathbb{R}$ , such that there exists  $\alpha > 0$  with*

$$\alpha |v|_X^2 \leq \langle v, v \rangle_F \quad \forall v \in X_{\mathcal{D}}, \quad (61)$$

defining  $\langle \cdot, \cdot \rangle_F$  by (52). We denote by

$$E(u) = \left( \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{d_{K, \sigma}}{|\sigma|} \left( F_{K, \sigma}(P_{\mathcal{D}, \mathcal{B}} u) + \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K, \sigma} d\gamma(\mathbf{x}) \right)^2 \right)^{1/2}. \quad (62)$$

Then the solution  $u_{\mathcal{D}}$  of (18) satisfies that there exists  $C_7$ , only depends on  $\alpha$  and on  $\theta$ , such that

$$\|\Pi_{\mathcal{M}} u_{\mathcal{D}} - P_{\mathcal{M}} u\|_{L^2(\Omega)} \leq C_7 E(u), \quad (63)$$

and satisfies that there exists  $C_8$ , only depending on  $\alpha$ ,  $\theta$  and  $u$  such that

$$\|\nabla_{\mathcal{D}} u_{\mathcal{D}} - \nabla u\|_{L^2(\Omega)^d} \leq C_8 (E(u) + h_{\mathcal{D}}). \quad (64)$$

Moreover, in the particular case where  $(F_{K, \sigma})_{K \in \mathcal{M}, \sigma \in \mathcal{E}}$  is defined by (28)-(31), there exists  $C_9$ , only depending on  $\alpha$ ,  $\theta$  and  $u$ , such that

$$E(u) \leq C_9 h_{\mathcal{D}}. \quad (65)$$

**Remark 4.2 (Extensions of the error estimate)** Note also that the extension of Theorem 4.2 to the case  $u \in H^2(\Omega)$  is possible for  $d = 2$  or  $d = 3$ . However it would demand a rather longer and more technical proof and is not expected to provide more information on the link between accuracy and the regularity of the mesh than the result presented here. In the case of the pure hybrid scheme (HFV,  $\mathcal{B} = \emptyset$ ), an error estimate could however be obtained by assuming  $u$  piece-wise to be  $H^2$ . Such error estimates were also obtained for pure hybrid schemes of the mimetic type by using the tools of the mixed finite element theory (see e.g. [10]). If  $\mathcal{B} \neq \emptyset$ , one must furthermore assume that the barycentric formulae (14)-(15) or (19)-(20) are written with unknowns located in the same regularity zone, as explained in Remark 2.6. Nevertheless such error estimates are not possible for general  $L^\infty$  diffusion operators, since in such a case the maximal regularity of the continuous solution is  $H_0^1(\Omega)$ . Then, by interpolation, one may get some error estimates if the continuous solution is in  $H_0^1(\Omega) \cap H^s(\Omega)$  as in the classical finite element framework.

PROOF. Let  $v \in X_{\mathcal{D}}$ , since  $-\Delta u = f$ , we get:

$$-\sum_{K \in \mathcal{M}} v_K \int_K \Delta u(\mathbf{x}) d\mathbf{x} = \int_{\Omega} f(\mathbf{x}) \Pi_{\mathcal{M}} v(\mathbf{x}) d\mathbf{x}. \quad (66)$$

Thanks to the following equality (recall that  $u \in C^2(\bar{\Omega})$  and therefore  $\nabla u \cdot \mathbf{n}_{K,\sigma}$  is defined on each edge  $\sigma$ )

$$-\sum_{K \in \mathcal{M}} v_K \int_K \Delta u(\mathbf{x}) d\mathbf{x} = -\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} (v_K - v_\sigma) \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}),$$

we get that

$$\langle P_{\mathcal{D},\mathcal{B}} u, v \rangle_F = \int_{\Omega} f(\mathbf{x}) \Pi_{\mathcal{M}} v(\mathbf{x}) d\mathbf{x} + \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \left( F_{K,\sigma}^{\mathcal{D}}(P_{\mathcal{D},\mathcal{B}} u) + \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right) (v_K - v_\sigma).$$

Taking  $v = P_{\mathcal{D},\mathcal{B}} u - u_{\mathcal{D}} \in X_{\mathcal{D},\mathcal{B}}$  in this latter equality and using (66) we get

$$\langle v, v \rangle_F = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \left( F_{K,\sigma}^{\mathcal{D}}(P_{\mathcal{D},\mathcal{B}} u) + \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right) (v_K - v_\sigma),$$

which leads, using (61) and the Cauchy-Schwarz inequality, to

$$\alpha |v|_X \leq E(u). \quad (67)$$

Using (44) and the Sobolev inequality (77) with  $p = 2$  provides the conclusion of (63). Let us now prove (64). We have

$$\|\nabla_{\mathcal{D}} u_{\mathcal{D}} - \nabla u\|_{L^2(\Omega)^d} \leq \|\nabla_{\mathcal{D}} u_{\mathcal{D}} - \nabla_{\mathcal{D}} P_{\mathcal{D},\mathcal{B}} u\|_{L^2(\Omega)^d} + \|\nabla_{\mathcal{D}} P_{\mathcal{D},\mathcal{B}} u - \nabla u\|_{L^2(\Omega)^d}.$$

The bound of the first term in the above right-hand side is bounded thanks to Lemma 4.1 and (67). The inequality  $\|\nabla_{\mathcal{D}} P_{\mathcal{D},\mathcal{B}} u - \nabla u\|_{L^2(\Omega)^d} \leq C_{10} h_{\mathcal{D}}$  is obtained thanks to Lemma 4.3 and using a similar inequality to (58), replacing  $\varphi$  by  $u$ .

Let us now turn to the proof of (65) in the particular case where the family of fluxes is defined by (28)-(31). Indeed, we get in this case that, for all  $v \in X_{\mathcal{D}}$ ,

$$F_{K,\sigma}(v) = -\sum_{\sigma' \in \mathcal{E}_K} (\nabla_K v + R_{K,\sigma'} v \mathbf{n}_{K,\sigma'}) \cdot \frac{|\sigma'| d_{K,\sigma'}}{d} \mathbf{y}^{\sigma'\sigma},$$

with

$$\mathbf{y}^{\sigma'\sigma} = \begin{cases} \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} + \frac{\sqrt{d}}{d_{K,\sigma}} \left( 1 - \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \right) \mathbf{n}_{K,\sigma} & \text{if } \sigma = \sigma' \\ \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} - \frac{\sqrt{d}}{d_{K,\sigma'}|K|} |\sigma| \mathbf{n}_{K,\sigma} \cdot (\mathbf{x}_{\sigma'} - \mathbf{x}_K) \mathbf{n}_{K,\sigma'} & \text{otherwise .} \end{cases}$$

Using (23), we get that

$$\sum_{\sigma' \in \mathcal{E}_K} \frac{|\sigma'| d_{K,\sigma'}}{d} \mathbf{y}^{\sigma'\sigma} = |\sigma| \mathbf{n}_{K,\sigma}.$$

Since there exists  $C_{11} \in \mathbb{R}_+$  such that  $|R_{K,\sigma'} P_{\mathcal{D},\mathcal{B}u}| \leq C_{11} h_K$ , there exists some  $C_{12} \in \mathbb{R}_+$  with

$$\left| F_{K,\sigma}(P_{\mathcal{D},\mathcal{B}u}) + \int_{\sigma} \nabla u(\mathbf{x}) \cdot \mathbf{n}_{K,\sigma} d\gamma(\mathbf{x}) \right| \leq C_{12} |\sigma| h_K.$$

This leads to the conclusion of (65).  $\square$

## 5 Discrete functional analysis

This section is devoted to some results of functional analysis that are useful for the proof of convergence of numerical schemes when the approximate solution is piece-wise constant on the mesh. Although some of the results presented here were already introduced in previous works of the authors, they were mostly presented (even when not needed, see [17, Remark 9.13 p. 793]) in the framework of “admissible” meshes, that is meshes with an orthogonality condition.

We recall that in the proof of the main convergence Theorem 4.1, we first obtain from the scheme some estimates on the approximate solutions in the discrete  $H^1$  norm. We now show how, from a general discrete  $W^{1,p}$  estimate (this generalisation to  $p \neq 2$  is useful in the case of nonlinear problems) we obtain a discrete  $L^q$  estimate for some  $q > p$  (Lemma 5.3). We then obtain a certain compactness result in  $L^1$  (Lemma 5.5 and therefore in  $L^p$  (Lemma 5.6), which in turn allows to show that the limit of the approximate solution is in  $W_0^{1,p}(\Omega)$  (Lemma 5.7).

### 5.1 Discrete Sobolev embeddings

#### 5.1.1 Discrete embedding of $W^{1,1}$ in $L^{1^*}$

The discrete Sobolev embedding of  $W^{1,1}$  in  $L^{1^*}$  requires less assumptions on the mesh than those given in Definition 2.1. We therefore introduce a larger class of meshes in the following definition.

**Definition 5.1 (Polyhedral partition of  $\Omega$ )** *Let  $d \geq 1$  and let  $\Omega$  be an open bounded set in  $\mathbb{R}^d$ , whose boundary is a finite union of part of hyperplanes. A polyhedral partition  $\mathcal{M}$  of  $\Omega$  is a finite partition of  $\Omega$  such that each element  $K$  of this partition is measurable and has a boundary  $\partial K$  that is composed of a finite union of parts of hyperplanes (the facets of  $K$ ) denoted by  $\sigma$ :  $\partial K = \cup_{\sigma \in \mathcal{E}_K} \sigma$ . Let  $\mathcal{E}$  be the set of the facets of all the elements of  $\mathcal{M}$ :  $\mathcal{E} = \cup_{K \in \mathcal{M}} \mathcal{E}_K$ . If  $\sigma \in \mathcal{E}$  is a facet of this partition, one denotes by  $|\sigma|$  the  $(d-1)$ -Lebesgue measure of  $\sigma$ . Let  $H_{\mathcal{M}}(\Omega)$  be the set of functions from  $\Omega$  to  $\mathbb{R}$ , constant on each element of  $\mathcal{M}$ . Let  $u \in H_{\mathcal{M}}(\Omega)$ . If  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_L$  (that is  $\sigma$  is a facet such that  $\sigma \subset \overline{K} \cap \overline{L}$ ), one sets  $D_{\sigma}u = |u_K - u_L|$ . If  $\sigma \in \mathcal{E}$  is on the boundary of  $\Omega$  and  $K \in \mathcal{M}$  (that is  $\sigma = \partial\Omega \cap \overline{K}$ ), one sets  $D_{\sigma}u = |u_K|$ . For  $u \in H_{\mathcal{M}}(\Omega)$ , one sets*

$$\|u\|_{1,1,\mathcal{M}} = \sum_{\sigma \in \mathcal{E}} |\sigma| D_{\sigma}u. \quad (68)$$

**Lemma 5.1** *Let  $d \geq 1$  and let  $\Omega$  be an open bounded set of  $\mathbb{R}^d$ , whose boundary is a finite union of parts of hyperplanes. Let  $\mathcal{M}$  be a polyhedral partition of  $\Omega$  in the sense of Definition 5.1. Then, with the notations of Definition 5.1,*

$$\|u\|_{L^{1^*}(\Omega)} \leq \frac{1}{2\sqrt{d}} \|u\|_{1,1,\mathcal{M}} \quad \forall u \in H_{\mathcal{M}}(\Omega), \text{ with } 1^* = \frac{d}{d-1}. \quad (69)$$

PROOF. Different proofs of this lemma are possible. A first proof consists in adapting to this discrete setting the classical proof of the Sobolev embedding due to L. Nirenberg (actually, it gives  $1/2$  instead of  $1/(2\sqrt{d})$  in (69)): it is based on an induction on  $d$ . This proof is essentially given in [17, Lemma 9.5 page 790], with slightly less general hypotheses; in fact the so called orthogonality assumption is not used in the proof of Lemma 9.5 of [17]. An easy adaptation of this proof leads to the present lemma (with  $1/2$  instead of  $1/(2\sqrt{d})$  in (69)).

The present proof makes direct use of L. Nirenberg's result, namely:

$$\|u\|_{L^{1^*}(\mathbb{R}^d)} \leq \frac{1}{2d} \|u\|_{W^{1,1}(\mathbb{R}^d)} \quad \forall u \in W^{1,1}(\mathbb{R}^d), \quad (70)$$

where  $\|u\|_{W^{1,1}(\mathbb{R}^d)} = \sum_{i=1}^d \|D_i u\|_{L^1(\mathbb{R}^d)}$  and  $D_i u$  is the weak derivative (or derivative in the sense of distributions) of  $u$  in the direction  $x_i$  (with  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ ).

For  $u \in L^1(\mathbb{R}^d)$ , one sets  $\|u\|_{BV} = \sum_{i=1}^d \|D_i u\|_M$  with, for  $i = 1, \dots, d$ ,  $\|D_i u\|_M = \sup\{\int u \frac{\partial \varphi}{\partial x_i} d\mathbf{x}, \varphi \in C_c^\infty(\mathbb{R}^d), \|\varphi\|_{L^\infty(\mathbb{R}^d)} \leq 1\}$ . The function  $u$  belongs to the space  $BV$  if  $u \in L^1(\mathbb{R}^d)$  and  $\|u\|_{BV} < \infty$ . We first remark that (70) is true with  $\|u\|_{BV}$  instead of  $\|u\|_{W^{1,1}(\mathbb{R}^d)}$ , and if  $u \in BV$  instead of  $W^{1,1}(\mathbb{R}^d)$ . Indeed, to prove this result (which is classical), let  $\rho \in C_c^\infty(\mathbb{R}^d, \mathbb{R}_+)$  with  $\int \rho d\mathbf{x} = 1$ . For  $n \in \mathbb{N}^*$ , define  $\rho_n = n^d \rho(n \cdot)$ . Let  $u \in BV$  and  $u_n = u \star \rho_n$  so that, with (70):

$$\|u_n\|_{L^{1^*}(\mathbb{R}^d)} \leq \frac{1}{2d} \sum_{i=1}^d \|D_i u_n\|_{L^1(\mathbb{R}^d)}. \quad (71)$$

Since  $u_n$  is regular,  $\|D_i u_n\|_{L^1(\mathbb{R}^d)} = \|D_i u_n\|_M$ , and, for  $\varphi \in C_c^\infty(\mathbb{R}^d)$ , using Fubini's theorem:

$$\int_{\mathbb{R}^d} u_n \frac{\partial \varphi}{\partial x_i} d\mathbf{x} = \int_{\mathbb{R}^d} u \frac{\partial}{\partial x_i} (\varphi \star \rho_n) d\mathbf{x} \leq \|D_i u\|_M \|\varphi\|_{L^\infty(\mathbb{R}^d)}.$$

This leads to  $\|D_i u_n\|_{L^1(\mathbb{R}^d)} \leq \|D_i u\|_M$ . Since  $u_n \rightarrow u$  a.e., as  $n \rightarrow \infty$ , at least for a sub-sequence, Fatou's lemma gives, from (71):

$$\|u\|_{L^{1^*}(\mathbb{R}^d)} \leq \frac{1}{2d} \|u\|_{BV} \quad \forall u \in BV. \quad (72)$$

Let  $u \in H_{\mathcal{M}}(\Omega)$ . One sets  $u = 0$  outside  $\Omega$  so that  $u \in L^1(\mathbb{R}^d)$ . One has  $\|u\|_{BV} = \sup\{\int_{\mathbb{R}^d} u \operatorname{div} \varphi d\mathbf{x}, \varphi \in C_c^\infty(\mathbb{R}^d, \mathbb{R}^d), \|\varphi\|_{L^\infty(\mathbb{R}^d)} \leq 1\}$ , with  $\|\varphi\|_{L^\infty(\mathbb{R}^d)} = \sup_{i=1, \dots, d} \|\varphi_i\|_{L^\infty(\mathbb{R}^d)}$  and  $\varphi = (\varphi_1, \dots, \varphi_d)$ . But, for  $\varphi \in C_c^\infty(\mathbb{R}^d, \mathbb{R}^d)$  such that  $\|\varphi\|_{L^\infty(\mathbb{R}^d)} \leq 1$ , an integration by parts on each element of  $\mathcal{M}$  gives (where  $\mathbf{n}_\sigma$  is a normal vector to  $\sigma$  and  $\gamma$  is the  $(d-1)$ -Lebesgue measure on  $\sigma$ ):

$$\int_{\mathbb{R}^d} u \operatorname{div} \varphi d\mathbf{x} = \sum_{\sigma \in \mathcal{E}} D_\sigma u \int_\sigma |\varphi \cdot \mathbf{n}_\sigma| d\gamma(\mathbf{x}) \leq \sqrt{d} \|u\|_{1,1,\mathcal{M}}.$$

Then, one has  $\|u\|_{BV} \leq \sqrt{d} \|u\|_{1,1,\mathcal{M}}$  and (72) leads to (69).

□

### 5.1.2 Discrete embedding of $W^{1,p}$ in $L^{p^*}$ , $1 < p < d$

We now prove a discrete Sobolev embedding for  $1 < p < d$  and for meshes in the sense of Definition 2.1.

**Lemma 5.2** *Let  $d > 1$ ,  $1 < p < d$  and let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$ . Let  $\mathcal{D}$  be a discretization on  $\Omega$  in the sense of Definition 2.1. Let  $\eta > 0$  be such that  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}$ , where  $\mathcal{M}_\sigma = \{K, L\}$ . Then, there exists  $C_{13}$ , only depending on  $d, p$  and  $\eta$  such that*

$$\|u\|_{L^{p^*}(\Omega)} \leq C_{13} \|u\|_{1,p,\mathcal{M}} \quad \forall u \in H_{\mathcal{D}}(\Omega), \quad (73)$$

where  $p^* = \frac{pd}{d-p}$  and

$$\|u\|_{1,p,\mathcal{M}}^p = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{D_\sigma u}{d_\sigma} \right)^p, \quad (74)$$

with  $d_\sigma = d_{K,\sigma} + d_{L,\sigma}$ , if  $\mathcal{M}_\sigma = \{K, L\}$ , and  $d_\sigma = d_{K,\sigma}$ , if  $\mathcal{M}_\sigma = \{K\}$ .

PROOF. We follow here L. Nirenberg's proof of the Sobolev embedding. Let  $\alpha$  be such that  $\alpha 1^* = p^*$  (that is  $\alpha = p(d-1)/(d-p) > 1$ ). Let  $u \in H_{\mathcal{D}}(\Omega)$ . Inequality (69) applied with  $|u|^\alpha$  instead of  $u$  leads to:

$$\left( \int_{\Omega} |u|^{p^*} d\mathbf{x} \right)^{\frac{d-1}{d}} \leq \sum_{\sigma \in \mathcal{E}} |\sigma| D_\sigma |u|^\alpha.$$

For  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\mathcal{M}_\sigma = \{K, L\}$ , one has  $D_\sigma |u|^\alpha \leq \alpha(|u_K|^{\alpha-1} + |u_L|^{\alpha-1}) D_\sigma u$ . For  $\sigma \in \mathcal{E}_{\text{ext}}$ ,  $\mathcal{M}_\sigma = \{K\}$ , one has  $D_\sigma |u|^\alpha \leq \alpha |u_K|^{\alpha-1} D_\sigma u$ . This yields:

$$\left( \int_{\Omega} |u|^{p^*} d\mathbf{x} \right)^{\frac{d-1}{d}} \leq \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| \alpha |u_K|^{\alpha-1} D_\sigma u, \quad (75)$$

For all  $\sigma \in \mathcal{E}$ , one has  $1 \leq \frac{1+\eta}{\eta} \frac{d_{K,\sigma}}{d_\sigma}$ , if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\mathcal{M}_\sigma = \{K, L\}$ , or if  $\sigma \in \mathcal{E}_{\text{ext}}$ ,  $\mathcal{M}_\sigma = \{K\}$ . Then, Hölder's inequality applied to (75) yields, with  $q = p/(p-1)$ :

$$\left( \int_{\Omega} |u|^{p^*} d\mathbf{x} \right)^{\frac{d-1}{d}} \leq \alpha \frac{1+\eta}{\eta} \left( \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} |u_K|^{(\alpha-1)q} \right)^{\frac{1}{q}} \|u\|_{1,p,\mathcal{M}}. \quad (76)$$

Since  $(\alpha-1)q = p^*$ , one has:

$$\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} |u_K|^{(\alpha-1)q} = d \int_{\Omega} |u|^{p^*} d\mathbf{x}.$$

Then, noticing that  $(d-1)/d - 1/q = 1/p^*$ , we deduce (73) follows from (76) with  $C_{13} = \alpha \frac{1+\eta}{\eta} d^{1/q}$  only depending on  $d, p$  and  $\eta$ .  $\square$

### 5.1.3 Discrete embedding of $W^{1,p}$ in $L^q$ , for some $q > p$

Let  $1 \leq p < \infty$ , we now deduce from Lemma 5.3 the following lemma, which gives the discrete embedding of  $W^{1,p}$  in  $L^q$ , for some  $q > p$ .

**Lemma 5.3** *Let  $d \geq 1$ ,  $1 \leq p < \infty$  and let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$ . Let  $\mathcal{D}$  be a mesh of  $\Omega$  in the sense of Definition 2.1. Let  $\eta > 0$  be such that  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}$ , where  $\mathcal{M}_\sigma = \{K, L\}$ . Then, there exists  $q > p$  only depending on  $p$  and there exists  $C_{14}$ , only depending on  $d, \Omega, p$  and  $\eta$  such that*

$$\|u\|_{L^q(\Omega)} \leq C_{14} \|u\|_{1,p,\mathcal{M}} \quad \forall u \in H_{\mathcal{D}}(\Omega), \quad (77)$$

where  $\|u\|_{1,p,\mathcal{M}}^p$  is defined in (74).

PROOF. If  $p = 1$ , one takes  $q = 1^*$  and the result follows from Lemma 5.1 (in this case  $C_{14}$  does not depend on  $\eta$ ). If  $1 < p < d$ , one takes  $q = p^*$  applies Lemma 5.2.

If  $p \geq d$ , one chooses any  $q \in ]p, \infty[$  and  $p_1 < d$  such that  $p_1^* = q$  (this is possible since  $p_1^*$  tends to  $\infty$  as  $p_1$  tends to  $d$ ). Lemma 5.2 gives, for some  $C_{13}$  only depending on  $p, d$  and  $\eta$ , that  $\|u\|_{L^q(\Omega)} \leq C_{13} \|u\|_{1,p_1,\mathcal{M}}$ . But, using Hölder's inequality, there exists  $C_{15}$ , only depending on  $d, p, \Omega$ , such that  $\|u\|_{1,p_1,\mathcal{M}} \leq C_{15} \|u\|_{1,p,\mathcal{M}}$ . Inequality (5.3) follows with  $C_{14} = C_{13} C_{15}$ .  $\square$

## 5.2 Compactness results for bounded families in the discrete $W^{1,p}$ norm

### 5.2.1 Compactness in $L^p$

We prove in this section that bounded families in the discrete  $W^{1,p}$  norms are relatively compact in  $L^p$ . We begin here also with the case  $p = 1$ , giving in this case a crucial inequality which holds for general polyhedral partitions of  $\Omega$ .

**Lemma 5.4** *Let  $d \geq 1$  and let  $\Omega$  be an open bounded set in  $\mathbb{R}^d$ , whose boundary is a finite union of parts of hyperplanes. Let  $\mathcal{M}$  be a polyhedral partition of  $\Omega$  in the sense of Definition 5.1. Then, with the notations of Definition 5.1,*

$$\|u(\cdot + \mathbf{y}) - u\|_{L^1(\mathbb{R}^d)} \leq |\mathbf{y}| \sqrt{d} \|u\|_{1,1,\mathcal{M}} \quad \forall u \in H_{\mathcal{M}}(\Omega), \forall \mathbf{y} \in \mathbb{R}^d, \quad (78)$$

where  $u$  is defined on the whole space  $\mathbb{R}^d$ , taking  $u = 0$  outside  $\Omega$ , and  $|h|$  is the Euclidean norm of  $h \in \mathbb{R}^d$ .

PROOF. One may prove this result in a similar way to that of [17, Lemma 9.3 p. 770] where an  $L^2$  estimate on the translations is proven. Indeed, the proof of Lemma 9.3 [17] holds in the case  $p = 1$  considered here for a general partition, while for  $p > 1$ , it requires the orthogonality condition satisfied by the admissible meshes of [17, Definition 9.1 p 762]. We give here a simpler proof dedicated to the case  $p = 1$ , using the  $BV$ -space, as in Lemma 5.1.

Let  $u \in C_c^\infty(\mathbb{R}^d)$ . For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , one has:

$$|u(\mathbf{x} + \mathbf{y}) - u(\mathbf{x})| = \left| \int_0^1 \nabla u(\mathbf{x} + t\mathbf{y}) \cdot \mathbf{y} dt \right| \leq |\mathbf{y}| \int_0^1 |\nabla u(\mathbf{x} + t\mathbf{y})| dt.$$

Integrating with respect to  $\mathbf{x}$  and using Fubini's Theorem gives the well-known result

$$\|u(\cdot + \mathbf{y}) - u\|_{L^1(\mathbb{R}^d)} \leq |\mathbf{y}| \int_{\mathbb{R}^d} |\nabla u| d\mathbf{x} \leq |\mathbf{y}| \sum_{i=1}^d \|D_i u\|_{L^1(\mathbb{R}^d)}, \quad (79)$$

where  $\nabla u = (D_1 u, \dots, D_d u)$ . By density of  $C_c^\infty(\mathbb{R}^d)$  in  $W^{1,1}(\mathbb{R}^d)$ , Inequality (79) is also true for  $u \in W^{1,1}(\mathbb{R}^d)$ .

We proceed now as in Lemma 5.1, using the same notations. Let  $u \in BV$  and  $u_n = u \star \rho_n$ . Since  $u_n \in W^{1,1}(\mathbb{R}^d)$ , Inequality (79) gives, for all  $\mathbf{y} \in \mathbb{R}^d$ ,  $\|u_n(\cdot + \mathbf{y}) - u_n\|_{L^1(\mathbb{R}^d)} \leq |\mathbf{y}| \sum_{i=1}^d \|D_i u_n\|_{L^1(\mathbb{R}^d)}$ . But, for  $i = 1, \dots, d$ , as in Lemma 5.1,  $\|D_i u_n\|_{L^1(\mathbb{R}^d)} \leq \|D_i u\|_M$ . Then, since  $u_n \rightarrow u$  in  $L^1(\mathbb{R}^d)$ , as  $n \rightarrow \infty$ , we obtain:

$$\|u(\cdot + \mathbf{y}) - u\|_{L^1(\mathbb{R}^d)} \leq |\mathbf{y}| \sum_{i=1}^d \|D_i u\|_M = |\mathbf{y}| \|u\|_{BV} \quad \forall u \in BV, \forall \mathbf{y} \in \mathbb{R}^d. \quad (80)$$

Let  $u \in H_{\mathcal{M}}(\Omega)$ . One sets  $u = 0$  outside  $\Omega$  so that  $u \in L^1(\mathbb{R}^d)$ ; thanks to lemma 5.1,  $\|u\|_{BV} \leq \sqrt{d} \|u\|_{1,1,\mathcal{M}}$  and thus:

$$\|u(\cdot + \mathbf{y}) - u\|_{L^1(\mathbb{R}^d)} \leq |\mathbf{y}| \sqrt{d} \|u\|_{1,1,\mathcal{M}} \quad \forall \mathbf{y} \in \mathbb{R}^d.$$

□

An easy consequence of Lemmas 5.1 and 5.4 is a compactness result in  $L^1$  given in the following lemma.

**Lemma 5.5** *Let  $d \geq 1$  and let  $\Omega$  be an open bounded set in  $\mathbb{R}^d$ , such that its boundary  $\partial\Omega$  is a finite union of parts of hyperplanes. Let  $\mathcal{F}$  be a family of polyhedral partitions of  $\Omega$  in the sense of Definition 5.1. For  $\mathcal{M} \in \mathcal{F}$ , let  $u_{\mathcal{M}} \in H_{\mathcal{M}}(\Omega)$  and assume that there exists  $C \in \mathbb{R}$  such that for all  $\mathcal{M} \in \mathcal{F}$ ,  $\|u_{\mathcal{M}}\|_{1,1,\mathcal{M}} \leq C$ . Then, the family  $(u_{\mathcal{M}})_{\mathcal{M} \in \mathcal{F}}$  is relatively compact in  $L^1(\Omega)$  and also in  $L^1(\mathbb{R}^d)$  taking  $u_{\mathcal{M}} = 0$  outside  $\Omega$ .*

PROOF. By Lemma 5.1, the family  $(u_{\mathcal{M}})_{\mathcal{M} \in \mathcal{F}}$  is bounded in  $L^{1^*}(\Omega)$ . Since  $\Omega$  is bounded, the family  $(u_{\mathcal{M}})_{\mathcal{M} \in \mathcal{F}}$  is bounded in  $L^1(\Omega)$  and also in  $L^1(\mathbb{R}^d)$ , taking  $u_{\mathcal{M}} = 0$  outside  $\Omega$ . Thanks to the Kolmogorov compactness theorem, Lemma 5.4 gives that the family  $(u_{\mathcal{M}})_{\mathcal{M} \in \mathcal{F}}$  is relatively compact in  $L^1(\Omega)$  and also in  $L^1(\mathbb{R}^d)$  taking  $u_{\mathcal{M}} = 0$  outside  $\Omega$ . □

Note that in fact, the above result also holds for general (non polyhedral) partitions of  $\Omega$ , for instance in the case of curved boundaries. In the case  $p > 1$ , we need an additional hypothesis on the meshes which we state in the following lemma.

**Lemma 5.6** *Let  $d \geq 1$ ,  $1 \leq p < \infty$  and  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$ . Let  $F$  be a family of meshes of  $\Omega$  in the sense of Definition 2.1. Let  $\eta > 0$  be such that, for all  $\mathcal{D} \in F$ , one has  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}$ , where  $\mathcal{M}_{\sigma} = \{K, L\}$ . For  $\mathcal{D} \in F$ , let  $u_{\mathcal{D}} \in H_{\mathcal{D}}(\Omega)$  and assume that there exists  $C \in \mathbb{R}$  such, for all  $\mathcal{D} \in F$ ,  $\|u_{\mathcal{D}}\|_{1,p,\mathcal{M}} \leq C$ . Then, the family  $(u_{\mathcal{D}})_{\mathcal{D} \in F}$  is relatively compact in  $L^p(\Omega)$  and also in  $L^p(\mathbb{R}^d)$  taking  $u_{\mathcal{D}} = 0$  outside  $\Omega$ .*

PROOF. Thanks to Lemma 5.3 and to the fact that  $\Omega$  is bounded, the family  $(u_{\mathcal{D}})_{\mathcal{D} \in F}$  is bounded in  $L^1(\Omega)$  and also in  $L^1(\mathbb{R}^d)$  taking  $u_{\mathcal{D}} = 0$  outside  $\Omega$ . Thanks once again to the fact that  $\Omega$  is bounded, the family  $(\|u_{\mathcal{D}}\|_{1,1,\mathcal{M}})_{\mathcal{D} \in F}$  is bounded in  $\mathbb{R}$ . Then, as in the previous lemma, the Kolmogorov compactness theorem gives that the family  $(u_{\mathcal{D}})_{\mathcal{D} \in F}$  is relatively compact in  $L^1(\Omega)$  and also in  $L^1(\mathbb{R}^d)$  taking  $u_{\mathcal{D}} = 0$  outside  $\Omega$ .

In order to conclude we use, once again, Lemma 5.3. It gives that the family  $(u_{\mathcal{D}})_{\mathcal{D} \in F}$  is bounded in  $L^q(\Omega)$  for some  $q > p$ . With the relative compactness in  $L^1(\Omega)$ , this leads to the fact that the family  $(u_{\mathcal{D}})_{\mathcal{D} \in F}$  is relatively compact in  $L^p(\Omega)$  (and then also in  $L^p(\mathbb{R}^d)$  taking  $u_{\mathcal{D}} = 0$  outside  $\Omega$ ). □

### 5.2.2 Regularity of the limit

With the hypotheses of Lemma 5.6, assume that  $u_{\mathcal{D}} \rightarrow u$  in  $L^p$  as  $\text{size}(\mathcal{D}) \rightarrow 0$  (Lemma 5.6 gives that this is possible, at least for subsequences of sequences of meshes with vanishing size). We prove below that  $u \in W_0^{1,p}(\Omega)$ .

**Lemma 5.7** *Let  $d \geq 1$ ,  $1 \leq p < \infty$  and let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$ . Let  $(\mathcal{D}_n)_{n \in \mathbb{N}}$  be a family of discretisations of  $\Omega$  in the sense of Definition 2.1. Let  $\eta > 0$  be such that, for any discretisation  $\mathcal{D}_n = (\mathcal{M}_n, \mathcal{E}_n, \mathcal{P}_n)$ , one has  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}$ , where  $\mathcal{M}_{\sigma} = \{K, L\}$ . For  $n \in \mathbb{N}$ , let  $u^{(n)} \in H_{\mathcal{D}_n}(\Omega)$  and assume that there exists  $C \in \mathbb{R}$  such, for all  $n \in \mathbb{N}$ ,  $\|u^{(n)}\|_{1,p,\mathcal{M}_n} \leq C$ . Assume also that  $\text{size}(\mathcal{D}_n) \rightarrow 0$  as  $n \rightarrow \infty$ . Then:*

1. *There exists a sub-sequence of  $(u^{(n)})_{n \in \mathbb{N}}$ , still denoted by  $(u^{(n)})_{n \in \mathbb{N}}$ , and  $u \in L^p(\Omega)$  such that  $u^{(n)} \rightarrow u$  in  $L^p(\Omega)$  as  $n \rightarrow \infty$ .*
2.  *$u \in W_0^{1,p}(\Omega)$  and*

$$\|\nabla u\|_{L^p(\Omega)^d} = \|\ |\nabla u|\ \|_{L^p(\Omega)} \leq \frac{(1+\eta)d^{\frac{p-1}{p}}}{\eta} C \quad (81)$$

(recall that  $|\nabla u|$  is the Euclidean norm of  $\nabla u$ ).

PROOF. The fact that there exists a subsequence of  $(u^{(n)})_{n \in \mathbb{N}}$ , still denoted by  $(u^{(n)})_{n \in \mathbb{N}}$ , and  $u \in L^p(\Omega)$  such that  $u^{(n)} \rightarrow u$  in  $L^p(\Omega)$  as  $n \rightarrow \infty$  is a consequence of the relative compactness of  $(u^{(n)})_{n \in \mathbb{N}}$  in  $L^p$  given in Lemma 5.6. There only remains to prove that  $u \in W_0^{1,p}(\Omega)$ .

Letting  $u^{(n)} = 0$  and  $u = 0$  outside  $\Omega$ , one also has  $u^{(n)} \rightarrow u$  in  $L^p(\mathbb{R}^d)$ . Let us now construct an approximate gradient, denoted by  $\tilde{\nabla}_{\mathcal{D}_n} u^{(n)}$ , bounded in  $L^p(\Omega)$ , equal to 0 outside  $\Omega$  and converging, at least in the distributional sense, to  $\nabla u$ .

**Step 1** Construction of  $\tilde{\nabla}_{\mathcal{D}} u$ , for  $u \in H_{\mathcal{D}}(\Omega)$ , and its properties.

Let  $n \in \mathbb{N}$  and  $\mathcal{D} = \mathcal{D}_n$ . For this step, one sets  $u = u^{(n)}$  (not to be confused with the limit of the sequence  $(u^{(n)})_{n \in \mathbb{N}}$ ). For  $\sigma \in \mathcal{E}$ , one sets  $u_\sigma = 0$  if  $\sigma$  is on the boundary of  $\Omega$ . Otherwise, one has  $\mathcal{M}_\sigma = \{K, L\}$  and we choose a value  $u_\sigma$  between  $u_K$  and  $u_L$  (it is possible to choose, for instance,  $u_\sigma = \frac{1}{2}(u_K + u_L)$  but any other choice between  $u_K$  and  $u_L$  is possible). Then, one defines  $\tilde{\nabla}_{\mathcal{D}} u$  on  $K \in \mathcal{D}$  in the following way:

$$\tilde{\nabla}_{\mathcal{D}} u = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (u_\sigma - u_K).$$

The function  $\tilde{\nabla}_{\mathcal{D}} u$  is constant on each  $K \in \mathcal{M}$  and, on  $K$ , using Hölder's inequality

$$|\tilde{\nabla}_{\mathcal{D}} u|^p \leq \frac{1}{(|K|)^p} \left( \sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} |u_\sigma - u_K| \right)^p \leq \frac{1}{(|K|)^p} \left( \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \right)^{p-1} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{D_\sigma u}{d_{K,\sigma}} \right)^p.$$

Since  $\sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} = d|K|$ , one deduces

$$\|\tilde{\nabla}_{\mathcal{D}} u\|^p \leq \frac{d^{p-1}}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{D_\sigma u}{d_{K,\sigma}} \right)^p.$$

This gives an  $L^p$ -estimate on  $\tilde{\nabla}_{\mathcal{D}} u$  in  $(L^p(\Omega))^d$  (or in  $(L^p(\mathbb{R}^d))^d$ , setting  $\tilde{\nabla}_{\mathcal{D}} u = 0$  outside  $\Omega$ ), in terms of  $\|u\|_{1,p,\mathcal{M}}$ , namely

$$\|\tilde{\nabla}_{\mathcal{D}} u\|_{L^p} \leq \frac{(1+\eta)d^{\frac{p-1}{p}}}{\eta} \|u\|_{1,p,\mathcal{M}}. \quad (82)$$

In order to prove, in the next step, the convergence of this approximate gradient, we now compute the integral of this gradient against a test function. Let  $\varphi \in C_c^\infty(\mathbb{R}^d; \mathbb{R}^d)$ ,  $\varphi_K$  the mean value of  $\varphi$  on  $K \in \mathcal{D}$ , and  $\varphi_\sigma$  the mean value of  $\varphi$  on  $\sigma$ . Then,

$$\int_{\mathbb{R}^d} \tilde{\nabla}_{\mathcal{D}} u \cdot \varphi \, d\mathbf{x} = \sum_{K \in \mathcal{D}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (u_\sigma - u_K) \varphi_K = \sum_{K \in \mathcal{D}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (-u_K) \varphi_\sigma + R(u, \varphi), \quad (83)$$

with

$$R(u, \varphi) = \sum_{K \in \mathcal{D}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (u_\sigma - u_K) (\varphi_K - \varphi_\sigma).$$

Then, there exists  $C_\varphi$  only depending on  $\varphi$ ,  $d$ ,  $p$ ,  $\Omega$  and  $\eta$  such that  $|R(u, \varphi)| \leq C_\varphi \text{size}(\mathcal{D}) \|u\|_{1,p,\mathcal{M}}$ . Equation (83) can also be written as

$$\int_{\mathbb{R}^d} \tilde{\nabla}_{\mathcal{D}} u \cdot \varphi \, d\mathbf{x} = \sum_{K \in \mathcal{D}} \int_K (-u_K) \operatorname{div}(\varphi) \, d\mathbf{x} + R(u, \varphi) = - \int_{\mathbb{R}^d} u \operatorname{div}(\varphi) \, d\mathbf{x} + R(u, \varphi). \quad (84)$$

**Step 2** Convergence of  $\tilde{\nabla}_{\mathcal{D}_n} u^{(n)}$  to  $\nabla u$  and proof of  $u \in W_0^{1,p}(\Omega)$ .

We consider now the sequence  $(u^{(n)})_{n \in \mathbb{N}}$ . Inequality (82) gives

$$\|\tilde{\nabla}_{\mathcal{D}} u^{(n)}\|_{L^p} \leq \frac{(1+\eta)d^{\frac{p-1}{p}}}{\eta} \|u^{(n)}\|_{1,p,\mathcal{M}}.$$

Then, the sequence  $(\tilde{\nabla}_{\mathcal{D}} u^{(n)})_{n \in \mathbb{N}}$  is bounded in  $L^p(\mathbb{R}^d)^d$  and we can assume, up to a subsequence, that  $\tilde{\nabla}_{\mathcal{D}} u^{(n)}$  converges to some  $w$  weakly in  $L^p(\mathbb{R}^d)^d$ , as  $n \rightarrow \infty$  and  $\| |w| \|_{L^p} \leq \frac{(1+\eta)d^{\frac{p-1}{p}}}{\eta} C$ .

Let  $\varphi \in C_c^\infty(\mathbb{R}^d; \mathbb{R}^d)$ , Equation (84) gives

$$\int_{\mathbb{R}^d} \tilde{\nabla}_{\mathcal{D}} u^{(n)} \cdot \varphi \, d\mathbf{x} = - \int_{\mathbb{R}^d} u^{(n)} \operatorname{div}(\varphi) \, d\mathbf{x} + R(u^{(n)}, \varphi). \quad (85)$$

Thanks to  $|R(u^{(n)}, \varphi)| \leq C_\varphi \operatorname{size}(\mathcal{D}_n) \|u^{(n)}\|_{1,p,\mathcal{M}_n}$ , one has  $R(u^{(n)}, \varphi) \rightarrow 0$ , as  $n \rightarrow \infty$ . Since  $u^{(n)} \rightarrow u$  in  $L^p(\mathbb{R}^d)$  as  $n \rightarrow \infty$ , passing to the limit in (85) gives:

$$\int_{\mathbb{R}^d} w \cdot \varphi \, d\mathbf{x} = - \int_{\mathbb{R}^d} u \operatorname{div}(\varphi) \, d\mathbf{x}.$$

Since  $\varphi$  is arbitrary in  $C_c^\infty(\mathbb{R}^d; \mathbb{R}^d)$ , one deduces that  $\nabla u = w$ . Then  $u \in W^{1,p}(\mathbb{R}^d)$  and  $\| |\nabla u| \|_{L^p} \leq \frac{(1+\eta)d^{\frac{p-1}{p}}}{\eta} C$ . Finally, since  $u = 0$  outside  $\Omega$ , one has  $u \in W_0^{1,p}(\Omega)$ .  $\square$

## 6 Conclusion and perspectives

A symmetric discretisation scheme was introduced for anisotropic heterogeneous problems on distorted nonconforming meshes. Although this scheme stems from the finite volume analysis, which was developed these past years, its formulation is actually derived from a discrete weak formulation; in this respect it may be seen as a nonconforming finite element method. Tools of functional analysis were obtained, which allow a mathematical analysis of the scheme; the convergence of the discrete solution to the exact solution of the continuous problem is shown with no regularity assumption on the solution (other than the natural assumption that it is in  $H_0^1(\Omega)$ ). Even though this convergence result yields no rate of convergence, it is probably more interesting than error estimates which require some assumptions on the diffusion tensor. Nevertheless, we show an order 1 estimate in the case of the Laplace operator, which is readily adaptable to regular (say piece-wise  $C^1$ ) isotropic diffusion operators. The numerical results presented here show the good performance of the scheme (in particular order 2 is obtained for the convergence in the  $L^2$  norm of the solution), and so do three dimensional experiments which were performed in [12] for the incompressible Navier–Stokes equations on general grids. Note that the convergence analysis which is performed here readily extends to the non-linear setting of Leray–Lions operators. This will be the subject of a future paper.

## References

- [1] I. Aavatsmark, T. Barkve, O. Boe, and T. Mannseth. Discretization on non-orthogonal, quadrilateral grids for inhomogeneous, anisotropic media. *J. Comput. Phys.*, 127(1):2–14, 1996.
- [2] I. Aavatsmark, T. Barkve, O. Boe, and T. Mannseth. Discretization on unstructured grids for inhomogeneous, anisotropic media. part i: Derivation of the methods. *SIAM Journal on Sc. Comp.*, 19:1700–1716, 1998.
- [3] I. Aavatsmark, T. Barkve, O. Boe, and T. Mannseth. Discretization on unstructured grids for inhomogeneous, anisotropic media. part ii: Discussion and numerical results. *SIAM Journal on Sc. Comp.*, 19:1717–1736, 1998.
- [4] L. Agelas, D.A. Di Pietro, and R. Masson. A symmetric and coercive finite volume scheme for multiphase porous media flow problems with applications in the oil industry. In R. Eymard and J.-M. Hérard, editors, *Finite Volumes for Complex Applications V*, pages 35–51. Wiley, 2008.

- [5] K. Aziz and A. Settari. *Petroleum reservoir simulation*. Applied Science, London, 1979.
- [6] E. Bertolazzi and G. Manzini. On vertex reconstructions for cell-centered finite volume approximations of 2D anisotropic diffusion problems. *Math. Models Methods Appl. Sci.*, 17(1):1–32, 2007.
- [7] F. Boyer and Hubert F. Finite volume method for 2d linear and nonlinear elliptic problems with discontinuities. *SIAM J. on Numer. Anal.*, 46(6):3032–3070, 2008.
- [8] F. Brezzi, M. Fortin *Mixed and Hybrid Finite Element Methods* Springer-Verlag, New York, 1991.
- [9] 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.
- [10] 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.
- [11] E. Chénier, R. Eymard, and R. Herbin. A collocated finite volume scheme for the incompressible Navier-Stokes equations on general non-matching grids. In R. Eymard and J.-M. Hérard, editors, *Finite Volumes for Complex Applications V*, pages 289–296. Wiley, 2008.
- [12] E. Chénier, R. Eymard, and Herbin R. A collocated finite volume scheme to solve free convection for general non-onorming grids. *J.Comput. Phys.*, under revision.
- [13] Y. Coudière, J.-P. Vila, and Ph. Villedieu. Convergence rate of a finite volume scheme for a two-dimensional convection-diffusion problem. *M2AN Math. Model. Numer. Anal.*, 33(3):493–516, 1999.
- [14] K. Domelevo and P. Omnes. A finite volume method for the laplace equation on almost arbitrary two-dimensional grids. *M2AN Math. Model. Numer. Anal.*, 39(6):1203–1249, 2005.
- [15] J. Droniou and R. Eymard. A mixed finite volume scheme for anisotropic diffusion problems on any grid. *Numer. Math.*, 105(1):35–71, 2006.
- [16] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Comparison between mimetic finite difference methods, hybrid finite volume methods and mixed finite volume methods. in preparation.
- [17] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. In P. G. Ciarlet and J.-L. Lions, editors, *Techniques of Scientific Computing, Part III*, Handbook of Numerical Analysis, VII, pages 713–1020. North-Holland, Amsterdam, 2000.
- [18] R. Eymard, T. Gallouët, and R. Herbin. A cell-centered finite-volume approximation for anisotropic diffusion operators on unstructured meshes in any space dimension. *IMA J. Numer. Anal.*, 26(2):326–353, 2006.
- [19] R. Eymard, T. Gallouët, and R. Herbin. A new finite volume scheme for anisotropic diffusion problems on general grids: convergence analysis. *C. R., Math., Acad. Sci. Paris*, 344(6):403–406, 2007.
- [20] R. Eymard, T. Gallouët, and R. Herbin. Benchmark on anisotropic problems, SUSHI: a scheme using stabilization and hybrid interfaces for anisotropic heterogeneous diffusion problems. In R. Eymard and J.-M. Hérard, editors, *Finite Volumes for Complex Applications V*, pages 801–814. Wiley, 2008.
- [21] R. Eymard and R. Herbin. A new collocated finite volume scheme for the incompressible Navier-Stokes equations on general non matching grids. *C. R. Math. Acad. Sci. Paris*, 344(10):659–662, 2007.

- [22] Ph. Guillaume and V. Latocha. Numerical convergence of a parametrisation method for the solution of a highly anisotropic two-dimensional elliptic problem. *J. Sci. Comput.*, 25(3):423–444, 2005.
- [23] R. Herbin. An error estimate for a finite volume scheme for a diffusion-convection problem on a triangular mesh. *Numer. Methods Partial Differential Equations*, 11(2):165–173, 1995.
- [24] R. Herbin. Finite volume methods for diffusion convection equations on general meshes. In F. Benkhaldoun and R. Vilsmeier, editors, *Finite volumes for complex applications, Problems and Perspectives*, pages 153–160. Hermès, 1996.
- [25] R. Herbin and F. Hubert. Benchmark on discretization schemes for anisotropic diffusion problems on general grids for anisotropic heterogeneous diffusion problems. In R. Eymard and J.-M. Hérard, editors, *Finite Volumes for Complex Applications V*, pages 659–692. Wiley, 2008.
- [26] F. Hermeline. Approximation of diffusion operators with discontinuous tensor coefficients on distorted meshes. *Comput. Methods Appl. Mech. Engrg.*, 192(16-18):1939–1959, 2003.
- [27] Y. Kuznetsov and S. Repin. Convergence analysis and error estimates for mixed finite element methods on distorted meshes. *Numer. Math.*, 13(1):33–51, 2005.
- [28] C. Le Potier. Schéma volumes finis monotone pour des opérateurs de diffusion fortement anisotropes sur des maillages de triangles non structurés. *C. R. Math. Acad. Sci. Paris*, 341(12):787–792, 2005.
- [29] S.V. Patankar. *Numerical heat transfer and fluid flow*. Series in Computational Methods in Mechanics and Thermal Sciences. Washington - New York - London: Hemisphere Publishing Corporation; New York etc.: McGraw-Hill Book Company. XIII, 197 p., 1980.
- [30] J. E. Roberts and J.-M. Thomas. Mixed and hybrid methods. In *Handbook of numerical analysis, Vol. II*, Handb. Numer. Anal., II, pages 523–639. North-Holland, Amsterdam, 1991.
- [31] M. Vohralík. Equivalence between mixed finite element and multi-point finite volume methods. *C. R. Acad. Sci. Paris., Ser. I*, 339:525–528, 2004.