

# A discrete kinetic energy preserving convection operator for variable density flows on locally refined staggered meshes

J.-C. LATCHÉ<sup>\*</sup>, B. PIAR<sup>\*</sup>, K. SALEH<sup>‡</sup>

<sup>\*</sup> Institut de Radioprotection et de Sûreté Nucléaire (IRSN),  
13115 St-Paul-Lez-Durance, France,

<sup>‡</sup> Université de Lyon, CNRS UMR 5208, Université Lyon 1, Institut Camille Jordan,  
43 bd 11 novembre 1918; F-69622 Villeurbanne cedex, France.

[jean-claude.latche, bruno.piar]@irsn.fr, saleh@math.univ-lyon1.fr

December 3, 2019

## Abstract

In this paper, we build and analyze a scheme for the time-dependent variable density Navier-Stokes equations, able to cope with unstructured non-conforming meshes, with hanging nodes. The time advancement relies on a pressure correction algorithm and the space approximation is based on a low-order staggered non-conforming finite element, the so-called Rannacher-Turek element. The convection term in the momentum balance equation is discretized by a finite volume technique, and a careful construction of the fluxes, especially through non-conforming faces, ensures that the solution obeys a discrete kinetic energy balance. Its consistency is addressed by analyzing a model problem, namely the convection-diffusion equation, for which we theoretically establish a first order convergence in space for energy norms. This convergence order is also observed in the numerical experiments for the Navier-Stokes equations. Navier-Stokes equations, pressure correction scheme, finite volumes, finite elements, stability, kinetic energy, non-conforming local refinement.

## 1 Introduction

Let  $\Omega$  be an open bounded connected subset of  $\mathbb{R}^d$ , with  $d \in \{2, 3\}$ , which is supposed to be polygonal if  $d = 2$  and polyhedral if  $d = 3$ . Let  $T \in \mathbb{R}^+$ . We address in this paper the following system of variable density unstationary Navier-Stokes equations:

$$\partial_t(\rho \mathbf{u}) + \mathbf{div}(\rho \mathbf{u} \otimes \mathbf{u}) - \mathbf{div}(\boldsymbol{\tau}(\mathbf{u})) + \nabla p = 0, \quad \text{on } \Omega \times (0, T), \quad (1.1a)$$

$$\partial_t \rho + \mathbf{div}(\rho \mathbf{u}) = 0, \quad \text{on } \Omega \times (0, T). \quad (1.1b)$$

The unknowns  $\mathbf{u} \in \mathbb{R}^d$  and  $p \in \mathbb{R}$  are the velocity and the pressure in the flow. The density  $\rho$  is assumed to be a *known* positive function of  $\Omega \times (0, T)$  (so Equation (1.1b) must be seen as a constraint on  $\mathbf{u}$ ). The shear stress tensor  $\boldsymbol{\tau}$  is given by:

$$\boldsymbol{\tau}(\mathbf{u}) = \mu(\nabla \mathbf{u} + \nabla^t \mathbf{u}) - \frac{2\mu}{3} \mathbf{div} \mathbf{u} \mathbf{I}, \quad (1.2)$$

where  $\mu$  is a positive parameter, possibly depending on  $\mathbf{x}$ . Consequently, we have:

$$\boldsymbol{\tau}(\mathbf{u}) : \nabla \mathbf{u} \geq 0, \quad \forall \mathbf{u} \in \mathbb{R}^d. \quad (1.3)$$

The two equations of system (1.1) respectively express the momentum balance and the mass conservation of the fluid. System (1.1) is supplemented with initial and boundary conditions:

$$\mathbf{u}|_{\partial\Omega} = \mathbf{u}_{\partial\Omega}, \quad \mathbf{u}|_{t=0} = \mathbf{u}_0. \quad (1.4)$$

The present work is a continuation of a research program undertaken to develop staggered schemes satisfying a discrete kinetic energy balance [1, 10]. This point is crucial with respect to many issues: such a relation readily provides stability estimates, may be seen as a prerequisite for LES applications [4], and, last but not least, is a starting point for the extension of the schemes to compressible flows (compressible Navier-Stokes and Euler equations [13]). The difficulty lies in the definition of the velocity convection operator, which must be in some sense consistent with the discrete mass balance. More specifically, let us define the convection operator  $\mathcal{C}$  as  $\mathcal{C} : v \mapsto \partial_t(\rho v) + \mathbf{div}(\rho v \mathbf{u})$ , in such a way that the  $i^{th}$  component of the convection term in the momentum balance equation reads:

$$\partial_t(\rho u_i) + \mathbf{div}(\rho u_i \mathbf{u}) = \mathcal{C}(u_i).$$

Then we observe that, thanks to the mass balance equation,  $\mathcal{C}$  may be recast as a transport operator:

$$\mathcal{C}(v) = \rho(\partial_t v + \mathbf{u} \cdot \nabla v),$$

and this equivalence between a divergence and transport form easily allows to show that (formally):

$$\mathcal{C}(v) \varphi'(v) = \mathcal{C}(\varphi(v)) = \partial_t(\rho \varphi(v)) + \mathbf{div}(\rho \varphi(v) \mathbf{u}),$$

for any regular real function  $\varphi$ . Taking  $\varphi(v) = v^2/2$  yields the kinetic energy identity. We show in [15] that a discrete analogue of this computation holds for finite volume operators, with a similar interplay between the momentum and mass balances, although with some numerical dissipation and only for convex functions  $\varphi$ . Hence, even if, for unstructured quadrangular ( $d = 2$ ) or hexahedric ( $d = 3$ ) meshes, we base our space approximation on the Rannacher-Turek element, we are led to implement a finite volume discretization of the convection term in order to obtain a scheme satisfying a discrete kinetic energy balance. The resulting scheme thus mixes finite element and finite volume techniques, in the spirit of the algorithm developed and analyzed for incompressible flows in [24, 25]. It is based on a pressure-correction method (see *e.g.* [14]); its description, for variable density flows and conforming meshes, may be found in [1].

The objective of the present paper is to extend this scheme to non-conforming meshes. To this purpose, we have to define a discretization for the diffusion and for the convection terms, able to deal with faces with hanging nodes. For the diffusion term, we adopt the same finite element approach as for a regular mesh, and the extension to non-conforming meshes is easily obtained by slightly relaxing the weak continuity constraint across a face [3] (mean continuity is required only across the whole face containing the hanging node and not across each sub-face having this node as vertex). The construction of the convection operator is more intricate, and is the essential difficulty tackled in this work. Basically, the required consistency of the discrete mass and momentum balance equations is sufficient by itself to suggest a definition; this latter is, in some way, abstract, in the sense that it does not rely on a standard finite volume technique based on control volumes of known shape, but on a set of necessary algebraic equation, in a spirit reminiscent of mimetic schemes. The consistency of this operator then remains to be proved; to this purpose, we perform the error analysis of the scheme on a model problem, namely the convection-diffusion equation.

As mentioned above, in their version working on regular meshes (*i.e.* without hanging nodes), these discrete operators have already been used as building bricks in algorithms tackling various problems:

variable density low Mach number flows [1] (as here), barotropic and non-barotropic compressible flows [9, 15], drift-flux two-phase flow model [11, 12, 16]. The present extension thus opens the road to the implementation of local refinement in these schemes (and, as a matter of fact, this potentiality is already offered by our development platform, the open-source software library for fluid applications CALIF<sup>3</sup>S [5], as a striking outcome of the use of Object Oriented Programming).

This paper is organized as follows. We first describe the space discretization (Section 2) then the proposed scheme (Section 3). Section 4 is devoted to the error analysis of the scheme on a model problem, and Section 5 gathers some numerical experiments which exemplifies this theoretical study. In Sections 2 and 3, we present the scheme as it is utilized in industrial applications at IRSN, through CALIF<sup>3</sup>S. In particular, we get rid of restrictions on the mesh which are necessary for the theoretical analysis, but often overcome in practice, and the introduction of which is consequently postponed to Section 4.

## 2 Meshes and discretization spaces

**Definition 2.1 (Unrefined mesh)** A mesh  $\mathcal{M}_0$  is said an unrefined mesh if it is a regular decomposition (in the usual sense of the finite element literature, see *e.g.* [7]) of the domain  $\Omega$  either in quadrilaterals ( $d = 2$ ) or hexahedra ( $d = 3$ ). Each cell  $K$  of  $\mathcal{M}_0$  is defined by the image by the standard  $Q_1$  mapping associated with its vertices, denoted by  $\mathcal{Q}_K$ , of the unit square or cube  $(0, 1)^d$ .

Let us now define a refinement process. In two dimensions, it consists in cutting a cell  $K$  in four sub-cells, which are defined as the image by  $\mathcal{Q}_K$  of the four sub-squares of the unit square  $(\alpha_1/2, (\alpha_1 + 1)/2) \times (\alpha_2/2, (\alpha_2 + 1)/2)$ , for  $(\alpha_1, \alpha_2) \in \{0, 1\}^2$ . In three dimensions, sub-cells are obtained by applying  $\mathcal{Q}_K$  to the eight subset of the unit cube  $(\alpha_1/2, (\alpha_1 + 1)/2) \times (\alpha_2/2, (\alpha_2 + 1)/2) \times (\alpha_3/2, (\alpha_3 + 1)/2)$ , for  $(\alpha_1, \alpha_2, \alpha_3) \in \{0, 1\}^3$ . The additional vertices produced by this process lie in the mid-point of a coarse edge in 2D, and at the center of a coarse face in 3D (precisely speaking, at the image by  $\mathcal{Q}_K$  of the mass center of the associated face of the reference unit cube). The following lemma is essential for the well-posedness of the refinement process, and is a (not so easy) consequence of the properties of the  $Q_1$  mapping.

**LEMMA 2.1** The sub-cells produced by the refinement process are themselves the image of the unit square or cube by the  $Q_1$  mapping associated with their vertices.

*Proof.* We give the proof in two dimensions, the extension to  $d = 3$  being easy although cumbersome. Let  $K$  be a quadrangle of vertices  $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$  and  $\mathbf{a}_4$ . The cell  $K$  is thus defined by:

$$K = \left\{ \sum_{i=1}^4 \varphi_i(\hat{\mathbf{x}}) \mathbf{a}_i, \hat{\mathbf{x}} \in (0, 1)^2 \right\},$$

with  $\varphi_1(\hat{\mathbf{x}}) = (1 - \hat{x}_1)(1 - \hat{x}_2)$ ,  $\varphi_2(\hat{\mathbf{x}}) = \hat{x}_1(1 - \hat{x}_2)$ ,  $\varphi_3(\hat{\mathbf{x}}) = \hat{x}_1\hat{x}_2$  and  $\varphi_4(\hat{\mathbf{x}}) = (1 - \hat{x}_1)\hat{x}_2$ . A fundamental property of the  $Q_1$  mapping is that it may be obtained by 2 (in fact  $d$ ) successive interpolations, performed in any order. Indeed, let us denote by  $\mathcal{I}(\mathbf{a}, \mathbf{b}, s)$  the following interpolated point between  $\mathbf{a}$  and  $\mathbf{b}$ :

$$\mathcal{I}(\mathbf{a}, \mathbf{b}, s) = (1 - s)\mathbf{a} + s\mathbf{b}.$$

Then, by the definition of  $\varphi$ , we get:

$$K = \left\{ \mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_4, \hat{x}_2), \mathcal{I}(\mathbf{a}_2, \mathbf{a}_3, \hat{x}_2), \hat{x}_1), \hat{\mathbf{x}} \in (0, 1)^2 \right\}.$$

In addition, we may check that a first property of the interpolation operator  $\mathcal{I}$  is that, for any vector  $\hat{\mathbf{x}}$  of  $\mathbb{R}^2$ :

$$\mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_4, \hat{x}_2), \mathcal{I}(\mathbf{a}_2, \mathbf{a}_3, \hat{x}_2), \hat{x}_1) = \mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_2, \hat{x}_1), \mathcal{I}(\mathbf{a}_4, \mathbf{a}_3, \hat{x}_1), \hat{x}_2).$$

Let  $\mathbf{a}_{1,4}$  be defined by  $\mathbf{a}_{1,4} = (\mathbf{a}_1 + \mathbf{a}_4)/2$ . Then, it is easy to see that  $\mathcal{I}$  also satisfies:

$$\{\mathcal{I}(\mathbf{a}_1, \mathbf{a}_4, s), s \in (0, 1/2)\} = \{\mathcal{I}(\mathbf{a}_1, \mathbf{a}_{1,4}, s), s \in (0, 1)\}.$$

We a similar definition for  $\mathbf{a}_{1,3}$ , we get for the first sub-cell of  $K$ , let us say  $K_1$ :

$$\begin{aligned} K_1 &= \{\mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_4, \hat{x}_2), \mathcal{I}(\mathbf{a}_2, \mathbf{a}_3, \hat{x}_2), \hat{x}_1), \hat{x}_1 \in (0, \frac{1}{2})^2\} \\ &= \{\mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_{1,4}, \xi), \mathcal{I}(\mathbf{a}_2, \mathbf{a}_{2,3}, \xi), \hat{x}_1), \hat{x}_1 \in (0, \frac{1}{2}), \xi \in (0, 1)\}. \end{aligned}$$

We now permute the interpolations to obtain:

$$K_1 = \{\mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_2, \hat{x}_1), \mathcal{I}(\mathbf{a}_{1,4}, \mathbf{a}_{2,3}, \hat{x}_1), \xi), \hat{x}_1 \in (0, \frac{1}{2}), \xi \in (0, 1)\}.$$

Finally, with  $\mathbf{a}_{1,2} = (\mathbf{a}_1 + \mathbf{a}_2)/2$  and  $\mathbf{a}_{1,2,3,4} = (\mathbf{a}_1 + \mathbf{a}_2 + \mathbf{a}_3 + \mathbf{a}_4)/4$ , we get:

$$K_1 = \{\mathcal{I}(\mathcal{I}(\mathbf{a}_1, \mathbf{a}_{1,2}, \eta), \mathcal{I}(\mathbf{a}_{1,4}, \mathbf{a}_{1,2,3,4}, \eta), \xi), (\eta, \xi) \in (0, 1)^2\}.$$

This means that  $K_1$  may be obtained from the unit square by the  $Q_1$  mapping associated with the vertices  $\mathbf{a}_1, \mathbf{a}_{1,2}, \mathbf{a}_{1,2,3,4}$  and  $\mathbf{a}_{1,4}$ , which is the result we are searching for. The same conclusion may be obtained for the other sub-cells of  $K$ , by simple change of axes for the unit square. In three dimensions, the proof follows similar lines, but the interpolation operator  $\mathcal{I}$  must be applied three times (along each of the coordinates) instead of twice.  $\square$

We are now in position to define a locally refined mesh.

**Definition 2.2 (Refined mesh)** A locally refined mesh is obtained from an unrefined one by recursively splitting some cells by the above defined refinement process, in such a way that the number of hanging nodes per face is at most one (which means that the difference of level of refinement between two adjacent cells is at most one).

Examples of locally refined meshes are given on Figure 1

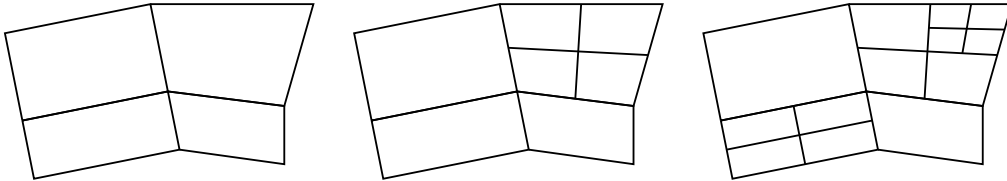


Figure 1: An exemple of admissible mesh refinement.

We denote by  $\mathcal{E}(K)$  the set of the faces of an element  $K \in \mathcal{M}$ . We exclude the presence of a node in the interior of a face, *i.e.* we split an initial face in  $2^{d-1}$  faces if one of the cells adjacent to the face is split. The number of faces,  $N_K^{\mathcal{E}}$ , of a cell  $K$  thus ranges between  $2d$  and  $2^d d$ . Let  $\mathcal{E} = \cup_{K \in \mathcal{M}} \mathcal{E}(K)$ ,  $\mathcal{E}_{\text{ext}} = \{\sigma \in \mathcal{E}, \sigma \subset \partial\Omega\}$  and  $\mathcal{E}_{\text{int}} = \mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ . A face  $\sigma \in \mathcal{E}_{\text{int}}$  separating the cells  $K$  and  $L$  is denoted by  $K|L$ . Hereafter,  $|\cdot|$  stands for the  $d$ - or  $(d-1)$ -dimensional measure of a subset of  $\mathbb{R}^d$  or  $\mathbb{R}^{d-1}$  respectively.

For  $\sigma \in \mathcal{E}(K)$ ,  $\mathbf{n}_{K,\sigma}$  stands for a unit normal vector to  $\sigma$  outward  $K$ . In two dimensions and for plane faces in three dimensions, its definition is clear. For non-plane faces, many definitions are possible; here, we split the face in four 2-dimensional simplices the boundary of which joins the center of the face (with

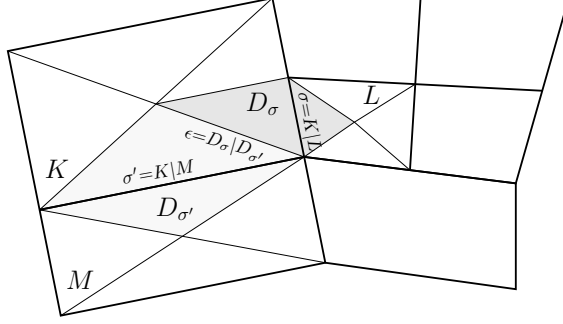


Figure 2: Notations for control volumes and diamond cells.

the same definition as before, namely the image by  $\mathcal{Q}_K$  of the mass center of the associated face of the unit cube) and two vertices of the face, and use for  $\mathbf{n}_{K,\sigma}$  the average of the unit normal vectors of these four simplices, weighted by their area.

We define a dual mesh associated with the faces  $\mathcal{E}$  as follows. When  $K \in \mathcal{M}$  is a rectangle or a cuboid, for  $\sigma \in \mathcal{E}(K)$ , we define the half-diamond cell  $D_{K,\sigma}$  as the cone with basis  $\sigma$  and with vertex the mass center of  $K$  (see Fig. 2). We thus obtain a partition of  $K$  in  $N_K^\mathcal{E}$  sub-volumes, each sub-volume having a measure  $|D_{K,\sigma}|$  equal to  $|K|/(2d)$ , when  $\sigma$  has not been split, or  $|K|/(2^d d)$  otherwise. We extend this definition to general quadrangles and hexahedra, by supposing that we have built a partition with the same connectivities and the same ratio between the volumes of the half-diamonds and of the cell. For  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , we now define the dual (or diamond) cell  $D_\sigma$  associated with  $\sigma$  by  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$ . For  $\sigma \in \mathcal{E}(K) \cap \mathcal{E}_{\text{ext}}$ , we define  $D_\sigma = D_{K,\sigma}$ . We denote by  $\tilde{\mathcal{E}}(D_\sigma)$  the set of faces of  $D_\sigma$ , and by  $\epsilon = D_\sigma | D_{\sigma'}$  the face separating two dual cells  $D_\sigma$  and  $D_{\sigma'}$  (see Fig. 2).

The space discretization is staggered in the sense that the pressure and the velocity unknowns are discretized as piecewise constant functions respectively on the primal and dual mesh. In addition, in order for the algorithm to be further suitable for more general models including the density as unknown (*e.g.* the asymptotic model for anisothermal low Mach number flows, see Section 5.2), we suppose that  $\rho$  is approximated by a discrete function which is also piecewise constant on the primal mesh. Hence, the degrees of freedom for the pressure and the density are associated with the cells of the primal mesh:  $\{p_K, K \in \mathcal{M}\}$  and  $\{\rho_K, K \in \mathcal{M}\}$  while a discrete velocity field is associated with degrees of freedom localized at the cells of the dual mesh:  $\{\mathbf{u}_\sigma, \sigma \in \mathcal{E}\}$ .

### 3 The pressure correction scheme

#### 3.1 General form of the scheme

Let us consider a uniform partition  $0 = t_0 < t_1 < \dots < t_N = T$  of the time interval  $(0, T)$ , and let  $\delta t = t_{n+1} - t_n$  for  $n = 0, \dots, N-1$  be the constant time step. The initial discrete pressure and velocity are defined as follows:

$$\begin{aligned} p_K^0 &= 0, & K \in \mathcal{M}, \\ \mathbf{u}_\sigma^0 &= \frac{1}{|\sigma|} \int_\sigma \mathbf{u}_0(\mathbf{x}) \, d\gamma(\mathbf{x}), & \sigma \in \mathcal{E}_{\text{int}}. \end{aligned} \quad (3.1)$$

The Dirichlet boundary condition is taken into account by setting  $\mathbf{u}_\sigma^n$  to the mean value of  $\mathbf{u}_{\partial\Omega}$  over  $\sigma$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$  and all  $n$  in  $\{0, 1, \dots, N\}$ .

We also define discrete values of the density at the discrete time  $n = -1$ , which are needed by the algorithm, in order for this latter to satisfy a discrete kinetic energy balance for the scheme (see Section 3.4 below) :

$$\rho_K^{-1} = \frac{1}{|K|} \int_K \rho(\mathbf{x}, 0) \, d\mathbf{x} \quad K \in \mathcal{M}. \quad (3.2)$$

As usual [26, 6, 14], the pressure correction scheme is a two-step algorithm. Let us assume that  $(\rho_K^{n-1})_{K \in \mathcal{M}} \subset \mathbb{R}$ ,  $(\rho_K^n)_{K \in \mathcal{M}} \subset \mathbb{R}$ ,  $(p_K^n)_{K \in \mathcal{M}} \subset \mathbb{R}$  and  $(\mathbf{u}_\sigma^n)_{\sigma \in \mathcal{E}_{\text{int}}} \subset \mathbb{R}^d$  are known families of real numbers. Then a time step computation consists in finding  $(p_K^{n+1})_{K \in \mathcal{M}} \subset \mathbb{R}$  and  $(\mathbf{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}} \subset \mathbb{R}^d$  through the following procedure:

*Prediction step* – Find  $(\mathbf{u}_\sigma^{n+\frac{1}{2}})_{\sigma \in \mathcal{E}_{\text{int}}}$  such that:

$$\begin{aligned} \frac{1}{\delta t} (\rho_\sigma^n \mathbf{u}_\sigma^{n+\frac{1}{2}} - \rho_\sigma^{n-1} \mathbf{u}_\sigma^n) + \frac{1}{|D_\sigma|} \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon}^n \mathbf{u}_\epsilon^{n+\frac{1}{2}} \\ - (\operatorname{div} \boldsymbol{\tau}(\mathbf{u}))_\sigma^{n+\frac{1}{2}} + \left( \frac{\rho_\sigma^n}{\rho_\sigma^{n-1}} \right)^{\frac{1}{2}} (\nabla p)_\sigma^n = 0, \quad \sigma \in \mathcal{E}_{\text{int}}. \end{aligned} \quad (3.3a)$$

*Correction step* – Find  $(\mathbf{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  and  $(p_K^{n+1})_{K \in \mathcal{M}}$  such that:

$$\frac{1}{\delta t} \rho_\sigma^n (\mathbf{u}_\sigma^{n+1} - \mathbf{u}_\sigma^{n+\frac{1}{2}}) + (\nabla p)_\sigma^{n+1} - \left( \frac{\rho_\sigma^n}{\rho_\sigma^{n-1}} \right)^{\frac{1}{2}} (\nabla p)_\sigma^n = 0, \quad \sigma \in \mathcal{E}_{\text{int}}, \quad (3.3b)$$

$$\frac{1}{\delta t} (\rho_K^{n+1} - \rho_K^n) + \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}^{n+1} = 0, \quad K \in \mathcal{M}. \quad (3.3c)$$

In the following sections, we give the definitions of the notations used in (3.3). In particular, Section 3.4 is devoted to the definition of the quantities  $\rho_\sigma$  and  $F_{\sigma,\epsilon}$  which are built so as to ensure a discrete kinetic energy balance for the scheme.

## 3.2 Mass balance equation and pressure gradient

Equation (3.3c) is a discretization of the mass balance over the primal mesh, and  $F_{K,\sigma}^{n+1}$  stands for the mass flux across  $\sigma$  outward  $K$ . On the primal mesh faces,  $F_{K,\sigma}^{n+1}$  is given by:

$$F_{K,\sigma}^{n+1} = |\sigma| \widehat{\rho}_\sigma^{n+1} \mathbf{u}_\sigma^{n+1} \cdot \mathbf{n}_{K,\sigma}, \quad \sigma \in \mathcal{E},$$

where  $\widehat{\rho}_\sigma^{n+1}$  stands for an approximation of the density at the face  $\sigma$ . Note that, since the density is not an unknown and we do not need that the discrete mass balance equation ensures its positivity, any reasonable approximation may be used. Here, we choose a centered approximation at the internal faces of the domain (more precisely speaking, we compute  $\widehat{\rho}_\sigma^{n+1}$  as the mean value of its approximation at the two neighbour cells of  $\sigma$ ). For inflow external faces (*i.e.* a face where the flow is entering the domain), the quantity  $\widehat{\rho}_\sigma^{n+1}$  is directly expressed from the data on  $\rho$ ; for an outflow one, we take  $\widehat{\rho}_\sigma^{n+1} = \rho_K^{n+1}$ , where  $K$  is the cell adjacent to  $\sigma$ .

The pressure gradient is built as the dual operator of the discrete divergence, itself expressed from the already defined approximation of  $\operatorname{div}(\rho \mathbf{u})$  by considering that the density is everywhere equal to 1. For  $k = n, n+1$ ,

$$(\nabla p)_\sigma^k = \frac{|\sigma|}{|D_\sigma|} (p_L^k - p_K^k) \mathbf{n}_{K,\sigma}, \quad \sigma = K|L.$$

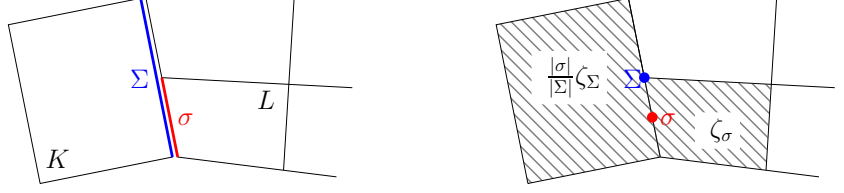


Figure 3: Piecewise definition of  $\zeta_\sigma$ .

### 3.3 Discretization of the diffusion term

The discretization of the diffusion term relies on the so-called "rotated bi-linear element" introduced by Rannacher and Turek [23]. The reference element  $\widehat{K}$  is the unit  $d$ -cube  $(0, 1)^d$ , and the discrete functional space is:

$$\widetilde{Q}_1(\widehat{K}) = \text{span} \{1, (\mathbf{x}_i)_{i=1, \dots, d}, (\mathbf{x}_i^2 - \mathbf{x}_{i+1}^2)_{i=1, \dots, d-1}\}. \quad (3.4)$$

When no vertex of the face  $\sigma$  is a hanging node, we impose the jump of a discrete function through the face to have a zero mean value. When one vertex is a hanging node, we only impose to zero the integral of the jump through the initial coarse face containing  $\sigma$ . Hence, the set  $\{\zeta_\sigma, \sigma \in \mathcal{E}\}$  of nodal functions associated with the Rannacher-Turek element is defined as follows.

Let  $K \in \mathcal{M}$  and  $\sigma$  be a face of  $K$ . If  $\sigma$  is a whole side of  $K$ , the standard definition applies:

$$\begin{aligned} (i) \quad & \zeta_\sigma|_K = \zeta_\sigma \circ \mathcal{Q}_K^{-1} \text{ where } \zeta_\sigma \text{ is some function of } \widetilde{Q}_1(\widehat{K}), \\ (ii) \quad & \frac{1}{|\sigma|} \int_\sigma \zeta_\sigma = 1 \text{ and, for all other sides } \Sigma' \text{ of } K, \int_{\Sigma'} \zeta_\sigma = 0. \end{aligned} \quad (3.5)$$

Let us now suppose that  $\sigma$  is only a subset of a side of  $K$ , which we denote by  $\Sigma$  (which occurs when  $\sigma$  separates  $K$  from a cell which refinement level is equal to the refinement level of  $K$  plus 1, see Figure 3). Let  $\zeta_\Sigma$  be the Rannacher-Turek usual shape function associated with  $\Sigma$  (*i.e.* the function satisfying the above definition (3.5), replacing  $\sigma$  by  $\Sigma$ ). Then, we define  $\zeta_\sigma$  on  $K$  by:

$$\zeta_\sigma(\mathbf{x}) = \frac{|\sigma|}{|\Sigma|} \zeta_\Sigma(\mathbf{x}).$$

Finally, of course, the nodal functions are local, in the sense that, for  $\sigma \in \mathcal{E}$ , the support of  $\zeta_\sigma$  is reduced to the (one or two) cells adjacent to  $\sigma$ .

Finally, the discretization of the  $i^{\text{th}}$  component of the diffusion term reads:

$$-(\text{div} \boldsymbol{\tau}(\mathbf{u}))_{\sigma, i}^{n+\frac{1}{2}} = \frac{1}{|D_\sigma|} \sum_{K \in \mathcal{M}} \int_K \sum_{\sigma' \in \mathcal{E}(K)} \sum_{j=1}^d \mathbf{u}_{\sigma', j}^{n+\frac{1}{2}} \boldsymbol{\tau}(\zeta_{\sigma'} \mathbf{e}^{(j)}) : \nabla(\zeta_\sigma \mathbf{e}^{(i)}) \, d\mathbf{x},$$

where, for  $1 \leq j \leq d$ ,  $\mathbf{e}^{(j)}$  stands for the  $j^{\text{th}}$  vector of the canonical basis of  $\mathbb{R}^d$ .

### 3.4 The discrete momentum convection operator

In this section, we describe the approximation of the convection operator  $\partial_t(\rho \mathbf{u}) + \text{div}(\rho \mathbf{u} \otimes \mathbf{u})$  which appears in the momentum balance equation. It is of finite volume type, and takes the general form given by the first two terms of (3.3a).

The quantities  $\rho_\sigma^n$  and  $\rho_\sigma^{n-1}$  are approximations of the density on the dual cell  $D_\sigma$  at time  $t_n$  and  $t_{n-1}$  respectively, while the quantities  $F_{\sigma,\epsilon}^n$  are mass fluxes across the edges of the dual cells. These quantities are built so that a finite volume discretization of the mass balance (1.1b) holds over the internal dual cells:

$$\frac{1}{\delta t} (\rho_\sigma^n - \rho_\sigma^{n-1}) + \frac{1}{|D_\sigma|} \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon}^n = 0, \quad \sigma \in \mathcal{E}_{\text{int}}. \quad (3.6)$$

This is crucial in order to reproduce, at the discrete level, the derivation of a kinetic energy balance equation, a consequence of which are discrete analogues of the usual  $L^\infty(L^2)$ - and  $L^2(H^1)$ - stability estimates for the velocity, if the density is assumed to be bounded by below by a positive constant. Indeed, the following result holds:

**PROPOSITION 3.1** (Discrete kinetic energy balance) Assume that the discrete densities  $\rho_\sigma^n, \rho_\sigma^{n-1}$  and the dual fluxes  $F_{\sigma,\epsilon}^n$  are built so as to satisfy equation (3.6) for all  $\sigma \in \mathcal{E}_{\text{int}}$  and all  $n$  in  $\{0, 1, \dots, N-1\}$ . Then, any solution to the scheme (3.3) satisfies the following equality, for all  $\sigma \in \mathcal{E}_{\text{int}}$  and  $n$  in  $\{0, 1, \dots, N\}$ :

$$\begin{aligned} \frac{1}{2\delta t} \left( \rho_\sigma^n |\mathbf{u}_\sigma^{n+1}|^2 - \rho_\sigma^{n-1} |\mathbf{u}_\sigma^n|^2 \right) + \frac{1}{2|D_\sigma|} \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma) \\ \epsilon = D_\sigma | D_{\sigma'}'}} F_{\sigma,\epsilon}^n \mathbf{u}_\sigma^{n+\frac{1}{2}} \cdot \mathbf{u}_{\sigma'}^{n+\frac{1}{2}} \\ - (\operatorname{div} \boldsymbol{\tau}(\mathbf{u}))_{\sigma}^{n+\frac{1}{2}} \cdot \mathbf{u}_\sigma^{n+\frac{1}{2}} + (\nabla p)_\sigma^{n+1} \cdot \mathbf{u}_\sigma^{n+1} = -R_\sigma^{n+1} - P_\sigma^{n+1}, \end{aligned} \quad (3.7)$$

where

$$R_\sigma^{n+1} = \frac{1}{2\delta t} \rho_\sigma^{n-1} |\mathbf{u}_\sigma^{n+\frac{1}{2}} - \mathbf{u}_\sigma^n|^2, \quad P_\sigma^{n+1} = \delta t \left( \frac{1}{\rho_\sigma^n} |(\nabla p)_\sigma^{n+1}|^2 - \frac{1}{\rho_\sigma^{n-1}} |(\nabla p)_\sigma^n|^2 \right).$$

Multiplying equation (3.7) by  $\delta t$  and summing over  $\sigma \in \mathcal{E}_{\text{int}}$  and  $n = 0, \dots, M-1$ , yields, for all  $M$  in  $\{0, 1, \dots, N\}$ :

$$\begin{aligned} \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma^{M-1} |\mathbf{u}_\sigma^M|^2 + \sum_{n=0}^{M-1} \delta t \sum_{K \in \mathcal{M}} \int_K \boldsymbol{\tau}(\tilde{\mathbf{u}}^{n+\frac{1}{2}}) : \nabla \tilde{\mathbf{u}}^{n+\frac{1}{2}} + \mathcal{R}^M + \mathcal{P}^M \\ = \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma^{-1} |\mathbf{u}_\sigma^0|^2, \end{aligned} \quad (3.8)$$

where for all  $n \in \mathbb{N}$ ,  $\tilde{\mathbf{u}}^{n+\frac{1}{2}} = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \mathbf{u}_\sigma^{n+\frac{1}{2}} \zeta_\sigma$ , and

$$\mathcal{R}^M = \sum_{n=0}^{M-1} \sum_{\sigma \in \mathcal{E}_{\text{int}}} \rho_\sigma^{n-1} |\mathbf{u}_\sigma^{n+\frac{1}{2}} - \mathbf{u}_\sigma^n|^2 \geq 0, \quad \mathcal{P}^M = \delta t^2 \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{1}{\rho_\sigma^M} |(\nabla p)_\sigma^M|^2 \geq 0.$$

*Proof.* We refer to [16, 13] where the proof of this result can be found in the context of a similar pressure correction scheme for Euler and Navier-Stokes equations.  $\square$

Let us now give the detailed construction of the dual densities and mass fluxes that ensures (3.6). For  $\sigma \in \mathcal{E}_{\text{int}}$  such that  $\sigma = K|L$ , the approximate densities on the dual cell  $D_\sigma$  are given by the following weighted average:

$$|D_\sigma| \rho_\sigma^k = \xi_K^\sigma |K| \rho_K^k + \xi_L^\sigma |L| \rho_L^k, \quad \text{for } k = n-1 \text{ and } k = n, \quad (3.9)$$

where

$$\xi_K^\sigma = \frac{|D_{K,\sigma}|}{|K|}, \quad K \in \mathcal{M}, \quad \sigma \in \mathcal{E}(K). \quad (3.10)$$



The set of dual fluxes  $F_{\sigma,\epsilon}^n$  with  $\epsilon$  included in the primal cell  $K$ , is computed by solving a linear system which right-hand side is a linear combination of the primal fluxes  $(F_{K,\sigma}^n)_{\sigma \in \mathcal{E}(K)}$ , appearing in the discrete mass balance (1.1b) at the previous time step. More precisely, we have the following definition for the dual fluxes, in which we omit for short the time dependence (*i.e.* the superscript  $n$ ).

DEFINITION 1 (Definition of the dual fluxes from the primal ones) The fluxes through the faces of the dual mesh are defined so as to satisfy the following three constraints:

- (H1) – For all primal cell  $K$  in  $\mathcal{M}$ , the set  $(F_{\sigma,\epsilon})_{\epsilon \subset K}$  of dual fluxes through faces included in  $K$  satisfies the following linear system:

$$F_{K,\sigma} + \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \epsilon \subset K} F_{\sigma,\epsilon} = \xi_K^\sigma \sum_{\sigma' \in \mathcal{E}(K)} F_{K,\sigma'}, \quad \forall \sigma \in \mathcal{E}(K). \quad (3.11)$$

- (H2) – The dual fluxes are conservative:  $F_{\sigma,\epsilon} = -F_{\sigma',\epsilon}$  for all  $\epsilon = D_\sigma | D_{\sigma'}$ .
- (H3) – The dual fluxes are bounded with respect to the primal ones  $(F_{K,\sigma})_{\sigma \in \mathcal{E}(K)}$ :

$$|F_{\sigma,\epsilon}| \leq C \max \{|F_{K,\sigma}|, \sigma \in \mathcal{E}(K)\}, \quad K \in \mathcal{M}, \sigma \in \mathcal{E}(K), \epsilon \in \tilde{\mathcal{E}}(D_\sigma), \epsilon \subset K.$$

Appendix A provides the detailed construction of the dual fluxes (*i.e.* the way to find a solution to (3.11) satisfying the other two constraints). Owing to these definitions of the dual densities and mass fluxes, the mass balance on the dual cells (3.6) is an easy consequence of the mass balance on the cells of the primal mesh [1]:

$$\frac{1}{\delta t} (\rho_K^n - \rho_K^{n-1}) + \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}^n = 0, \quad K \in \mathcal{M}.$$

Since the mass balance is not yet solved when performing the prediction step, the mass fluxes  $F_{K,\sigma}^{n+1}$  are not yet known. This is what leads us to perform a backward shift of the density in the prediction step, thus using  $\rho^n$  and  $\rho^{n-1}$  rather than the apparently more natural choice  $\rho^{n+1}$  and  $\rho^n$ .

To make the description of the scheme complete, it remains to define the velocity interpolates  $\mathbf{u}_\epsilon^{n+\frac{1}{2}}$  at the internal dual faces. We chose a centered approximation:

$$\mathbf{u}_\epsilon^{n+\frac{1}{2}} = \frac{\mathbf{u}_\sigma^{n+\frac{1}{2}} + \mathbf{u}_{\sigma'}^{n+\frac{1}{2}}}{2}, \quad \text{for } \epsilon = D_\sigma | D_{\sigma'}.$$

## 4 Study of the prediction step: Error analysis for the advection-diffusion equation

For the sake of simplicity, we propose to perform an error analysis of the scheme on the simplified model of the advection-diffusion equation. This scalar equation can be seen as a model problem for the prediction step of the method (for one component of the velocity), where the known density is assumed to be constant and the pressure gradient and the advective velocity are known given functions from the previous time step. Since the density is considered constant, for the kinetic energy identity to hold, we need to suppose that the given advection field  $\mathbf{w}$  is divergence free; we assume in addition that  $\mathbf{w}$  is a regular function, namely  $\mathbf{w} \in C^1(\bar{\Omega})^d$ , and vanishes at the boundary  $\partial\Omega$  of the computational domain. The problem we consider thus consists in finding a function  $u$  such that:

$$u + \text{div}(u\mathbf{w}) - \Delta u = f, \quad \text{on } \Omega, \quad (4.1a)$$

$$u = 0, \quad \text{on } \partial\Omega, \quad (4.1b)$$

where  $f \in L^2(\Omega)$ . In order to write a weak formulation of (4.1), let us define the following forms corresponding respectively to the diffusion and convection terms:

$$\begin{aligned} a(u, v) &= \int_{\Omega} \nabla u \cdot \nabla v, & u, v \in H_0^1(\Omega), \\ b(\mathbf{w}, u, v) &= \int_{\Omega} \operatorname{div}(u\mathbf{w})v, & u, v \in H_0^1(\Omega). \end{aligned}$$

We also introduce the inner product of  $L^2(\Omega)$  denoted  $(\cdot, \cdot)$ , and we are now in position to state the weak formulation of problem (4.1).

**DEFINITION 2** A weak solution of the advection-diffusion problem (4.1) is a function  $u \in H_0^1(\Omega)$  such that:

$$(u, v) + b(\mathbf{w}, u, v) + a(u, v) = (f, v), \quad \forall v \in H_0^1(\Omega). \quad (4.2)$$

It is well known that, by the Lax-Milgram theorem, there exists a unique solution to (4.2).

#### 4.1 Regularity of the mesh and approximation space

In addition to the definition of the mesh given in Section 2, we suppose that any quadrilaterals ( $d = 2$ ) or hexahedra ( $d = 3$ ) of the mesh are convex, which implies that their faces are hyperplanes of  $\mathbb{R}^d$ . For  $K \in \mathcal{M}$ , we denote by  $h_K$  the diameter of  $K$ . Similarly, we denote by  $h_\sigma$  the diameter of an edge  $\sigma \in \mathcal{E}$ . The size of the discretization is defined as usual by:

$$h = \max\{h_K, K \in \mathcal{M}\}.$$

For the consistency of the Rannacher-Turek finite element approximation of the diffusion term, we need a measure of the difference between the cells of  $\mathcal{M}$  and parallelograms ( $d = 2$ ) or parallelotopes ( $d = 3$ ), as defined in [23]. For  $K \in \mathcal{M}$ , we denote by  $\bar{\alpha}_K$  the maximum of the angles between the normal vectors of opposite faces, choosing for the latter the orientation which maximizes the angle, and set  $\alpha_K = \pi - \bar{\alpha}_K$  (so  $\alpha_K = 0$  if  $K$  is a parallelogram or a parallelotope, and  $\alpha_K > 0$  otherwise). Then we define  $\alpha_h$  as:

$$\alpha_h = \max\{\alpha_K, K \in \mathcal{M}\}. \quad (4.3)$$

For  $K \in \mathcal{M}$ , let  $\{\mathbf{a}_{i,K}, i = 1, \dots, 2^d\}$  denote the vertices of  $K$ . Let  $\mathbf{a}_{i,K}$  be one of these vertices and let  $S_{i,K}$  be the simplex whose vertices are  $\mathbf{a}_{i,K}$  and the  $d$  adjacent vertices to  $\mathbf{a}_{i,K}$ . We denote by  $r_{i,K}$  the diameter of the largest ball included in  $S_{i,K}$  and by  $r_K$  the real number given by  $r_K = \min\{r_{i,K}, i = 1, \dots, 2^d\}$ . We define the real number  $\theta_h$  by:

$$\theta_h = \max\left\{\frac{h_K}{r_K}, K \in \mathcal{M}\right\}. \quad (4.4)$$

In accordance with the velocity discretization described in Section 3, a discrete function  $v$  is associated with degrees of freedom  $\{v_\sigma, \sigma \in \mathcal{E}_{\text{int}}\}$  located at the internal faces, the values associated with the boundary faces  $\mathcal{E} \in \mathcal{E}_{\text{ext}}$  being set to zero, consistently with the boundary conditions (4.1b). For  $K \in \mathcal{M}$ , the restriction of a discrete function  $v$  to  $K$  belongs to the local Rannacher-Turek space:  $v|_K = \hat{v} \circ \mathcal{Q}_K^{-1}$  where  $\hat{v} \in \tilde{Q}_1(\hat{K})$  (see (3.4)). The last step to obtain a definition of the discrete approximation space is to state the continuity constraints satisfied by the discrete functions at the faces. Once again, when no vertex of the face is a hanging node, this issue is clear: as usual for the Rannacher-Turek element, we impose the jump through the face to have a zero mean value. Otherwise, we only impose to zero the integral of the jump through the initial coarse face. Hence the approximation space is given by

$$V_h = \left\{v(\mathbf{x}) = \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_\sigma \zeta_\sigma(\mathbf{x}), \quad (v_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}} \subset \mathbb{R}\right\}, \quad (4.5)$$

where the shape functions  $\zeta_\sigma$  are defined in Section 3.1.

With the continuity requirements described above, the functions of  $V_h$  are discontinuous across each edge; the discretization is thus non-conforming in the sense that  $V_h \not\subset H_0^1(\Omega)$ . Therefore we define the "broken" gradient  $\nabla_h v$  on  $V_h$  as the function of  $L^2(\Omega)^d$  which is equal to  $\nabla v|_K$  for all  $K \in \mathcal{M}$ . The corresponding broken Sobolev  $H^1$  semi-norm is defined for  $v \in V_h$  by:

$$\|v\|_{h,b}^2 = \int_{\Omega} |\nabla_h v|^2 = \sum_{K \in \mathcal{M}} \int_K |\nabla v|^2.$$

Thanks to the homogeneous Dirichlet boundary conditions, this defines a norm on  $V_h$  which is known to control the  $L^2$ -norm by a Poincaré inequality:

$$\|v\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|v\|_{h,b}, \quad \forall v \in V_h. \quad (4.6)$$

## 4.2 The staggered scheme for the advection-diffusion equation

Let us now adapt to the advection-diffusion problem the first step of the staggered scheme proposed in Section 3. For this purpose, we assume that the divergence of the velocity field  $\mathbf{w}$  is discretized on the cells of the primal mesh by

$$\text{div}(\mathbf{w})|_K \approx \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}(\mathbf{w}), \quad \text{with } F_{K,\sigma}(\mathbf{w}) = \int_{\sigma} \mathbf{w} \cdot \mathbf{n}_{K,\sigma} d\gamma, \quad (4.7)$$

and we define the dual fluxes  $F_{\sigma,\epsilon}(\mathbf{w})$  from the primal ones following exactly the procedure described in Section 3.4. Adapting the staggered scheme to the advection-diffusion problem thus consists in finding a discrete function  $u_h \in V_h$  such that:

$$(u_h, v)_h + b_h(\mathbf{w}, u_h, v) + a_h(u_h, v) = (f, v), \quad \text{for all } v \in V_h, \quad (4.8)$$

where, for any  $u, v \in V_h$ ,

$$(u, v)_h = \sum_{\sigma \in \mathcal{E}} |D_\sigma| u_\sigma v_\sigma, \quad (4.9)$$

$$b_h(\mathbf{w}, u, v) = \sum_{\sigma \in \mathcal{E}} v_\sigma \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon}(\mathbf{w}) u_\epsilon, \quad (4.10)$$

$$a_h(u, v) = \int_{\Omega} \nabla_h u \cdot \nabla_h v. \quad (4.11)$$

## 4.3 Stability properties of the scheme

One first important property of the scheme follows from the particular construction of the discrete convection term  $b_h(\mathbf{w}, u, v)$ . Indeed, since  $\mathbf{w}$  is divergence-free, one has

$$\sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}(\mathbf{w}) = 0,$$

which, by construction, implies a similar discrete divergence-free property on the cells of the dual mesh:

$$\sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} F_{\sigma,\epsilon}(\mathbf{w}) = 0, \quad \forall \sigma \in \mathcal{E}_{\text{int}}. \quad (4.12)$$

This is crucial in order to reproduce, at the discrete level, the derivation of a stability estimate for the scheme, analogous to the kinetic energy balance equation in the context of the full Navier-Stokes equations. Indeed, by a computation which may be found in [1], Equation (4.12) yields:

$$b_h(\mathbf{w}, v, v) = 0, \quad \forall v \in V_h. \quad (4.13)$$

REMARK 4.1 At first glance, this relation is the usual well known antisymmetry of the centered convection operator with a divergence free advection field; note however that the scheme is posed on a dual mesh, which, if we restrict the set of dual cells to those associated with the internal faces, does not cover  $\Omega$ , and that no balance equation is written on the remainder of the domain (*i.e.* the half-diamond-shaped volumes associated with the external faces). This necessitates a slight adaptation of the usual proof, which is performed in [1].

A first consequence of this antisymmetry property is the existence of a unique solution to the scheme (4.8). Let  $\mathcal{A}_h$  be the following form defined on the finite dimensional space  $V_h \times V_h$  by

$$\mathcal{A}_h(u, v) = (u, v)_h + b_h(\mathbf{w}, u, v) + a_h(u, v), \quad u, v \in V_h.$$

The mapping  $\Phi$  from  $V_h$  to  $V'_h$  which maps an element  $u$  to the form  $\Phi(u) := \mathcal{A}_h(u, \cdot)$  is clearly linear and satisfies  $(\Phi(u) = 0 \Rightarrow u = 0)$  since  $\mathcal{A}_h(u, u) = (u, u)_h + \|u\|_{h,b}^2$  for all  $u \in V_h$  by (4.13). Hence,  $\Phi$  is a one-to-one linear mapping which, in a finite dimensional context, is equivalent to the existence of a unique solution  $u_h \in V_h$  which satisfies  $\Phi(u_h) = f$  *i.e.*  $\mathcal{A}_h(u_h, v) = (f, v)$  for all  $v \in V_h$ .

A second straightforward consequence is a stability estimate on the solution  $u_h$ . Taking  $u_h$  as a test function in (4.8) and using again (4.13), one gets  $(u_h, u_h)_h + \|u_h\|_{h,b}^2 = (f, u_h)$ . Applying Young's inequality in the right hand side term and recalling the Poincaré inequality (4.6), we obtain the stability estimate

$$2(u_h, u_h)_h + \|u_h\|_{h,b}^2 \leq \text{diam}(\Omega)^2 \|f\|_{L^2(\Omega)}^2. \quad (4.14)$$

#### 4.4 Error estimate

We may now state the main result of this section, which is an error estimate for the scheme (4.8).

THEOREM 4.1 Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  (in the sense of Definition 2.2) such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). Let  $u_h \in V_h$  be the solution to the scheme (4.8). We assume that the solution  $u$  of the continuous problem (4.2) belongs to  $H_0^1(\Omega) \cap H^2(\Omega)$ . Then  $u_h$  satisfies:

$$\|u_h - u\|_{h,b} \leq C(h + \alpha_h) \|u\|_{H^2(\Omega)},$$

where  $C$  only depends on  $\mathbf{w}$ ,  $\Omega$  and  $\theta_0$ .

REMARK 4.2 (A particular construction of a regular sequence of discretizations) For  $d = 2$ , a sequence of discretizations satisfying  $h \rightarrow 0$  and  $\alpha_h \rightarrow 0$  is obtained by successively dividing each quadrangle in four sub-quadrangles, splitting it along the lines joining the mid-points of opposite faces. Unfortunately, the extension of this construction to the three-dimensional case is not straightforward, since this subdivision process may generate non-plane faces.

REMARK 4.3 Theorem 4.1 shows that the accuracy of the scheme for the energy norm  $\|\cdot\|_{h,b}$  is the same as that of the usual Rannacher-Turek approximation of the Stokes problem on non-refined meshes [23]. In particular, there is no loss in the convergence rate due to the non-conforming local refinement or the particular discretization of the convection term.

#### 4.4.1 Preliminary lemmas

We begin with stating some technical lemmas, which will be useful in the proof of Theorem 4.1. We first introduce the following discrete  $H^1$ -norm on the space  $V_h$ :

$$\|v\|_{h,\text{fv}}^2 = \sum_{K \in \mathcal{M}} h_K^{d-2} \sum_{\sigma, \sigma' \in \mathcal{E}(K)} |v_\sigma - v_{\sigma'}|^2, \quad (4.15)$$

which, by an easy computation, may be shown to be equivalent, over a regular sequence of discretizations such that  $\max_h \theta_h \leq \theta_0$  for some  $\theta_0 > 0$ , to the usual finite volume  $H^1$ -norm. Lemma 4.1 shows that this  $H^1$ -norm is controlled by the broken-Sobolev  $H^1$ -norm.

**LEMMA 4.1** Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). Then, there exists  $C$  only depending on  $\theta_0$  such that:

$$\|v\|_{h,\text{fv}} \leq C \|v\|_{h,\text{b}}, \quad \forall v \in V_h.$$

*Proof.* For  $K \in \mathcal{M}$ , let  $\hat{v}$  be the function defined over  $\hat{K}$  by  $\hat{v}(\hat{\mathbf{x}}) = v(\mathbf{x})$ , where  $\mathbf{x} \in K$  stands for the image of  $\hat{\mathbf{x}}$  by the  $Q_1$  mapping  $\mathcal{Q}_K$ . By definition of the discretization space, we have  $\hat{v} \in \tilde{Q}_1(\hat{K})$ . Now, since there exists only a bounded number of possible configurations for  $K$  (depending on the fact that its sides are split in  $2^{d-1}$  faces or not), a finite dimensional argument for norms acting on  $\tilde{Q}_1(\hat{K})$  shows that there exists a constant  $C$  such that:

$$\sum_{\sigma, \sigma' \in \mathcal{E}(K)} |v_\sigma - v_{\sigma'}|^2 \leq C \int_{\hat{K}} |\nabla \hat{v}|^2.$$

We conclude the proof by invoking standard properties of the  $Q_1$  mapping which enable to write:

$$\int_{\hat{K}} |\nabla \hat{v}|^2 \leq C(\theta_0) \frac{h_K^2}{|K|} \int_K |\nabla v|^2 \leq C'(\theta_0) h_K^{2-d} \int_K |\nabla v|^2.$$

□

The following lemma compares piecewise constant approximates of a discrete function with the function itself. It is an easy consequence of the previous one.

**LEMMA 4.2** Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). For  $v \in V_h$ , let  $v_c$  stand for the piecewise constant function over each diamond cell  $D_\sigma$  and equal to  $v_\sigma$ . Then, there exists  $C$ , only depending on  $\theta_0$ , such that:

$$\|v - v_c\|_{L^2(\Omega)} \leq C h \|v\|_{h,\text{b}}.$$

Let now  $v_m$  be defined as a piecewise constant function over the primal mesh, the value over a cell  $K$  being a convex combination of  $(v_\sigma)_{\sigma \in \mathcal{E}(K)}$ . Then, once again, there exists  $C$ , only depending on  $\theta_0$ , such that:

$$\|v - v_m\|_{L^2(\Omega)} \leq C h \|v\|_{h,\text{b}}.$$

The following lemma is the analogue, for the Rannacher-Turek element and a locally refined mesh, of a well known result for the Crouzeix-Raviart element and regular meshes.

LEMMA 4.3 Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). Let  $v \in V_h$ . For  $\sigma \in \mathcal{E}_{\text{int}}$ , let  $[v]_\sigma$  denote the jump of  $v$  across  $\sigma$  (the orientation of which we do not need to precise), and, for  $\sigma \in \mathcal{E}_{\text{ext}}$ , let  $[v]_\sigma = v$ . Then the following inequality holds:

$$\left( \sum_{\sigma \in \mathcal{E}} \frac{1}{h_\sigma} \int_\sigma [v]_\sigma^2 \right)^{1/2} \leq C \|v\|_{h,b}.$$

where  $C$  only depends on  $\theta_0$ .

*Proof.* Let  $\sigma$  be a face of the mesh. By assumption, the integral of the jump of a discrete function  $v$  vanishes either through  $\sigma$  or through a coarse face  $\Sigma$  including  $\sigma$ . In any case, there exists  $\mathbf{x}_\sigma \in \Sigma$  such that  $[v]_\sigma(\mathbf{x}_\sigma) = 0$  and the distance between  $\mathbf{x}_\sigma$  and any point of  $\sigma$  is lower than  $C(\theta_0)h_\sigma$ . We thus get:

$$\begin{aligned} \sum_{\sigma \in \mathcal{E}} \frac{1}{h_\sigma} \int_\sigma [v]_\sigma^2 &= \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \frac{1}{h_\sigma} \int_\sigma \left( \int_{\mathbf{x}_\sigma}^{\mathbf{x}} (\nabla v|_K - \nabla v|_L) \cdot \frac{\mathbf{x} - \mathbf{x}_\sigma}{|\mathbf{x} - \mathbf{x}_\sigma|} \right)^2 \\ &\quad + \sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}(K)} \frac{1}{h_\sigma} \int_\sigma \left( \int_{\mathbf{x}_\sigma}^{\mathbf{x}} \nabla v|_K \cdot \frac{\mathbf{x} - \mathbf{x}_\sigma}{|\mathbf{x} - \mathbf{x}_\sigma|} \right)^2. \end{aligned}$$

The Cauchy-Schwarz inequality yields:

$$\begin{aligned} \sum_{\sigma \in \mathcal{E}} \frac{1}{h_\sigma} \int_\sigma [v]_\sigma^2 &\leq \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \frac{1}{h_\sigma} \int_\sigma |\mathbf{x} - \mathbf{x}_\sigma| \int_{\mathbf{x}_\sigma}^{\mathbf{x}} |\nabla v|_K|^2 + |\nabla v|_L| \\ &\quad + \sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}(K)} \frac{1}{h_\sigma} \int_\sigma |\mathbf{x} - \mathbf{x}_\sigma| \int_{\mathbf{x}_\sigma}^{\mathbf{x}} |\nabla v|_K|^2. \end{aligned}$$

Bounding  $|\mathbf{x} - \mathbf{x}_\sigma|$  by  $C(\theta_0)h_\sigma$  and switching the order of integration, we get:

$$\sum_{\sigma \in \mathcal{E}} \frac{1}{h_\sigma} \int_\sigma [v]_\sigma^2 \leq C(\theta_0) \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \int_{\mathbf{x}_\sigma}^{\mathbf{x}} \int_\sigma |\nabla v|_K|^2 + |\nabla v|_L| + C(\theta_0) \sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}(K)} \int_{\mathbf{x}_\sigma}^{\mathbf{x}} \int_\sigma |\nabla v|_K|^2,$$

and so, finally:

$$\sum_{\sigma \in \mathcal{E}} \frac{1}{h_\sigma} \int_\sigma [v]_\sigma^2 \leq C(\theta_0)^2 \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} h_\sigma \int_\sigma |\nabla v|_K|^2.$$

We conclude the proof as for Lemma 4.1 by transporting the quantity  $\int_\sigma |\nabla v|_K|^2$  to the reference element  $\widehat{K}$ , invoking finite dimensional arguments to relate the integral of the squared gradient norm over the faces and over the element and, finally, standard properties of the  $Q_1$  mapping.  $\square$

The proof of the following trace lemma is an easy adaptation of a result which can be found in [8, appendix A].

LEMMA 4.4 Let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  and  $K$  be a control volume of  $\mathcal{M}$ , and let  $\sigma$  be one of its faces. Then there exists  $C$ , only depending on  $d$ , such that the following inequality holds:

$$\|v\|_{L^2(\sigma)} \leq C \frac{1}{r_K^{1/2}} \left( \|v\|_{L^2(K)} + h_K \|\nabla v\|_{L^2(K)^d} \right), \quad \forall v \in H^1(K).$$

We will also need the following Poincaré-Wirtinger inequality, which is proven for any convex domain  $K$  in [21, 2].

LEMMA 4.5 For all convex domain  $K$  of  $\mathbb{R}^d$ ,  $1 \leq d \leq 3$ :

$$\|v - m_K(v)\|_{L^2(K)} \leq \frac{1}{\pi} h_K \|\nabla v\|_{L^2(K)^d}, \quad \forall v \in H^1(K), \quad (4.16)$$

where  $m_K(v)$  stands for the mean value of  $v$  over  $K$ .

We now give two corollaries of Lemmas 4.4 and 4.5; they both compare different mean values of functions of  $H^1(\Omega)$ .

LEMMA 4.6 Let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$ ; let  $K \in \mathcal{M}$  and let  $D$  be a subset of  $K$ . For all  $v \in H^1(K)$ , we denote  $m_K(v)$  and  $m_D(v)$  the mean values of  $v$  on  $K$  and  $D$  respectively. Then:

$$|m_K(v) - m_D(v)| \leq \frac{1}{\pi |D|^{\frac{1}{2}}} h_K \|\nabla v\|_{L^2(K)^d}, \quad \forall v \in H^1(K).$$

*Proof.* By the Cauchy-Schwarz inequality, we have:

$$|m_K(v) - m_D(v)| = \frac{1}{|D|} \left| \int_D (v - m_K(v)) \right| \leq \frac{1}{|D|^{\frac{1}{2}}} \|v - m_K(v)\|_{L^2(D)}.$$

Thus, since  $D$  is a subset of  $K$ , we have:

$$|m_K(v) - m_D(v)| \leq \frac{1}{|D|^{\frac{1}{2}}} \|v - m_K(v)\|_{L^2(K)},$$

and we conclude by (4.16), since  $K$  is convex.  $\square$

LEMMA 4.7 Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). Let  $K \in \mathcal{M}$  be a control volume and let  $\sigma$  be a face of  $K$ . For all  $v \in H^1(K)$ , we denote  $m_\sigma(v)$  and  $m_K(v)$  the mean values of  $v$  on  $\sigma$  and  $K$  respectively. Then, there exists  $C$ , only depending on  $d$  and  $\theta_0$  such that:

$$|m_\sigma(v) - m_K(v)| \leq \frac{C}{r_K^{d/2}} h_K \|\nabla v\|_{L^2(K)^d}, \quad \forall v \in H^1(K). \quad (4.17)$$

*Proof.* We have  $|m_\sigma(v) - m_K(v)| \leq |\sigma|^{-1} \int_\sigma |v - m_K(v)| \leq |\sigma|^{-1/2} \|v - m_K(v)\|_{L^2(\sigma)}$  by the Cauchy-Schwarz inequality. Invoking successively Lemma 4.4 and Lemma 4.5 and using the regularity of the mesh yields the result.  $\square$

We are now in position to prove the following technical lemma.

LEMMA 4.8 Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \geq \theta_0$ , where  $\theta_h$  is defined by (4.4). We define  $\mathcal{E}_c$  as the set containing:

- the internal faces of the mesh which do not have as vertex a hanging node,
- the coarse faces, *i.e.* for any hanging node  $\mathbf{b}$ , the subset of a hyperplane made of the union of the faces having  $\mathbf{b}$  as vertex.

Let  $(a_\Sigma)_{\Sigma \in \mathcal{E}_c}$  be a family of real numbers such that for all  $\Sigma \in \mathcal{E}_c$ ,  $|a_\Sigma| \leq 1$ . Let  $v$  be a function of the Rannacher-Turek space  $V_h$  associated with  $\mathcal{M}$ , and, for  $\Sigma \in \mathcal{E}_c$  let  $[v]_\Sigma$  be the jump of  $v$  through  $\Sigma$ . Then the following bound holds:

$$\sum_{\Sigma \in \mathcal{E}_c} \left| \int_\Sigma a_\Sigma [v]_\Sigma g \right| \leq C h \|v\|_{h,\mathbf{b}} \|g\|_{H_0^1(\Omega)}, \quad \forall g \in H_0^1(\Omega).$$

where the real number  $C$  only depends on  $\theta_0$  and  $d$ .

*Proof.* Let  $\Sigma \in \mathcal{E}_c$  and  $v \in V_h$ . Since the integral of the jump of  $v$  through  $\Sigma$  is zero, we have:

$$\int_{\Sigma} a_{\Sigma} [v]_{\Sigma} g = \int_{\Sigma} a_{\Sigma} [v]_{\Sigma} (g - g_{\Sigma}),$$

where  $g_{\Sigma}$  is any real number. By the Cauchy-Schwarz inequality in  $L^2(\Sigma)$ , we thus get:

$$\sum_{\Sigma \in \mathcal{E}_c} \left| \int_{\Sigma} a_{\Sigma} [v]_{\Sigma} g \right| \leq \sum_{\Sigma \in \mathcal{E}_c} \left( \int_{\Sigma} [v]_{\Sigma}^2 \right)^{\frac{1}{2}} \left( \int_{\Sigma} (g - g_{\Sigma})^2 \right)^{\frac{1}{2}}.$$

Let us now decompose the integral over  $\Sigma$  in integrals over the faces (in fact, for the coarse faces only), then use the concavity of the square root function, to obtain:

$$\sum_{\Sigma \in \mathcal{E}_c} \left| \int_{\Sigma} a_{\Sigma} [v]_{\Sigma} g \right| \leq \sum_{\Sigma \in \mathcal{E}_c} \left( \sum_{\sigma \subset \Sigma} \left( \int_{\sigma} [v]_{\sigma}^2 \right)^{\frac{1}{2}} \right) \left( \int_{\Sigma} (g - g_{\Sigma})^2 \right)^{\frac{1}{2}}.$$

The discrete Cauchy-Schwarz inequality now yields:

$$\sum_{\Sigma \in \mathcal{E}_c} \left| \int_{\Sigma} a_{\Sigma} [v]_{\Sigma} g \right| \leq \left( \sum_{\Sigma \in \mathcal{E}_c} \frac{1}{h_{\Sigma}} \int_{\Sigma} [v]_{\Sigma}^2 d\gamma \right)^{\frac{1}{2}} \underbrace{\left( \sum_{\Sigma \in \mathcal{E}_c} h_{\Sigma} \int_{\Sigma} (g - g_{\Sigma})^2 \right)^{\frac{1}{2}}}_{T_1},$$

where  $h_{\Sigma}$  stands for the sum of the diameters of the faces included in  $\Sigma$  (which is equal to the diameter of  $\Sigma$  in two dimensions, and lower than four times this diameter in three dimensions). By Lemma 4.3, the first term of the latter product is bounded by  $C(\theta_0) \|v\|_{h,b}$ . For the second one, for  $\Sigma \in \mathcal{E}_c$ , let  $K_{\Sigma}$  be a cell of the mesh having  $\Sigma$  as a whole side (two choices are possible for a standard face, and only one for a coarse one). Applying the trace lemma 4.4, we get:

$$T_1^2 \leq C(d) \sum_{\Sigma \in \mathcal{E}_c} \frac{h_{\Sigma}}{r_{K_{\Sigma}}} \left( \|g - g_{\Sigma}\|_{L^2(K_{\Sigma})}^2 + h_{K_{\Sigma}}^2 \|\nabla g\|_{L^2(K_{\Sigma})}^2 \right).$$

Choosing for  $g_{\Sigma}$  the mean value of  $g$  on  $K_{\Sigma}$  and using (4.16), we thus get:

$$T_1^2 \leq C(d) \sum_{\Sigma \in \mathcal{E}_c} \left( 1 + \frac{1}{\pi^2} \right) \frac{h_{\Sigma}}{r_{K_{\Sigma}}} h_{K_{\Sigma}}^2 \|\nabla g\|_{L^2(K_{\Sigma})}^2 \leq C(d, \theta_0) h^2 \sum_{\Sigma \in \mathcal{E}_c} \|\nabla g\|_{L^2(K_{\Sigma})}^2.$$

The result follows by observing that the  $H^1$  semi-norm of  $g$  on a given cell  $K$  of the mesh appears at most  $2d$  (the maximum number of sides of a cell  $K$ ) times in the summation.  $\square$

We define by  $r_h$  the following interpolation operator:

$$r_h : \begin{cases} H_0^1(\Omega) & \longrightarrow V_h \\ v & \mapsto r_h v(\mathbf{x}) = \sum_{\sigma \in \mathcal{E}} |\sigma|^{-1} \left( \int_{\sigma} v d\gamma \right) \zeta_{\sigma}(\mathbf{x}). \end{cases} \quad (4.18)$$

The stability and approximation properties of  $r_h$  are given in the following lemma.

**LEMMA 4.9** Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). There exists  $C_1$  and  $C_2$ , only depending on  $\theta_0$  such that



1. Stability:

$$\forall v \in \mathbf{H}_0^1(\Omega), \quad \|r_h v\|_{h,b} \leq C_1 \|\nabla v\|_{\mathbf{L}^2(\Omega)^d}.$$

2. Approximation properties:

$$\begin{aligned} \forall v \in \mathbf{H}_0^1(\Omega) \cap \mathbf{H}^2(\Omega), \quad \forall K \in \mathcal{M}, \\ \|v - r_h v\|_{\mathbf{L}^2(K)} + h_K \|\nabla(v - r_h v)\|_{\mathbf{L}^2(K)^d} \leq C_2 h_K (h_K + \alpha_K) |v|_{\mathbf{H}^2(K)}. \end{aligned}$$

*Proof.* The stability property follows from usual estimates on the shape functions and from the trace lemma 4.4. We prove the approximation property. If there is no hanging node on the faces of  $K$ ,  $r_h$  is the usual Rannacher-Turek interpolation operator, and the result is known. In the other case, for  $d = 2$ , let us suppose that an initial face  $\Sigma$  of  $K$  has been split:  $\Sigma = \sigma_1 \cup \sigma_2$ . Let  $v \in \mathbf{H}_0^1(\Omega)$ . Then we get for the part of the expansion of  $r_h v$  associated with  $\sigma_1$  and  $\sigma_2$ , let us say  $r_h v^{\sigma_1 \cup \sigma_2}$ :

$$\begin{aligned} r_h v^{\sigma_1 \cup \sigma_2}|_K &= \frac{1}{|\sigma_1|} \left( \int_{\sigma_1} v \, d\gamma \right) \frac{|\sigma_1|}{|\Sigma|} \zeta_\Sigma(\mathbf{x}) + \frac{1}{|\sigma_2|} \left( \int_{\sigma_2} v \, d\gamma \right) \frac{|\sigma_2|}{|\Sigma|} \zeta_\Sigma(\mathbf{x}) \\ &= \frac{1}{|\Sigma|} \left( \int_{\Sigma} v \, d\gamma \right) \zeta_\Sigma(\mathbf{x}). \end{aligned}$$

Once again, we recognize the usual Rannacher-Turek interpolation operator. The same arguments readily extend to the 3D case.  $\square$

REMARK 4.4 (*inf-sup* condition on locally refined meshes) The same computation shows that, as for the usual Rannacher-Turek approximation on regular meshes, the operator  $r_h$ , or, more precisely, its natural extension to vector-valued functions, is a Fortin operator (*i.e.* continuous from  $\mathbf{H}_0^1(\Omega)^d$  to  $V_h^d$  endowed with the  $H^1$ -broken norm and such that  $\int_\Omega q \operatorname{div}(\mathbf{u} - r_h \mathbf{u}) = 0$  for any discrete pressure function  $q$ ); this implies that the *inf-sup* condition is satisfied by the pair of velocity and pressure approximation spaces also on locally refined meshes.

#### 4.4.2 Estimates on the discrete convective term

The form  $b_h$  is a discretization of the convection term on the cells of the dual mesh. The analysis of the scheme actually requires an equivalent (or nearly equivalent) re-formulation of the form  $b_h$  on the cells of the primal mesh  $\mathcal{M}$ , that makes use of the primal fluxes  $F_{K,\sigma}(\mathbf{w})$ . Indeed, contrary to the dual fluxes  $F_{\sigma,\epsilon}(\mathbf{w})$ , the expression of  $F_{K,\sigma}(\mathbf{w})$  with respect to the convection field  $\mathbf{w}$  is quite simple (see (4.7)). This motivates the introduction of the following auxiliary form:

$$\tilde{b}_h(\mathbf{w}, u, v) = \sum_{K \in \mathcal{M}} v_K \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}(\mathbf{w}) u_\sigma, \quad u, v \in V_h, \quad (4.19)$$

where  $v_K = \sum_{\sigma \in \mathcal{E}(K)} \xi_K^\sigma v_\sigma$  is a convex combination of  $(v_\sigma)_{\sigma \in \mathcal{E}(K)}$ , where the coefficient  $\xi_K^\sigma$  is equal to  $1/(2d)$ , when  $\sigma$  has not been split, or  $1/(2^d d)$  otherwise. The following lemma provides a bound of the error made when replacing  $b_h$  by  $\tilde{b}_h$ . It may be seen as a simplified version of a slightly more general result, dealing with density-dependent fluxes  $F_{K,\sigma}$ , which may be found in [18]; we however give its proof, for the sake of completeness.

LEMMA 4.10 Let  $\theta_0 > 0$  and let  $\mathcal{M}$  be a locally refined mesh of the computational domain  $\Omega$  such that  $\theta_h \leq \theta_0$ , with  $\theta_h$  defined by (4.4). There exists  $C$ , only depending on  $\theta_0$  such that:

$$|\tilde{b}_h(\mathbf{w}, u, v) - b_h(\mathbf{w}, u, v)| \leq C h \|\mathbf{w}\|_{\mathbf{L}^\infty(\Omega)^d} \|u\|_{h,b} \|v\|_{h,b}, \quad \forall u, v \in V_h. \quad (4.20)$$

*Proof.* Denote  $R = b_h(\mathbf{w}, u, v) - \tilde{b}_h(\mathbf{w}, u, v)$ . In the expression (4.10) of  $b_h(\mathbf{w}, u, v)$ , for  $\sigma = K|L$ , let us split the sum over the fluxes through the faces of  $D_\sigma$  in the sum over the dual faces, on one side, included in  $K$  and, on the other side, included in  $L$ . We get by conservativity (*i.e.* using  $F_{K,\sigma}(\mathbf{w}) = -F_{L,\sigma}(\mathbf{w})$ ):

$$b_h(\mathbf{w}, u, v) = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} v_\sigma \left( F_{K,\sigma}(\mathbf{w}) u_\sigma + \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \\ \epsilon \subset K, \epsilon = D_\sigma | D'_\sigma}} F_{\sigma,\epsilon}(\mathbf{w}) \frac{u_\sigma + u_{\sigma'}}{2} \right).$$

Let us write  $b_h(\mathbf{w}, u, v) = T_1 + T_2$  with:

$$T_1 = \sum_{K \in \mathcal{M}} v_K \sum_{\sigma \in \mathcal{E}(K)} \left( F_{K,\sigma}(\mathbf{w}) u_\sigma + \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \\ \epsilon \subset K, \epsilon = D_\sigma | D'_\sigma}} F_{\sigma,\epsilon}(\mathbf{w}) \frac{u_\sigma + u_{\sigma'}}{2} \right),$$

$$T_2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} (v_\sigma - v_K) \left( F_{K,\sigma}(\mathbf{w}) u_\sigma + \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \\ \epsilon \subset K, \epsilon = D_\sigma | D'_\sigma}} F_{\sigma,\epsilon}(\mathbf{w}) \frac{u_\sigma + u_{\sigma'}}{2} \right).$$

By the conservativity assumption of the dual fluxes (H2) (see Definition 1), we remark that  $T_1 = \tilde{b}_h(\mathbf{w}, u, v)$  so that  $R = T_2$ . Using now (H1), we write  $R = R_1 + R_2$  with:

$$R_1 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} (v_\sigma - v_K) \left( \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \\ \epsilon \subset K, \epsilon = D_\sigma | D'_\sigma}} F_{\sigma,\epsilon}(\mathbf{w}) \frac{u_{\sigma'} - u_\sigma}{2} \right),$$

$$R_2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} (v_\sigma - v_K) u_\sigma \xi_K^\sigma \left( \sum_{\sigma' \in \mathcal{E}(K)} F_{K,\sigma'}(\mathbf{w}) \right).$$

The assumption (H3) yields  $|F_{\sigma,\epsilon}(\mathbf{w})| \leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h_K^{d-1}$ . As a consequence, since  $v_K$  is a convex combination of the  $(v_\sigma)_{\sigma \in \mathcal{E}(K)}$ , we have for any  $K \in \mathcal{M}$ :

$$\left| \sum_{\sigma \in \mathcal{E}(K)} (v_\sigma - v_K) \left( \sum_{\substack{\epsilon \in \tilde{\mathcal{E}}(D_\sigma), \\ \epsilon \subset K, \epsilon = D_\sigma | D'_\sigma}} F_{\sigma,\epsilon}(\mathbf{w}) \frac{u_{\sigma'} - u_\sigma}{2} \right) \right| \leq$$

$$C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h \sum_{\sigma, \sigma', \sigma'', \sigma''' \in \mathcal{E}(K)} h_K^{d-2} |v_\sigma - v_{\sigma'}| |u_{\sigma''} - u_{\sigma'''}|,$$

and, for  $\sigma, \sigma' \in \mathcal{E}(K)$ , the quantity  $|u_\sigma - u_{\sigma'}|$  (or  $|v_\sigma - v_{\sigma'}|$ ) appears in the sum a finite number of times which depends of the dimension  $d$ . Hence, by the Cauchy-Schwarz inequality:

$$|R_1| \leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} \|u\|_{h,\text{fv}} \|v\|_{h,\text{fv}} h \leq C' \|\mathbf{w}\|_{L^\infty(\Omega)^d} \|u\|_{h,\text{b}} \|v\|_{h,\text{b}} h,$$

by Lemma 4.1. Let us now turn to  $R_2$ . By definition of  $v_K$ ,  $\sum_{\sigma \in \mathcal{E}(K)} \xi_K^\sigma (v_\sigma - v_K) = 0$ , and we obtain that:

$$R_2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} (v_\sigma - v_K) \xi_K^\sigma (u_\sigma - u_K) \left( \sum_{\sigma' \in \mathcal{E}(K)} F_{K,\sigma'}(\mathbf{w}) \right),$$

so, once again:

$$|R_2| \leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h \sum_{K \in \mathcal{M}} h_K^{d-2} \sum_{\sigma \in \mathcal{E}(K)} |v_\sigma - v_K| |u_\sigma - u_K|,$$

and, by the same arguments as for the term  $R_1$ , we conclude the proof invoking Lemma 4.1.  $\square$

### 4.4.3 Proof of Theorem 4.1

Let  $u$  be the solution of the continuous problem (4.2) which we assume to belong to  $H_0^1(\Omega) \cap H^2(\Omega)$  and let  $u_h$  be the solution of the scheme (4.8). By the triangle inequality, we have  $\|u_h - u\|_{h,b} \leq \|u_h - r_h u\|_{h,b} + \|r_h u - u\|_{h,b}$  where  $r_h$  is the interpolation operator defined in (4.18). The approximation property stated in Lemma 4.9 yields:

$$\|r_h u - u\|_{h,b} = \left( \sum_{K \in \mathcal{M}} \|\nabla(u - r_h u)\|_{L^2(K)^d}^2 \right)^{\frac{1}{2}} \leq C_2(h + \alpha_h)|u|_{H^2(\Omega)}. \quad (4.21)$$

To complete the proof, we now have to estimate the quantity  $\|u_h - r_h u\|_{h,b}$ . To this purpose, we introduce  $\|v\|_h^2 = \mathcal{A}_h(v, v) = (v, v)_h + \|v\|_{h,b}^2$ , which defines a norm on  $V_h$  that controls the broken Sobolev  $H^1$ -norm  $\|v\|_{h,b}$ . Hence, we have:

$$\|u_h - r_h u\|_{h,b} \leq \|u_h - r_h u\|_h \leq \sup_{v \in V_h} \frac{\mathcal{A}_h(u_h - r_h u, v)}{\|v\|_h}.$$

Since  $u_h$  is the solution of the scheme (4.8), we have  $\mathcal{A}_h(u_h, v) = (f, v)$  for all  $v \in V_h$ . In addition,  $u$  belongs to  $H_0^1(\Omega) \cap H^2(\Omega)$ , which implies that the strong form of the continuous problem, *i.e.* Equation (4.1a), holds in  $L^2(\Omega)$  and thus, since  $V_h \subset L^2(\Omega)$ :

$$(u, v) + (\operatorname{div}(u\mathbf{w}), v) - (\mathbf{\Delta}u, v) = (f, v), \quad \forall v \in V_h.$$

As a consequence, we have  $\mathcal{A}_h(u_h - r_h u, v) = T_1 + T_2 + T_3$  where:

$$T_1 = (u, v) - (r_h u, v)_h, \quad T_2 = (\operatorname{div}(u\mathbf{w}), v) - b_h(\mathbf{w}, r_h u, v), \quad T_3 = (-\mathbf{\Delta}u, v) - a_h(r_h u, v).$$

**Reaction term** - The term  $T_1$  may be split as follows:

$$T_1 = \underbrace{(u, v - v_c)}_{T_{1,1}} - \underbrace{\sum_{\sigma \in \mathcal{E}} |D_\sigma| v_\sigma \left( (r_h u)_\sigma - \frac{1}{|D_\sigma|} \int_{D_\sigma} u \right)}_{T_{1,2}},$$

where  $v_c$  stands for the piecewise constant function over each diamond cell and which takes the value  $v_c$  over  $D_\sigma$ , for  $\sigma \in \mathcal{E}$ . By Lemma 4.2, we get for the first term at the right-hand side of the previous relation:

$$|T_{1,1}| = |(u, v - v_c)| \leq C h \|u\|_{L^2(\Omega)} \|v\|_{h,b} \leq C h \|u\|_{L^2(\Omega)} \|v\|_h.$$

Let us now turn to the second term. We adopt the notation  $m_\omega(u) = |\omega|^{-1} \int_\omega u d\omega$  where  $\omega$  is any measurable subset of  $\mathbb{R}^d$  or  $\mathbb{R}^{d-1}$ . By definition, we have  $(r_h u)_\sigma = m_\sigma(u)$  and the second term may be decomposed as follows (invoking the boundary conditions to exclude the external faces from the summation):

$$\begin{aligned} |T_{1,2}| &= \left| \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| v_\sigma (m_\sigma(u) - m_{D_\sigma}(u)) \right| \\ &= \frac{1}{2} \left| \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} |D_\sigma| v_\sigma (m_\sigma(u) - m_{D_{K,\sigma}}(u) + m_\sigma(u) - m_{D_{L,\sigma}}(u)) \right| \\ &\leq \frac{1}{2} \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} |D_\sigma| |v_\sigma| |m_\sigma(u) - m_{D_{K,\sigma}}(u)| + \frac{1}{2} \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} |D_\sigma| |v_\sigma| |m_\sigma(u) - m_{D_{L,\sigma}}(u)|. \end{aligned}$$

We only treat the first term since the following calculations are similar for the second one. We have  $|m_\sigma(u) - m_{D_{K,\sigma}}(u)| \leq |m_\sigma(u) - m_K(u)| + |m_K(u) - m_{D_{K,\sigma}}(u)|$ . Applying Lemma 4.6, we obtain  $|m_K(u) - m_{D_{K,\sigma}}(u)| \leq 2 r_K^{-d/2} \pi^{-1} h_K \|\nabla u\|_{L^2(K)^d}$ . Lemma 4.7 provides a similar control on  $|m_\sigma(u) - m_K(u)|$  and we get, by the regularity of the mesh:

$$\begin{aligned} |T_{1,2}| &\leq C h \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma|^{1/2} |v_\sigma| \|\nabla u\|_{L^2(K)^d} \\ &\leq C h \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| |v_\sigma|^2 \right)^{\frac{1}{2}} \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} \|\nabla u\|_{L^2(K)^d}^2 \right)^{\frac{1}{2}} \\ &= C h (v, v)_h^{\frac{1}{2}} \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} \|\nabla u\|_{L^2(K)^d}^2 \right)^{\frac{1}{2}}. \end{aligned}$$

Observing that the  $H^1$  semi-norm of  $u$  on  $K$  appears at most  $2^d$  times in the summation, we obtain  $|T_{1,2}| \leq C h \|\nabla u\|_{L^2(\Omega)^d} \|v\|_h$ .

**Convection term** - Let us now turn to  $T_2$ . We may write  $T_2 = T_{2,1} + R$  where:

$$T_{2,1} = (\text{div}(u\mathbf{w}), v) - \tilde{b}_h(\mathbf{w}, r_h u, v), \quad R = \tilde{b}_h(\mathbf{w}, r_h u, v) - b_h(\mathbf{w}, r_h u, v).$$

Applying Lemma 4.10 and invoking the stability property of Lemma 4.9, we get

$$|R| \leq C h \|\mathbf{w}\|_{L^\infty(\Omega)^d} \|\nabla u\|_{L^2(\Omega)^d} \|v\|_h.$$

For the term  $T_{2,1}$ , we have:

$$T_{2,1} = (\text{div}(u\mathbf{w}), v - v_m) - \sum_{K \in \mathcal{M}} v_K \sum_{\sigma \in \mathcal{E}(K)} \int_\sigma (u - (r_h u)_\sigma) \mathbf{w} \cdot \mathbf{n}_{K,\sigma},$$

where  $v_m$  stands for the piecewise constant function over each cell  $K$  and equal to  $v_K$ . By the Cauchy-Schwarz inequality and Lemma 4.2, the first term at the right-hand side of this relation may be bounded as follow:

$$|(\text{div}(u\mathbf{w}), v - v_m)| \leq C h \|\mathbf{w}\|_{H^1(\Omega)^d} \|u\|_{H^1(\Omega)} \|v\|_h.$$

For the second term, we first remark that:

$$\begin{aligned} \left| \int_\sigma (u - (r_h u)_\sigma) \mathbf{w} \cdot \mathbf{n}_{K,\sigma} \right| &\leq \|\mathbf{w}\|_{L^\infty(\Omega)^d} |\sigma|^{1/2} \|u - m_\sigma(u)\|_{L^2(\sigma)} \\ &\leq \|\mathbf{w}\|_{L^\infty(\Omega)^d} (|\sigma|^{1/2} \|u - m_K(u)\|_{L^2(\sigma)} + |\sigma| |m_K(u) - m_\sigma(u)|) \\ &\leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h_K^{d/2} \|\nabla u\|_{L^2(K)^d}, \end{aligned}$$

by the regularity of the mesh and Lemmas 4.4, 4.5 and 4.7. Hence, reordering the summations, we thus

get (invoking the boundary conditions to exclude the external faces from the summation):

$$\begin{aligned}
& \left| \sum_{K \in \mathcal{M}} v_K \sum_{\sigma \in \mathcal{E}(K)} \int_{\sigma} (u - (r_h u)_{\sigma}) \mathbf{w} \cdot \mathbf{n}_{K,\sigma} \right| \\
& \leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} h_K^{d/2-1} \|\nabla u\|_{L^2(K)^d} |v_K - v_L| \\
& \leq C \|\mathbf{w}\|_{L^\infty(\Omega)^d} h \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} \|\nabla u\|_{L^2(K)^d}^2 \right)^{1/2} \left( \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} h_K^{d-2} |v_K - v_L|^2 \right)^{\frac{1}{2}} \\
& \leq C 2^d d \|\mathbf{w}\|_{L^\infty(\Omega)^d} \|\nabla u\|_{L^2(\Omega)^d} h \left( \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} h_K^{d-2} |v_K - v_L|^2 \right)^{\frac{1}{2}}.
\end{aligned}$$

Since  $v_K$  is a convex combination of  $(v_{\sigma})_{\sigma \in \mathcal{E}(K)}$ , we have, possibly splitting a difference  $v_{\sigma} - v_{\sigma'}$  as  $v_{\sigma} - v_{\sigma'} = (v_{\sigma} - v_{K|L}) - (v_{\sigma'} - v_{K|L})$  if  $\sigma$  and  $\sigma'$  are not faces of the same element:

$$|v_K - v_L| \leq 2 \sum_{\sigma, \sigma' \in \mathcal{E}(K)} |v_{\sigma} - v_{\sigma'}| + 2 \sum_{\sigma, \sigma' \in \mathcal{E}(L)} |v_{\sigma} - v_{\sigma'}|.$$

Therefore, since the number of edges of an element is bounded, there exists a fixed real number  $C$  such that:

$$\left( \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} h_K^{d-2} |v_K - v_L|^2 \right)^{1/2} \leq C \left( \sum_{K \in \mathcal{M}} h_K^{d-2} \sum_{\sigma, \sigma' \in \mathcal{E}(K)} |v_{\sigma} - v_{\sigma'}|^2 \right)^{1/2},$$

and, invoking Lemma 4.1, this sum is bounded by  $\|v\|_{h,b}$ , and thus by  $\|v\|_h$ . We finally get:

$$|T_{2,1}| \leq C (\|\mathbf{w}\|_{L^\infty(\Omega)^d} + \|\mathbf{w}\|_{H^1(\Omega)^d}) \|u\|_{H^1(\Omega)} \|v\|_h,$$

where  $C$  only depends on  $\theta_0$ . Combining this relation with the bound obtained for  $R$  yields that  $T_2$  satisfies the same inequality.

**Diffusion term** - We finally turn to  $T_3$ . Integrating by parts, we get:

$$T_3 = \sum_{K \in \mathcal{M}} \int_K (\nabla u - \nabla r_h u) \cdot \nabla v - \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} \int_{\sigma} v \nabla u \cdot \mathbf{n}_{K,\sigma}.$$

Thanks to Lemma 4.9, we obtain for the first term, using first the Cauchy-Schwarz inequality in  $L^2(K)$  and then the discrete Cauchy-Schwarz inequality:

$$\left| \sum_{K \in \mathcal{M}} \int_K (\nabla u - \nabla r_h u) \cdot \nabla v \right| \leq C (h + \alpha_h) \|u\|_{H^2(\Omega)} \|v\|_{h,b} \leq C (h + \alpha_h) \|u\|_{H^2(\Omega)} \|v\|_h.$$

Reordering the sums in the second term yields:

$$\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} \int_{\sigma} v \nabla u \cdot \mathbf{n}_{K,\sigma} = \sum_{\sigma \in \mathcal{E}} \int_{\sigma} [v]_{\sigma} \nabla u \cdot \mathbf{n}_{\sigma},$$

where  $[v]_{\sigma}$  and  $\mathbf{n}_{\sigma}$  stand for the jump of  $v$  through  $\sigma$  and a normal vector to  $\sigma$ , with the same orientation. Since, on a coarse face (*i.e.* the subset of a hyperplane which consists in the union of the faces sharing the

same hanging node), the normal is the same, we can group the terms in the above sum to obtain, with the notations of Lemma 4.8:

$$\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} \int_{\sigma} v \nabla u \cdot \mathbf{n}_{K,\sigma} = \sum_{\Sigma \in \mathcal{E}_c} \int_{\Sigma} [v]_{\Sigma} \nabla u \cdot \mathbf{n}_{\Sigma}.$$

We thus may apply Lemma 4.8 for  $i = 1, \dots, d$  with  $a_{\Sigma} = \mathbf{n}_{\Sigma}^i$  and  $g = \partial_i u$ , and we get:

$$\left| \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}(K)} \int_{\sigma} v \nabla u \cdot \mathbf{n}_{K,\sigma} \right| \leq C h \|u\|_{\mathbb{H}^1(\Omega)} \|v\|_h,$$

which provides the bound for  $T_3$  which we are seeking.

**Conclusion** - Collecting the bounds for  $T_1$ ,  $T_2$  and  $T_3$ , we obtain that  $\|u_h - r_h u\|_{h,b} \leq C (h + \alpha_h) \|u\|_{\mathbb{H}^2(\Omega)}$  where  $C$  only depends on  $\mathbf{w}$ ,  $\Omega$  and  $\theta_0$ . Combining this with equation (4.21) concludes the proof of Theorem 4.1.

## 5 Numerical tests

All computations presented in this section are performed with the CALIF<sup>3</sup>S free component library for fluid flows computation developed at IRSN [5].

### 5.1 A stationary incompressible flow

We first assess the behavior of the proposed numerical scheme on an exact analytical solution to the stationary incompressible Navier-Stokes equations known as the Kovasznay flow [17]. The velocity and pressure fields are given by:

$$\mathbf{u} = \begin{bmatrix} 1 - e^{\lambda x} \cos(2\pi y) \\ \frac{\lambda}{2\pi} e^{\lambda x} \sin(2\pi y) \end{bmatrix}, \quad p = \frac{1}{2} (1 - e^{2\lambda x}), \quad \lambda = \frac{1}{2\mu} - \left(\frac{1}{4\mu^2} + 4\pi^2\right)^{1/2},$$

where  $\mu$  stands for the viscosity of the flow, taken here as  $\mu = 1/40$ . The computational domain is  $\Omega = (-0.5, 1) \times (-0.5, 1.5)$ . The mesh is built from a regular  $n \times n$  grid, where we refine the sub-domain  $\Omega_f = (-0.5, 0.5) \times (-0.5, 0.5) \cup (0.5, 1) \times (0.5, 1.5)$  by splitting each (square) cell included in  $\Omega_f$  in four sub-squares. The solution is computed by the projection scheme, by letting a computed fictitious transient tend to the desired steady state. Boundary conditions are given by the analytical solution. The obtained numerical errors for various values of  $n$  are gathered in the following table, where  $\mathbf{u}_{\text{exact}}$  and  $p_{\text{exact}}$  stand for the exact velocity and pressure, respectively.

$n$	$\ \mathbf{u} - \mathbf{u}_{\text{exact}}\ _{\mathbb{L}^2(\Omega)^d}$	$\ p - p_{\text{exact}}\ _{\mathbb{L}^2(\Omega)}$
20	0.0384	0.0334
40	0.00825	0.0158
80	0.00211	0.00782
160	0.000544	0.00390

The observed order of convergence (in  $\mathbb{L}^2$ -norm) is approximately 2 for the velocity and 1 for the pressure.

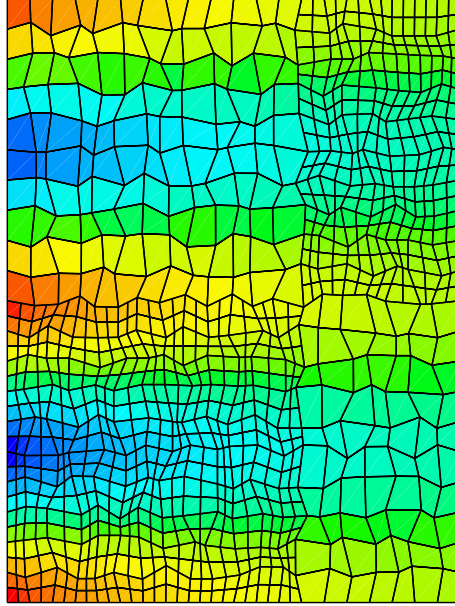


Figure 4: The unstructured mesh built by perturbation of the  $20 \times 20$  uniform grid. The background is coloured as a function of the first component of the velocity.

We now confirm this behaviour on unstructured grids, which are obtained as follows. First, we perturb the  $20 \times 20$  uniform grid, by moving each internal vertex of the mesh to a random position on a circle centered on its initial location and the radius of which is equal to 0.3 times the smallest distance between the considered vertex and its neighbours (so, for the specific unperturbed mesh used here,  $1.5/20$ ). Secondly, we build 3 refined meshes, by splitting each cell in 2, 4 and 8 respectively; by this process, the deviation from a parallelogram of each cell is divided by the same ratio, which ensures optimal convergence properties for the parametric version of the Rannacher-Turek element (see Remark 4.2). Finally, we apply local refinement to the same zones as previously. The coarsest of the obtained four meshes is plotted on Figure 4.

The obtained numerical errors are given in the following table. The orders of convergence are the same as in the uniform case.

$n$	$\ \mathbf{u} - \mathbf{u}_{\text{exact}}\ _{L^2(\Omega)^d}$	$\ p - p_{\text{exact}}\ _{L^2(\Omega)}$
20	0.0617	0.0406
40	0.0119	0.0179
80	0.00281	0.0087
160	0.000718	0.0043

The contour lines of the first component of the velocity are drawn on Fig. 5, for both meshes obtained from the  $80 \times 80$  uniform grid. We may check that no spurious perturbation appears along the lines separating the refined and non-refined parts of the computational domain (in other words, the lines composed by the union of the faces including a hanging node). A careful examination is needed to observe that, as expected, contour lines are slightly more irregular in the unstructured case.

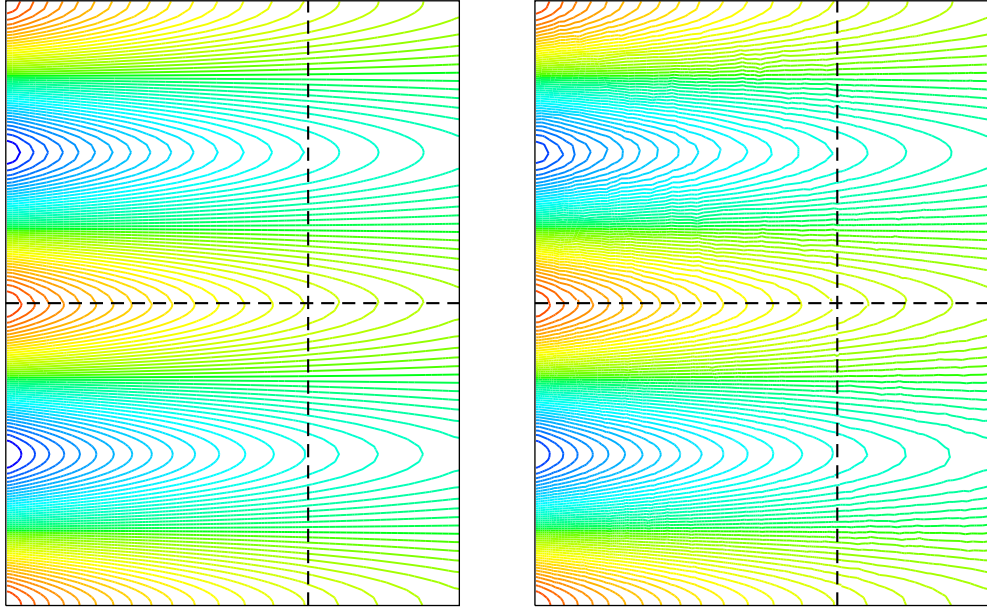


Figure 5: Contour lines of the field  $\mathbf{u}_1$ :  $80 \times 80$  uniform grid (left) and perturbed mesh (right) built by perturbation of the  $80 \times 80$  grid. The dashed lines materialize the boundary of the refined area (bottom-left and top-right sub-domains).

## 5.2 Low-Mach buoyant flows

We now turn to a natural convection flow, supposed to be represented by the asymptotic model for low Mach number flows [20], which combines System (1.1) with a balance equation for the energy:

$$\partial_t(\rho c_p \vartheta) + \operatorname{div}(\rho c_p \vartheta \mathbf{u}) - \operatorname{div}(\lambda \nabla \vartheta) = 0, \quad (5.1)$$

where  $\vartheta$  stands for the temperature, and  $c_p$  and  $\lambda$  are known positive real numbers (the heat capacity and the temperature diffusion coefficient, respectively). This system is complemented by an equation of state:

$$P_{th} = \rho R \vartheta, \quad (5.2)$$

where  $R$  stands for a constant specific to the gas under consideration (see its value below) and  $P_{th}$ , the so-called thermodynamical pressure, is supposed to depend only on time, or, in other words, to be constant in space, which has for effect to filter out the acoustic contributions from the flow field.

The computational domain is  $\Omega = (0, L)^2$ , the velocity and temperature are prescribed on the whole boundary; the velocity is set to zero and the temperature is given by:

$$\vartheta(\mathbf{x}) = \frac{L - x_1}{L} \vartheta_h + \frac{x_1}{L} \vartheta_c, \quad \vartheta_h = (1 + \varepsilon) \vartheta_0, \quad \vartheta_c = (1 - \varepsilon) \vartheta_0,$$

with  $\vartheta_0 = 600$  and  $\varepsilon = 0.6$ . The left and right vertical boundaries are thus set at a constant (hot and cold, *i.e.*  $\vartheta_h$  and  $\vartheta_c$ , respectively) temperature, while this latter varies linearly along the horizontal



boundaries, from  $\vartheta_h$  at  $x_1 = 0$  to  $\vartheta_c$  at  $x_1 = L$ . Variations of temperature are too large for the Boussinesq approximation of (1.1)-(5.1) to be valid [19].

The initial thermodynamical pressure and temperature are set to  $P_{th}(0) = 101325$  and  $\vartheta(\mathbf{x}, 0) = \vartheta_0$ ,  $\forall \mathbf{x} \in \Omega$ , and the fluid is supposed to be initially at rest. The evolution of  $P_{th}$  with time must be given by an additional relation, which, in the case of impermeability boundaries as here, may just be the conservation of the total mass in the domain:

$$\int_{\Omega} \rho(\mathbf{x}, t) d\mathbf{x} = \frac{P_{th}(t)}{R} \int_{\Omega} \frac{1}{\vartheta(\mathbf{x}, t)} d\mathbf{x} = |\Omega| \rho_0, \text{ with } \rho_0 = \frac{P_{th}(0)}{R \vartheta_0}. \quad (5.3)$$

A dimensional analysis shows that the flow is governed by two non-dimensional numbers, namely the Prandtl and the Rayleigh numbers, defined respectively by:

$$\text{Pr} = \frac{\mu c_p}{\lambda}, \quad \text{Ra} = \frac{\rho_0^2 c_p g (\vartheta_h - \vartheta_c) L^3}{\mu \lambda \vartheta_0}.$$

For the practical application performed here, we choose  $g = 9.81$  and physical properties close to those of the air:  $\mu = 1.68 \cdot 10^{-5}$ ,  $c_p = \gamma R / (\gamma - 1)$ , with  $R = 287$ . The Prandtl number is  $\text{Pr} = 0.71$ , which allows to compute the thermal diffusion  $\lambda$ , and the size of the domain  $L$  will be determined as a function of the chosen Rayleigh number.

A stability analysis of this natural convection flow is given in [27]. It shows that the flow reaches a steady state up to a critical value of the Rayleigh number approximately equal to  $\text{Ra} = 2.1 \cdot 10^6$ . Beyond this value, the flow remains time-dependent, with traveling waves along the boundaries, issued from exiting corners of the vertical boundary layers. Our aim here is to assess the capability of the proposed scheme to confirm this behaviour in the non-Boussinesq approach; an accurate determination of the critical Rayleigh number is however beyond the scope of this section.

For the solution (1.1)-(5.1), we implement a four step algorithm [4, Section 4]. We first solve (5.1) by a finite volume method, with an explicit MUSCL discretization of the convection term and an implicit approximation of the diffusion [22]. Then the thermodynamical pressure and the density are updated by (5.3) and (5.2). Finally, the velocity and the pressure are computed by the two-step pressure correction scheme (3.3).

The mesh used for this study is built as follows: we start from a  $80 \times 80$  uniform grid, perform a first refinement step splitting each cell located in  $\Omega \setminus (0.15 L, 0.85 L)^2$ , a second one by subdividing once-again the cells in  $\Omega \setminus (0.1 L, 0.9 L)^2$ , and, finally, a third one by splitting the resulting cells in  $\Omega \setminus (0.05 L, 0.95 L)^2$ ; at each refinement step, the cells are cut in 4, so *in fine* the characteristic sizes of the cells near the boundary are  $\delta x_1 = \delta x_2 = L/640$ . A mesh obtained by the same process from the  $10 \times 10$  uniform grid is shown on Figure 6. The background color in this figure is the temperature; one can see that steep variations of the solution are concentrated in the boundary layers, which justifies the chosen refinement.

At  $\text{Ra} = 2 \cdot 10^6$ , the flow after 20 time units is stationary (at least approximately, see Figure 8 below). The obtained temperature and density fields are plotted on Figure 7. One can observe that the non-linearity of the equation of state makes that the flow loses its symmetry, and generates steep density gradients at the right boundary.

We then plot the evolution with time of the temperature at the location  $\mathbf{x} = (0.1 L, 0.92 L)^t$ , for  $\text{Ra} = 2 \cdot 10^6$  and  $\text{Ra} = 2.2 \cdot 10^6$ . We observe that, for  $\text{Ra} = 2 \cdot 10^6$ , the flow tends to a steady state, while an oscillatory (quasi-periodic) behaviour subsists at  $\text{Ra} = 2.2 \cdot 10^6$ , which is consistent with the value for the critical Rayleigh known under the Boussinesq approximation [27].

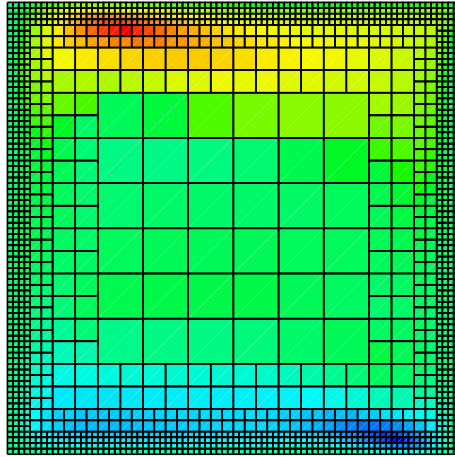


Figure 6: Typical (coarse) mesh and temperature.

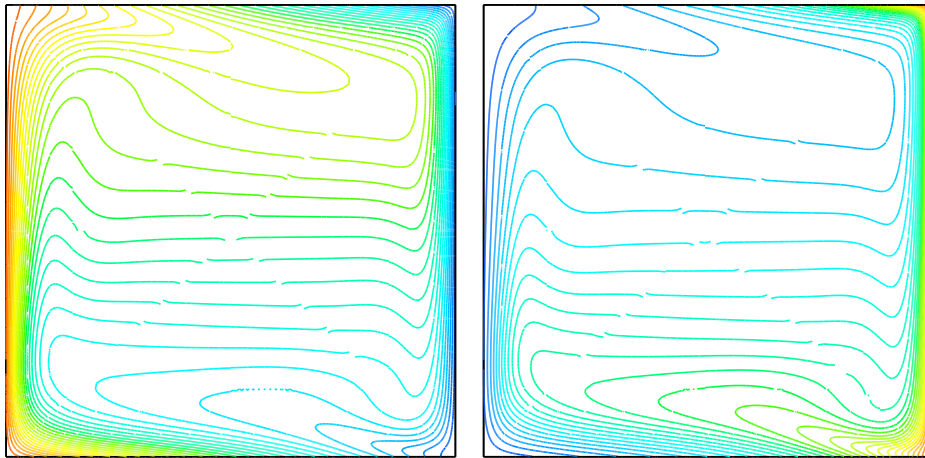


Figure 7: Temperature ( $\vartheta \in (240, 960)$ ) and density ( $\rho \in (0.33, 1.32)$ ) at the steady state, for  $Ra = 2 \cdot 10^6$ .

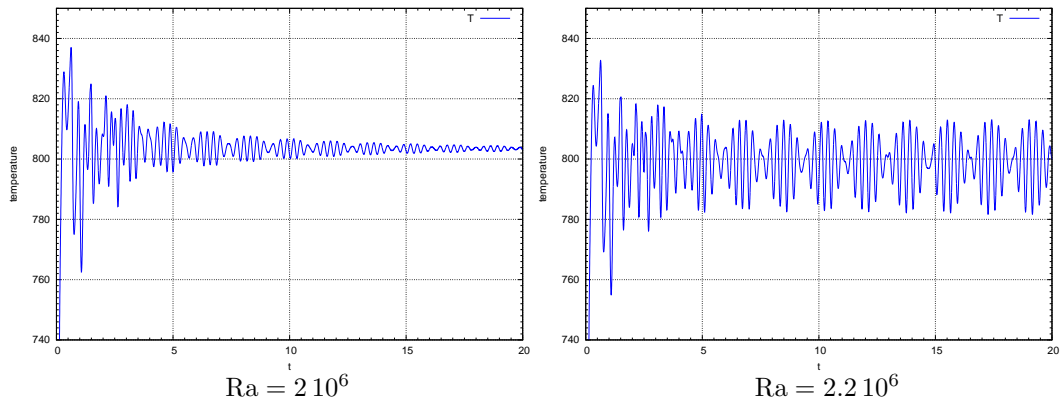


Figure 8: Temperature as a function of time at the location  $\mathbf{x} = (0.1 L, 0.92 L)^t$ , for different Rayleigh values.

## A Construction of the dual fluxes

### A.1 Dual fluxes for non-refined meshes

The system of equations (3.11) has an infinity of solutions, which makes necessary to impose in addition the constraint (H3). Since (3.11) is linear with respect to the  $F_{\sigma,\epsilon}$ ,  $\sigma \in \mathcal{E}(K)$ ,  $\epsilon \in \tilde{\mathcal{E}}(D_\sigma)$ ,  $\epsilon \subset K$ , a solution of (3.11) may thus be expressed as:

$$F_{\sigma,\epsilon} = \sum_{\sigma' \in \mathcal{E}(K)} (\alpha_K)_\sigma^{\sigma'} F_{K,\sigma'}, \quad \sigma \in \mathcal{E}(K), \epsilon \in \tilde{\mathcal{E}}(D_\sigma) \text{ and } \epsilon \subset K,$$

and (H3) is equivalent to requiring bounded coefficients  $((\alpha_K)_\sigma^{\sigma'})_{\sigma,\sigma' \in \mathcal{E}(K)}$ . In addition, since  $\xi_K^\sigma = 1/(2d)$  for all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}(K)$  (we recall that the mesh is not refined here), system (3.11) is completely independent from the cell  $K$  under consideration. We may thus consider a particular geometry for  $K$ , let us say  $K = (0, 1)^d$ , and find an expression for the coefficients  $((\alpha_K)_\sigma^{\sigma'})_{\sigma,\sigma' \in \mathcal{E}(K)}$  which we will apply to all the cells, thus automatically satisfying the constraint (H3). A technique for this computation is described in [1, Section 3.2]. The idea is to build a momentum field  $\mathbf{w}$  with a constant divergence and such that

$$\int_\sigma \mathbf{w} \cdot \mathbf{n}_{K,\sigma} = F_{K,\sigma}, \quad \forall \sigma \in \mathcal{E}(K).$$

Then, an easy computation shows that the following fluxes satisfy (3.11):

$$F_{\sigma,\epsilon} = \int_\epsilon \mathbf{w} \cdot \mathbf{n}_{\sigma,\epsilon}. \quad (\text{A.1})$$

For  $d = 2$ , using the notations introduced in Figure 9, such a momentum field  $\mathbf{w}$  is given by:

$$\mathbf{w}(x, y) = \begin{bmatrix} (1-x)(-F_W) + xF_E \\ (1-y)(-F_S) + yF_N \end{bmatrix}.$$

Using (A.1), we obtain:

$$F_{\sigma,\epsilon} = \alpha_W F_W + \alpha_E F_E + \alpha_S F_S + \alpha_N F_N,$$

with the coefficients  $\alpha_W$ ,  $\alpha_E$ ,  $\alpha_S$  and  $\alpha_N$  given in table 1. The notation  $F_{W|S}$  for the dual flux means that one calculates the flux from the western (W) to the southern (S) region with this orientation.

$F_{\sigma,\epsilon}$	$\alpha_W$	$\alpha_E$	$\alpha_S$	$\alpha_N$
$F_{W S}$	$-3/8$	$1/8$	$3/8$	$-1/8$
$F_{S E}$	$-1/8$	$3/8$	$-3/8$	$1/8$
$F_{E N}$	$1/8$	$-3/8$	$-1/8$	$3/8$
$F_{N W}$	$3/8$	$-1/8$	$1/8$	$-3/8$

Table 1: Expression of the dual fluxes in 2D.

For  $d = 3$ , using the notations introduced in Figure 10, we may choose, for the constant divergence momentum field  $\mathbf{w}$ , the following expression:

$$\mathbf{w}(x, y, z) = \begin{bmatrix} (1-x)(-F_W) + xF_E \\ (1-y)(-F_S) + yF_N \\ (1-z)(-F_B) + zF_F \end{bmatrix}.$$

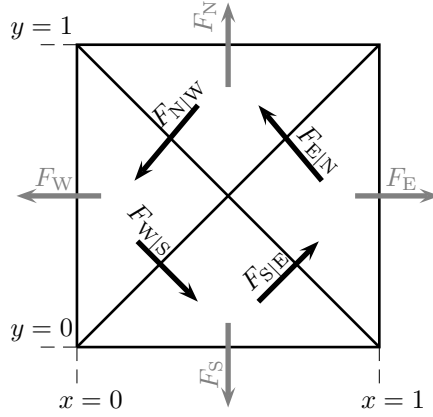


Figure 9: Notations for the primal and dual fluxes in 2D.

The dual fluxes may be expressed as linear combinations of the primal ones:

$$F_{\sigma,\epsilon} = \alpha_W F_W + \alpha_E F_E + \alpha_S F_S + \alpha_N F_N + \alpha_B F_B + \alpha_F F_F$$

where the coefficients  $\alpha_W, \alpha_E, \alpha_S, \alpha_N, \alpha_B, \alpha_F$  are given in table 2.

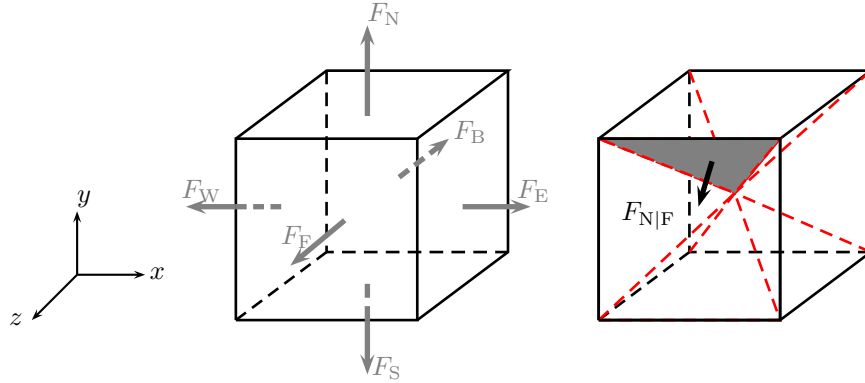


Figure 10: Notations for the primal and dual fluxes in 3D.

### A.1.1 Dual fluxes for 2D-refined meshes

Here again, we may restrict the computation to square cells. In 2D, if a primal cell is surrounded with four refined cells, the half-diamond cells are obtained by splitting the cell in four sub-squares, each one being split in two triangles. Hence, eight dual fluxes must be computed; if some of the neighboring cells are not refined, one uses a *coarsening* procedure. We begin with computing the dual fluxes across the four sub-squares faces (solid gray color in Fig. 11) so that (3.11) holds, with  $(F_{K,\sigma})_{\sigma \in \mathcal{E}(K)}$  denoted here by  $F_i$  ( $4 \leq i \leq 11$ ) and  $F_{\sigma,\epsilon}$ ,  $\sigma \in \mathcal{E}(K), \epsilon \subset K$  denoted here by  $\tilde{F}_i$  ( $4 \leq i \leq 7$ ). The linear system to solve has a one dimensional kernel and a particular solution satisfying (H3) is given by:

$F_{\sigma,\epsilon}$	$\alpha_W$	$\alpha_E$	$\alpha_S$	$\alpha_N$	$\alpha_B$	$\alpha_F$
$F_{F S}$	0	0	5/24	-1/24	1/24	-5/24
$F_{S B}$	0	0	-5/24	1/24	5/24	-1/24
$F_{B N}$	0	0	-1/24	5/24	-5/24	1/24
$F_{N F}$	0	0	1/24	-5/24	-1/24	5/24
$F_{W S}$	-5/24	1/24	5/24	-1/24	0	0
$F_{S E}$	-1/24	5/24	-5/24	1/24	0	0
$F_{E N}$	1/24	-5/24	-1/24	5/24	0	0
$F_{N W}$	5/24	-1/24	1/24	-5/24	0	0
$F_{F E}$	-1/24	5/24	0	0	1/24	-5/24
$F_{E B}$	1/24	-5/24	0	0	5/24	-1/24
$F_{B W}$	5/24	-1/24	0	0	-5/24	1/24
$F_{W F}$	-5/24	1/24	0	0	-1/24	5/24

Table 2: Expression of the dual fluxes in 3D.

$$\begin{aligned}\tilde{F}_4 &= \frac{3}{8} ( F_5 + F_6 - F_{11} - F_4 ) + \frac{1}{8} ( F_7 + F_8 - F_9 - F_{10} ), \\ \tilde{F}_5 &= \frac{3}{8} ( F_7 + F_8 - F_5 - F_6 ) + \frac{1}{8} ( F_9 + F_{10} - F_{11} - F_4 ) \\ \tilde{F}_6 &= \frac{3}{8} ( F_9 + F_{10} - F_7 - F_8 ) + \frac{1}{8} ( F_4 + F_{11} - F_5 - F_6 ) \\ \tilde{F}_7 &= \frac{3}{8} ( F_4 + F_{11} - F_9 - F_{10} ) + \frac{1}{8} ( F_5 + F_6 - F_7 - F_8 )\end{aligned}$$

Then, the dual fluxes across the diagonal faces  $\bar{F}_i$  ( $0 \leq i \leq 3$ ) (dashed gray color in Fig. 11) are computed by isolating the sub-squares and applying the procedure described above for the non-refined case. For instance,  $\bar{F}_1 = F_{E|N} - F_{W|S}$ , where  $F_E := F_7$ ,  $F_N := F_8$ ,  $F_W := \tilde{F}_6$ , and  $F_S = -\tilde{F}_5$ .

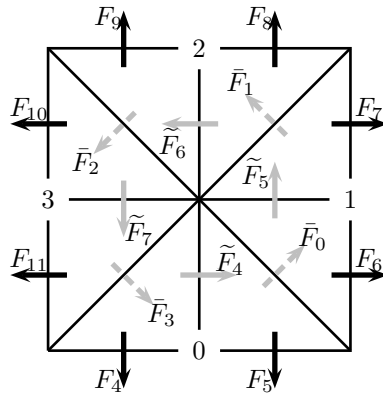


Figure 11: Dual fluxes for the neighboring cell of refined cells (2D case).

### A.1.2 Dual fluxes for 3D-refined meshes

The procedure is the same as in the 2D-case. The first step consists in splitting the cube in eight sub-cubes and computing the dual fluxes across the faces of these sub-cubes. The formula of one of these intermediate fluxes  $\tilde{F}$  (see Fig. 12) is given by:

$$\begin{aligned}
 24 \tilde{F} &= 7 ( F_{14} + F_{21} + F_{27} ) - 7 ( F_{15} + F_{19} + F_{25} ) \\
 &+ 2 ( F_7 + F_{20} + F_{26} ) - 2 ( F_6 + F_{18} + F_{24} ) \\
 &+ 2 ( F_{11} + F_{16} + F_{29} ) - 2 ( F_{13} + F_{17} + F_{23} ) \\
 &+ ( F_9 + F_{10} + F_{28} ) - ( F_8 + F_{12} + F_{22} )
 \end{aligned}$$

The computation of the other fluxes across the faces separating two sub-cubes is deduced by permutations of the indices.

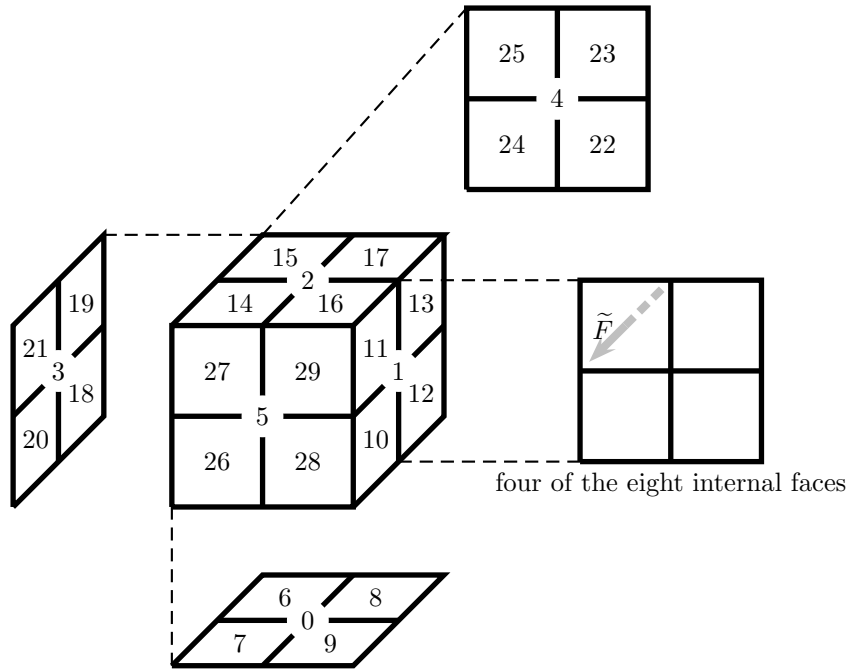


Figure 12: Intermediate dual fluxes for the neighboring cell of refined cells (3D case).

In the second step, each sub-cube is split in 3 half-diamonds of equal volumes. One obtains 24 half-diamonds and 48 internal half-diamond faces of two possible types (see Fig.13). The dual fluxes across these faces are obtained by isolating the sub-cubes and applying the procedure described above for the non-refined case, consisting in integrating the momentum field  $\mathbf{w}$  over the faces of interest.

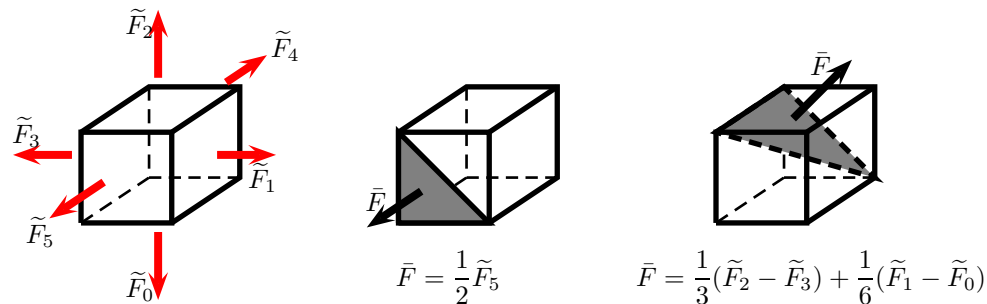


Figure 13: Two possible types of internal half-diamond faces (3D case).

## References

- [1] G. Ansanay-Alex, F. Babik, J.-C. Latché, and D. Vola. An  $L^2$ -stable approximation of the Navier-Stokes convection operator for low-order non-conforming finite elements. *International Journal for Numerical Methods in Fluids*, 66:555–580, 2011.
- [2] M. Bebendorf. A note on the Poincaré inequality for convex domains. *Journal for Analysis and its Applications*, 22:751–756, 2003.
- [3] R. Becker and S. Mao. Quasi-optimality of adaptative nonconforming finite element methods for the Stokes equations. *SIAM Journal on Numerical Analysis*, 49:970–991, 2011.
- [4] F. Boyer, F. Dardalhon, C. Lapuerta, and J.-C. Latché. Stability of a Crank-Nicolson pressure correction scheme based on staggered discretizations. *International Journal for Numerical Methods in Fluids*, 74:34–58, 2014.
- [5] CALIF<sup>3</sup>S. A software components library for the computation of reactive turbulent flows. <https://gforge.irsnn.fr/gf/project/isis>.
- [6] A.J. Chorin. Numerical solution of the Navier-Stokes equations. *Mathematics of Computation*, 22:745–762, 1968.
- [7] P. G. Ciarlet. Basic error estimates for elliptic problems. In P. Ciarlet and J.L. Lions, editors, *Handbook of Numerical Analysis, Volume II*, pages 17–351. North Holland, 1991.
- [8] R. Eymard, R. Herbin, J.-C. Latché, and B. Piar. Convergence analysis of a locally stabilized collocated finite volume scheme for incompressible flows. *Mathematical Modelling and Numerical Analysis*, 43:889–927, 2009.
- [9] T. Gallouët, L. Gastaldo, R. Herbin, and J.-C. Latché. An unconditionally stable pressure correction scheme for compressible barotropic Navier-Stokes equations. *Mathematical Modelling and Numerical Analysis*, 42:303–331, 2008.
- [10] L. Gastaldo, R. Herbin, W. Kheriji, C. Lapuerta, and J.-C. Latché. Staggered discretizations, pressure correction schemes and all speed barotropic flows. In *Finite Volumes for Complex Applications VI - Problems and Perspectives - Prague, Czech Republic*, volume 2, pages 39–56, 2011.
- [11] L. Gastaldo, R. Herbin, and J.-C. Latché. An unconditionally stable finite element-finite volume pressure correction scheme for the drift-flux model. *Mathematical Modelling and Numerical Analysis*, 44:251–287, 2010.



- [12] L. Gastaldo, R. Herbin, and J.-C. Latché. A discretization of phase mass balance in fractional step algorithms for the drift-flux model. *IMA Journal of Numerical Analysis*, 31:116–146, 2011.
- [13] D. Grapsas, R. Herbin, W. Kheriji, and J.-C. Latché. An unconditionally stable finite element-finite volume pressure correction scheme for the compressible navier-stokes equations. *submitted*, 2014.
- [14] J.L. Guermond, P. Mineev, and J. Shen. An overview of projection methods for incompressible flows. *Computer Methods in Applied Mechanics and Engineering*, 195:6011–6045, 2006.
- [15] R. Herbin, W. Kheriji, and J.-C. Latché. Consistent pressure correction staggered schemes for the shallow water and euler equations. *Under revision*, 2013.
- [16] R. Herbin, W. Kheriji, and J.-C. Latché. Pressure correction staggered schemes for barotropic one-phase and two-phase flows. *Computers and Fluids*, 88:524–542, 2013.
- [17] L.-I.-G. Kovasznay. Laminar flow behind a two-dimensional grid. *Mathematical Proceedings of the Cambridge Philosophical Society*, 44:[58], 1948.
- [18] J.-C. Latché and K. Saleh. A convergent staggered scheme for variable density incompressible Navier-Stokes equations. *submitted*, 2014.
- [19] P. Le Quéré, C. Weisman, H. Paillère, J. Vierendeels, E. Dick, R. Becker, M. Braack, and J. Locke. Modelling of natural convection flows with large temperature differences: a benchmark problem for low Mach number solvers. Part 1. Reference solutions. *Mathematical Modelling and Numerical Analysis*, 2005.
- [20] A. Majda and J. Sethian. The derivation and numerical solution of the equations for zero Mach number solution. *Combustion Science and Techniques*, 42:185–205, 1985.
- [21] L.E. Payne and H.F. Weinberger. An optimal Poincaré-inequality for convex domains. *Archive for Rational Mechanics and Analysis*, 5:286–292, 1960.
- [22] L. Piar, F. Babik, R. Herbin, and J.-C. Latché. A formally second order cell centered scheme for convection-diffusion equations on general grids. *International Journal for Numerical Methods in Fluids*, 71:873–890, 2013.
- [23] R. Rannacher and S. Turek. Simple nonconforming quadrilateral Stokes element. *Numerical Methods for Partial Differential Equations*, 8:97–111, 1992.
- [24] F. Schieweck and L. Tobiska. A nonconforming finite element method of upstream type applied to the stationary Navier-Stokes equation. *Mathematical Modelling and Numerical Analysis*, 23:627–647, 1989.
- [25] F. Schieweck and L. Tobiska. An optimal order error estimate for an upwind discretization of the Navier-Stokes equations. *Numerical Methods for Partial Differential Equations*, 12:407–421, 1996.
- [26] R. Temam. Sur l’approximation de la solution des équations de Navier-Stokes par la méthode des pas fractionnaires II. *Arch. Rat. Mech. Anal.*, 33:377–385, 1969.
- [27] S. Xin and P. Le Quéré. Linear stability analyses of natural convection flows in a differentially heated square cavity with conducting horizontal walls. *Physics of Fluids*, 13:2529–2542, 2001.