

A discretization of phase mass balance in fractional step algorithms for the drift-flux model

LAURA GASTALDO¹, RAPHAÈLE HERBIN², JEAN-CLAUDE LATCHÉ¹

¹ *Institut de Radioprotection et de Sûreté Nucléaire (IRSN).*

BP3 - 13115 St Paul-lez-Durance cedex, France.

email: [laura.gastaldo, jean-claude.latche]@irsn.fr

² *Université de Provence.*

CMI, 39 rue Frédéric Joliot-Curie, 13453 Marseille cedex 13, France.

email: herbin@cmi.univ-mrs.fr

[Received on xx October 2007]

We address in this paper a parabolic equation used to model the phase mass balance in two-phase flows, which differs from the mass balance for chemical species in compressible multi-component flows by the addition of a nonlinear term of the form $\nabla \cdot \rho \phi(y) u_r$, where y is the unknown mass fraction, ρ stands for the density, $\phi(\cdot)$ is a regular function such that $\phi(0) = \phi(1) = 0$ and u_r is a (not necessarily divergence free) velocity field. We propose a finite-volume scheme for the numerical approximation of this equation, with a discretization of the nonlinear term based on monotone flux functions. Under the classical assumption that the discretization of the convection operator must be such that it vanishes for constant y , we prove the existence and uniqueness of the solution, together with the fact that it remains within its physical bounds, *i.e.* within the interval $[0, 1]$. Then this scheme is combined with a pressure correction method to obtain a semi-implicit fractional-step scheme for the so-called drift-flux model. To satisfy the above-mentioned assumption, a specific time-stepping algorithm with particular approximations for the density terms is developed. Numerical tests are performed to assess the convergence and stability properties of this scheme.

Keywords: Two-phase flows, drift-flux model, finite volume methods, monotone schemes

1. Introduction

This paper addresses a class of physical problems which can be stated in the form of the Navier-Stokes equations, supplemented by the balance equation of an independent unknown field y :

$$\partial_t \rho + \nabla \cdot (\rho u) = 0, \quad (1.1a)$$

$$\partial_t (\rho u) + \nabla \cdot (\rho u \otimes u) + \nabla p - \nabla \cdot \tau = f, \quad (1.1b)$$

$$\partial_t (\partial \rho y) + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \phi(y) u_r) = \nabla \cdot (D \nabla y), \quad (1.1c)$$

$$\rho = \eta(p, y), \quad (1.1d)$$

where t stands for time, u for the fluid velocity, p for the pressure and ρ for the fluid density. The tensor τ is the viscous part of the stress tensor, given by the following expression:

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

where μ is the fluid viscosity and I stands for the identity tensor. For a constant viscosity, this relation yields:

$$-\nabla \cdot \tau = -\mu \left[\Delta u + \frac{1}{3} \nabla \nabla \cdot u \right]. \quad (1.3)$$

and, in this case, this term is coercive. The function η , which gives the density as an explicit function of y and the pressure, is obtained from the equation of state of the considered fluid. The nonlinear function ϕ is such that $\phi \in C^1([0, 1], \mathbb{R})$ and $\phi(0) = \phi(1) = 0$; for physical reasons, y is supposed to satisfy $0 \leq y \leq 1$, so that $\phi(\cdot)$ can be extended by continuity to $\mathbb{R} \setminus [0, 1]$ by 0 without altering the model. The volumic diffusion coefficient D and the velocity field u_r are known quantities. The problem is posed over an open bounded connected polygonal subset Ω of \mathbb{R}^d , $d \leq 3$, and over a finite time interval $(0, T)$. It must be supplemented by suitable boundary conditions, and initial conditions for ρ , u and y . In the sequel, we shall assume that $f \in L^2(\Omega)$.

Several physical problems enter this framework. For instance, taking for y the gas mass fraction, and for ρ the mixture density of dispersed two-phase flows yields the so-called drift-flux model, in the isothermal case. In this case, u_r is the relative velocity between the two phases and ϕ is given by $\phi(y) = y(1 - y)$. Dispersed two-phase flows and, in particular, bubbly flows are widely encountered in industrial applications; one may think, in particular, about bubble columns and airlift reactors, where the agitation due the gaseous phase is used to promote the contact and consequently the chemical reactions between chemical species in the flow. They are also of wide concern in the framework of nuclear safety studies, either for the modelling of boiling of water in the primary coolant circuit in case of an accidental depressurization or for the simulation of the late phases of a core-melt accident, when the flow of molten core and vessel structures comes to chemically interact with the concrete of the containment floor. This is the context of the present work.

When designing a numerical scheme for the solution of system (1.1), one faces at least two difficulties. First, the unknown y is expected, from both physical and mathematical reasons, to remain in the interval $[0, 1]$. Similarly, the unknown ρ should remain positive at all times. This suggests to build a numerical scheme that reproduces these features at the discrete level. This is performed by discretizing the mass balance equation (1.1a) with an upwind finite volume scheme and combining a monotone flux approach [13, section 21] for the term $\nabla \cdot (\rho \phi(y) u_r)$ in the advection diffusion equation (1.1c) with the argument introduced by Larrouturou [18]: the discrete counterpart of the advection operator $\partial \rho y / \partial t + \nabla \cdot (\rho y u)$ satisfies a maximum principle provided that this operator applied to a constant value of y vanishes, *i.e.* that a discrete version of the mass balance is satisfied. We first prove that, with the proposed scheme, the variable y is kept within its expected range $[0, 1]$. By a topological degree argument, this yields the existence of a discrete solution, which is then shown to be unique by a duality argument.

Now even if the model at hand represents a compressible flow, the liquid is in fact almost incompressible, so that zones may appear in the flow where the velocity of acoustic waves is very large, and the Mach number is accordingly very small. We thus need to design a numerical method which is stable in the low Mach number limit, and therefore able to deal with incompressible flows. For this purpose, we use a fractional step algorithm of the class from finite element projection methods, which are widely used for incompressible flows, see *e.g.* [17, 19] and references herein. An extension to the barotropic Navier-Stokes equations of a scheme closely related to the one developed here can be found in [14], together with references to related works. For stability reasons, the spatial discretization must preferably be based on pairs of velocity and pressure approximation spaces satisfying the so-called *inf-sup* or Babuška-Brezzi condition (*e.g.* [4]). Among these elements, nonconforming approximations with

degrees of freedom for the velocity located at the center of the faces seem to be well suited to coupling with a finite volume method for the advection-diffusion of y ; this is the choice made here. The fractional step approach is extended to the entire scheme, and the whole set of equations is thus solved in sequence. As a consequence, to ensure both conservativity and the above-mentioned monotonicity condition for the computation of y , a particular time stepping must be developed.

This paper is organized as follows. The finite volume scheme for the transport equations (1.1a) and (1.1c) are described and analysed in section 2. The fractional step algorithm for the solution of the whole problem is presented in section 3, along with the nonconforming finite element discretization of equation (1.1b). We show how the compatibility between the discretization of the different steps of the algorithm allows one to prove the well-posedness and the stability properties of the fully discrete algorithm. Numerical tests are reported in section 4; first a problem exhibiting an analytical solution allows us to assess convergence properties of the discretization; then a physical situation consisting of a phase separation problem under gravity is addressed.

2. A finite volume scheme for the nonlinear advection-diffusion equation

In this section, we present a finite volume discretization of the advection diffusion equation (1.1c). More precisely, the problem that we address here is the following: supposing that the velocity field u and the density ρ are known and satisfy a discrete mass balance equation of a particular form, we build a scheme for the computation of y which enjoys the property that the unknown y stays in the interval $[0, 1]$. We thus obtain a "building block" of a fractional step algorithm, which may be implemented for the solution of a wide range of problems, and exhibits the practically relevant property of keeping y within its physical bounds.

This section is structured as follows. We present in Paragraph 2.1 the finite volume mesh and discretization space considered. In Paragraph 2.2, we give the discretization of the mass balance equation (1.1a) which is used in Section 3 for the solution of the complete problem (1.1). It also serves as an example of a general form of the mass balance equation for ρ and u which is required for the stability analysis of the discretization of the nonlinear advection diffusion equation (1.1c); in particular, we show that it preserves the positivity of the density. We then address in Paragraph 2.3 the stability analysis of the finite volume discretization of the advection diffusion (1.1c) provided a specified general discrete mass balance is satisfied.

2.1 Discretization mesh and spaces

In this section, we denote by \mathcal{T} a set of nonintersecting convex subdomains of Ω , such that:

- (i) $\bar{\Omega} = \bigcup_{K \in \mathcal{T}} \bar{K}$.
- (ii) There exists a set \mathcal{E} of bounded subsets of hyperplanes of \mathbb{R}^d included in $\bar{\Omega}$, which are the edges (in 2D) or faces (in 3D) of the cells $K \in \mathcal{T}$. The set of boundary edges or faces (*i.e.* the edges or faces included in the boundary $\partial\Omega$ of the domain Ω) is denoted by \mathcal{E}_{ext} and the set of internal ones (*i.e.* $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$) is denoted by \mathcal{E}_{int} . If $K, L \in \mathcal{T}$, we suppose that either $\bar{K} \cap \bar{L} = \emptyset$, $\bar{K} \cap \bar{L}$ is a vertex or $\bar{K} \cap \bar{L} \in \mathcal{E}_{\text{int}}$, and, in the latter case, this common edge or face of K and L is denoted by $K|L$.

For the discretization of the advection diffusion equation (1.1c), the following additional orthogonality condition is assumed to hold:

- (iii) There exists a family $\mathcal{P} = (x_K)_{K \in \mathcal{T}}$ of points of Ω such that $x_K \in \bar{K}$ for all $K \in \mathcal{T}$ and, if $\sigma = K|L$, $x_K \neq x_L$ and the straight line passing through x_K and x_L is orthogonal to σ .

Even though this condition is not needed in the present analysis because we are only concerned here with stability issues, it enables us to perform a consistent approximation of the normal diffusion fluxes at the boundaries of the cells, which in turn is known, in simpler cases, to yield the convergence of the scheme (see e.g. [13, Chapter 2] for the convergence proof for a linear steady diffusion problem).

The set of edges (in 2D) or faces (in 3D) of a cell K of \mathcal{T} is denoted by $\mathcal{E}(K)$. For each internal edge or face of the mesh $\sigma = K|L$, n_{KL} stands for the normal vector to σ , oriented from K to L (so $n_{KL} = -n_{LK}$). By $|K|$ and $|\sigma|$, we denote the measure of the control volume K and of the edge or face σ , respectively. For any $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}(K)$, we denote by $d_{K,\sigma}$ the Euclidean distance between x_K and σ . For any $\sigma \in \mathcal{E}$, we define $d_\sigma = d_{K,\sigma} + d_{L,\sigma}$, if $\sigma \in \mathcal{E}_{\text{int}}$ (in which case d_σ is the Euclidean distance between x_K and x_L) and $d_\sigma = d_{K,\sigma}$ if $\sigma \in \mathcal{E}_{\text{ext}}$.

We denote by $X_{\mathcal{T}}$ the space of piecewise constant functions on each control volume $K \in \mathcal{T}$:

$$X_{\mathcal{T}} = \{q \in L^2(\Omega) : q|_K = \text{constant}, \forall K \in \mathcal{T}\}. \quad (2.1)$$

Let $N = \text{card}(\mathcal{T})$; any function $q \in X_{\mathcal{T}}$ can be defined by the data of the N values of q over the elements of \mathcal{T} , and hereafter we somewhat improperly identify the function itself and this family of real numbers, therefore allowing such expressions as $q \in X_{\mathcal{T}}$, $q = (q_K)_{K \in \mathcal{T}}$. The gas mass fraction y is approximated by functions from $X_{\mathcal{T}}$, so $y = (y_K)_{K \in \mathcal{T}}$.

Finally, we define $M = \text{card}(\mathcal{E}_{\text{int}})$ and, throughout this paper, for any real number a , we define $a^+ = \max(a, 0)$ and $a^- = -\min(a, 0)$, so that $a = a^+ - a^-$ with $a^+ \geq 0$ and $a^- \geq 0$.

2.2 Space discretization of the mass balance equation (1.1a)

We consider here a first-order backward Euler time discretization of the mass balance equation (1.1a), which reads in the semi-discrete form:

$$\frac{\rho - \rho^*}{\delta t} + \nabla \cdot (\rho u) = 0. \quad (2.2)$$

Assuming a known velocity field u and initial density ρ^* , we discretize the density field ρ by a finite volume method on the mesh \mathcal{T} described above. We assume that from the discrete velocity field u , we are able to derive a family of values representing the velocity on each internal edge or face $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^{M \times d}$ (for the sake of simplicity, we suppose that the velocity obeys homogeneous Dirichlet boundary conditions, and thus that its value over the external edges vanishes); note that these quantities are readily obtained for instance when using a Crouzeix-Raviart or Rannacher-Turek discretization as in Section 3.2. Hence, for the known families $\rho^* \in X_{\mathcal{T}}$ and $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^{M \times d}$, we look for $\rho \in X_{\mathcal{T}}$ that is a solution of the following upwind finite volume scheme:

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} [\rho_K - \rho_K^*] + \sum_{\sigma=K|L} F_{\sigma,K}^{up}(\rho, u) = 0, \quad (2.3)$$

where $F_{\sigma,K}^{up}(\rho, u)$ is taken to be the upwind mass flux with respect to u through the interface σ and is defined by:

$$\begin{aligned} \forall u = (u_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}} \in (\mathbb{R}^d)^M, \quad \forall \rho = (\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N, \\ F_{\sigma,K}^{up}(\rho, u) = (|\sigma| u_\sigma \cdot n_{KL})^+ \rho_K - (|\sigma| u_\sigma \cdot n_{KL})^- \rho_L, \quad \forall K \in \mathcal{T}, \quad \forall \sigma = K|L. \end{aligned} \quad (2.4)$$

Note that (2.4) can be rewritten as $F_{\sigma,K}^{up}(\rho, u) = |\sigma| u_\sigma \cdot n_{KL} \rho_\sigma^{up}$, with

$$\rho_\sigma^{up} = \begin{cases} \rho_K & \text{if } u_\sigma \cdot n_{KL} \geq 0, \\ \rho_L & \text{otherwise.} \end{cases} \quad (2.5)$$

The flux thus defined satisfies the so-called local conservativity property, which is one of the key features of a finite volume method, and which we can express here as:

$$F_{\sigma,K}^{up}(\rho, u) + F_{\sigma,L}^{up}(\rho, u) = 0. \quad (2.6)$$

Moreover, this scheme ensures that condition (2.8a) below holds, *i.e.* that the density stays positive at all times, as we now show.

LEMMA 2.1 Let $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^{M \times d}$ be two given families such that $\rho_K^* > 0$ for any $K \in \mathcal{T}$. Then the linear system (2.3) has a unique solution that satisfies $\rho_K > 0$ for any $K \in \mathcal{T}$.

Proof. With a natural equation ordering, the matrix of the linear system (2.3) is of the form $I + A$ where I denotes the (N, N) identity matrix and $A = (a_{i,j})_{1 \leq i,j \leq N}$. It is easy to check that thanks to the upwind choice (2.4), one has: $a_{i,i} \geq 0$ for $1 \leq i \leq N$, $a_{i,j} \leq 0$ for $1 \leq i, j \leq N$, $i \neq j$, and $\sum_{i=1}^N a_{i,j} = 0$ for $1 \leq j \leq N$ (note that the sum is over the lines and not the columns). Hence $I + A$ is an M-matrix. This yields that the matrix $I + A$ is invertible and that if $\rho_K^* > 0$ for all $K \in \mathcal{T}$ then $\rho_K > 0$. \square

2.3 Discretization of the nonlinear advection-diffusion equation (1.1c)

We now turn to the discretization of the balance equation (1.1c); in a semi-discrete form obtained by a first order backward Euler time discretization, it reads:

$$\frac{\rho y - \rho^* y^*}{\delta t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \phi(y) u_r^*) = \nabla \cdot (D \nabla y), \quad \text{in } \Omega \times (0, T), \quad (2.7)$$

where the density fields ρ and ρ^* , the beginning-of-step mass fraction y^* and the velocity fields u and u_r^* are supposed to be known quantities here. For the sake of simplicity, we assume in this section that both u and u_r^* vanish on the boundary $\partial\Omega$ of the computational domain, and that y obeys a homogeneous Neumann condition on the whole boundary; however, it is clear from the subsequent developments that similar results would hold for a Dirichlet boundary condition, provided the boundary datum lies in the interval $[0, 1]$.

Now at the fully discrete level, we assume that $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and a family of fluxes $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ are given, and the following relations are satisfied:

$$\forall K \in \mathcal{T}, \quad \rho_K^* > 0, \quad \rho_K > 0, \quad (2.8a)$$

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} [\rho_K - \rho_K^*] + \sum_{\sigma=K|L} F_{\sigma,K} = 0, \quad (2.8b)$$

$$\forall \sigma = K|L \in \mathcal{E}_{\text{int}}, \quad F_{\sigma,K} = -F_{\sigma,L}. \quad (2.8c)$$

This set of relations may for instance be derived by the upwind finite volume discretization of the mass balance equation (2.3)-(2.5) with $F_{\sigma,K} = F_{\sigma,K}^{up}(\rho, u)$, as in Section 3 below, thanks to (2.6) and Lemma 2.1. However, this is by far not the only way (2.8) may be obtained. For instance, the density

ρ may be not a natural unknown for the problem: indeed, if instead of considering a compressible (gaseous) dispersed phase, we consider an incompressible (solid suspension) phase, the density is then a function of y alone, and its positivity is ensured from the expression of the equation of state, provided that $0 \leq y \leq 1$; the complete problem essentially becomes an incompressible one, and the mass balance may be seen as a constraint which requires the presence of a (dynamic) pressure in the momentum balance to be satisfied, and for which a centered discretization is natural (see [1] for the treatment of such a system). In some cases (and especially for incompressible flow problems where the mass balance, *i.e.* the divergence free constraint, must have a specific discretization to ensure the stability of the scheme), the finite volume mesh for the computation of y may even be different from the mesh used for the discretization of the mass balance; for instance, one will find in [5] a way to derive from a (possibly high-order) mixed finite element solution of the momentum and mass balance equations an approximation for the velocity satisfying the mass conservation over a dual vertex-centered mesh.

With the previous notations, the discrete problem with unknown $y = (y_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$ considered in this section reads:

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} y_\sigma + \sum_{\sigma=K|L} G_{\sigma,K} \Phi_\sigma(y_K, y_L) + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (y_K - y_L) = 0, \quad (2.9)$$

where

- $y_\sigma = y_K$ if $F_{\sigma,K} \geq 0$ and $y_\sigma = y_L$ otherwise (upwind choice),
- the quantity $G_{\sigma,K}$ stands for the mass flux (*i.e.* the analogue to $F_{\sigma,K}$) associated to the relative velocity u_r^* :

$$G_{\sigma,K} = \rho_\sigma \int_\sigma u_r^* \cdot n_{KL}$$

with any reasonable approximation for the density ρ_σ on σ ; for instance $\rho_\sigma = \frac{1}{2}(\rho_K + \rho_L)$ (centered choice), or $\rho_\sigma = \rho_K$ if $F_{\sigma,K} \geq 0$ and $\rho_\sigma = \rho_L$ otherwise (upwind choice with respect to $F_{\sigma,K}$),

- $\Phi_\sigma(y_K, y_L)$ stands for $g(y_K, y_L)$ if $G_{\sigma,K} \geq 0$ and for $g(y_L, y_K)$ otherwise, g being a monotone numerical flux function with respect to ϕ in the sense of the following definition (see [13] for the theory and some examples).

Definition 2.1 [Monotone numerical flux function] Let the function $g \in C(\mathbb{R}^2, \mathbb{R})$ satisfy the following assumptions:

1. $g(a, b)$ is nondecreasing with respect to a and nonincreasing with respect to b , for any real numbers a and b ;
2. g is Lipschitz-continuous with respect to both variables over \mathbb{R} ;
3. $g(a, a) = \phi(a)$, for any $a \in \mathbb{R}$.

Then g is said to be a monotone numerical flux function with respect to the function ϕ .

Note that thanks to the definition of $G_{\sigma,K} \Phi_{\sigma}(y_K, y_L)$, the conservativity property holds, that is:

$$\forall \sigma = K|L \in \mathcal{E}_{\text{int}}, \quad G_{\sigma,K} \Phi_{\sigma}(y_K, y_L) = -G_{\sigma,L} \Phi_{\sigma}(y_K, y_L). \quad (2.10)$$

Moreover, it is easily seen that $G_{\sigma,K} \Phi_{\sigma}(y_K, y_L)$ is nondecreasing with respect to y_K and nonincreasing with respect to y_L , Lipschitz-continuous with respect to both variables over \mathbb{R} , and that, for any real number a :

$$G_{\sigma,K} \Phi_{\sigma}(a, a) = \rho_{\sigma} \int_{\sigma} u_r^* \cdot n_{KL} \phi(a). \quad (2.11)$$

The result proven in this section is the following.

THEOREM 2.2 (EXISTENCE, UNIQUENESS AND L^{∞} BOUNDS FOR A DISCRETE SOLUTION)

Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(\Phi_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Let g be a monotone numerical flux function such that $\phi(x) = g(x, x)$ vanishes for $x \leq 0$ and $x \geq 1$. Then, if $y_K^* \in [0, 1]$, $\forall K \in \mathcal{T}$, there exists a unique solution to the discrete problem (2.9), which satisfies $y_K \in [0, 1]$, $\forall K \in \mathcal{T}$.

This theorem summarizes a series of lemmata, which are detailed in the following subsections: first (section 2.3.1), we prove an *a priori* L^{∞} estimate for y , more precisely, we show the inequalities $0 \leq y(x) \leq 1$, $\forall x \in \Omega$; then, on the basis of this bound, we apply a topological degree technique to obtain the existence of a solution (section 2.3.2); finally, the solution is shown to be unique (section 2.3.3).

2.3.1 An L^{∞} stability property From a physical point of view, for instance thinking of the field y as mass fraction, it seems natural for y to satisfy an “ L^{∞} stability property”, more specifically to remain at any time in the interval $[0, 1]$. The aim of this section is to prove that this property holds for the solution of the scheme (2.9), provided that (2.8) holds at each time step and the initial condition for y takes its values in $[0, 1]$.

Let us first review the proof for the continuous problem, assuming all functions to be regular enough for the following calculations to make sense. Starting from the equation satisfied by y :

$$\partial_t(\rho y) + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \phi(y) u_r) = \nabla \cdot (D \nabla y),$$

we first prove that $y \geq 0$. Multiplying the previous equation by $-y^-$ and integrating over Ω yields:

$$-\int_{\Omega} \partial_t(\rho y) y^- - \int_{\Omega} \nabla \cdot (\rho y u) y^- - \int_{\Omega} \nabla \cdot (\rho \phi(y) u_r) y^- - \int_{\Omega} D \nabla y \cdot \nabla y^- = 0. \quad (2.12)$$

Consider the first two terms of the previous relation, *i.e.* the terms associated to the advection operator:

$$T_{\text{adv}} = -\int_{\Omega} \partial_t(\rho y) y^- - \int_{\Omega} \nabla \cdot (\rho y u) y^-.$$

Expanding the derivatives, we obtain the so-called non-conservative form of the equation, and using the fact that when y^- is non-zero, $y = -y^-$, we obtain:

$$T_{\text{adv}} = \int_{\Omega} [\partial_t \rho + \nabla \cdot (\rho u)] [y^-]^2 - \int_{\Omega} [\rho y^- \partial_t y + (\rho u y^-) \cdot \nabla y].$$

The first term vanishes because of the mass balance equation, and the second term reads:

$$-\int_{\Omega} [\rho y^- \partial_t y + (\rho u y^-) \cdot \nabla y] = \frac{1}{2} \int_{\Omega} [\rho \partial_t ((y^-)^2) + (\rho u) \cdot \nabla (y^-)^2].$$

Hence, integrating by parts and using once again the mass balance equation, we get:

$$\begin{aligned} T_{\text{adv}} &= \frac{1}{2} \int_{\Omega} [\rho \partial_t ((y^-)^2) - (y^-)^2 \nabla \cdot (\rho u)] \\ &= \frac{1}{2} \int_{\Omega} \left[\rho \partial_t ((y^-)^2) + (y^-)^2 \frac{\partial \rho}{\partial t} \right] = \frac{1}{2} \partial_t \left(\int_{\Omega} \rho (y^-)^2 \right). \end{aligned} \quad (2.13)$$

Substituting the term T_{adv} in the relation (2.12) yields:

$$\frac{1}{2} \partial_t \left(\int_{\Omega} \rho (y^-)^2 \right) - \int_{\Omega} \nabla \cdot (\rho \phi(y) u_r) y^- + \int_{\Omega} D |\nabla y^-|^2 = 0.$$

Since $\phi(x)$ vanishes for $x \leq 0$, the second integral vanishes and we have:

$$\frac{1}{2} \partial_t \left(\int_{\Omega} \rho (y^-)^2 \right) = -D \int_{\Omega} |\nabla y^-|^2 \leq 0.$$

Thus y is nonnegative, provided that the initial condition for y is non-negative. Considering the equation satisfied by $\tilde{y} = 1 - y$, one may prove similarly that $y \leq 1$.

The proof we give in the discrete setting closely follows this calculation. The first step is thus to obtain an estimate for the terms related to the advection operator, which can be seen as a discrete counterpart of relation (2.13); this is achieved by the following lemma. As the mass balance plays a central role in the continuous setting, it is natural that the discrete mass balance (2.8b) plays a similarly important role in the discrete setting.

LEMMA 2.2 Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Then, for any family $(y_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, the following property holds:

$$- \sum_{K \in \mathcal{T}} y_K^- \left[\frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} y_{\sigma} \right] \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* ((y_K^*)^-)^2], \quad (2.14)$$

where $y_{\sigma} = y_K$ if $F_{\sigma,K} \geq 0$ and $y_{\sigma} = y_L$ otherwise.

Proof. The left-hand side of (2.14) may be written as:

$$T = T_1 + T_2 \text{ with } T_1 = - \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} y_K^- (\rho_K y_K - \rho_K^* y_K^*) \text{ and } T_2 = - \sum_{K \in \mathcal{T}} y_K^- \left[\sum_{\sigma=K|L} F_{\sigma,K} y_{\sigma} \right].$$

We first remark that when y_K^- is non-zero, $y_K = -y_K^-$; hence T_1 may be written as:

$$T_1 = \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) + \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} \rho_K^* y_K^- (y_K^- + y_K^*).$$

Using the fact that $\rho_K^* y_K^- y_K^* \geq -\rho_K^* y_K^- (y_K^*)^-$, we then obtain:

$$T_1 \geq \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) + \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} \rho_K^* y_K^- [y_K^- - (y_K^*)^-];$$

thanks to the fact that $ab \leq \frac{1}{2}(a^2 + b^2)$, this yields:

$$T_1 \geq \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) + \frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} \rho_K^* [(y_K^-)^2 - ((y_K^*)^-)^2],$$

which in turn gives:

$$T_1 \geq \underbrace{\frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*)}_{T_{1,1}} + \underbrace{\frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* ((y_K^*)^-)^2]}_{T_{1,2}}. \quad (2.15)$$

We now turn to T_2 . By conservativity (Relation (2.8c)), we have:

$$T_2 = - \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} (y_K^- - y_L^-) F_{\sigma,K} y_{\sigma}.$$

Let us assume, without loss of generality that the orientation of the edges $\sigma = K|L$ is chosen so that $F_{\sigma,K} \geq 0$. Then, for any $\sigma = K|L$, we have $y_{\sigma} = y_K$, and we get that:

$$T_2 = - \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} (y_K^- - y_L^-) y_K F_{\sigma,K} = \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} (y_K^-)^2 F_{\sigma,K} + \sum_{\sigma=K|L} y_L^- y_K F_{\sigma,K}.$$

Let us then write $y_L^- y_K = \frac{1}{2}(y_L^- + y_K)^2 - \frac{1}{2}(y_L^-)^2 - \frac{1}{2}(y_K)^2$, which yields:

$$T_2 = \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} ((y_K^-)^2 - (y_L^-)^2) F_{\sigma,K} + \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} ((y_K^-)^2 - y_K^2 + (y_L^- + y_K)^2) F_{\sigma,K}. \quad (2.16)$$

The term $(y_K^-)^2 - y_K^2 + (y_L^- + y_K)^2$ is always nonnegative; indeed:

$$(y_K^-)^2 - y_K^2 + (y_L^- + y_K)^2 = \begin{cases} (y_L^- + y_K)^2 & \text{if } y_K \leq 0, \\ (y_L^-)^2 + 2y_K y_L^- & \text{otherwise.} \end{cases}$$

Hence, reordering the first summation in (2.16) and using (2.8b), we get:

$$T_2 \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma=K|L} (y_K^-)^2 F_{\sigma,K} = -\frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) = -T_{1,1}. \quad (2.17)$$

Summing (2.15) and (2.17) then yields (2.14). \square

We prove in the next two lemmas that y remains in the interval $[0, 1]$.

LEMMA 2.3 (LOWER BOUND ON y) Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Let g be a monotone numerical flux function such that $\phi(x) = g(x, x)$ vanishes for $x \leq 0$. Then, if $y_K^* \geq 0$, $\forall K \in \mathcal{T}$, the discrete solution of (2.9) also satisfies $y_K \geq 0$, $\forall K \in \mathcal{T}$.

Proof. As in the continuous case, the starting point is to multiply the equation by $-y^-$, which, in the discrete case, amounts to multiplying relation (2.9) by $-y_K^-$ and summing over the control volumes. We

get $T_{\text{adv}} + T_{\text{nl}} + T_{\text{dif}} = 0$ with:

$$\begin{aligned} T_{\text{adv}} &= \sum_{K \in \mathcal{T}} -y_K^- \left[\frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} y_\sigma \right], \\ T_{\text{nl}} &= \sum_{K \in \mathcal{T}} -y_K^- \left[\sum_{\sigma=K|L} G_{\sigma,K} \Phi_\sigma(y_K, y_L) \right], \\ T_{\text{dif}} &= D \sum_{K \in \mathcal{T}} -y_K^- \left[\sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (y_K - y_L) \right]. \end{aligned}$$

By Lemma 2.2:

$$T_{\text{adv}} \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* ((y_K^*)^-)^2].$$

Reordering the sum in the term T_{nl} by conservativity (2.10), we have:

$$T_{\text{nl}} = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \sum_{(\sigma=K|L)} T_{\text{nl},K|L} \quad \text{with} \quad T_{\text{nl},K|L} = -G_{\sigma,K} \Phi_\sigma(y_K, y_L) (y_K^- - y_L^-).$$

Let K and L be two neighbouring control volumes. If both y_K and y_L are nonnegative, the term $T_{\text{nl},K|L}$ vanishes. If $y_L \leq 0$, we get from (2.11) that $\Phi_\sigma(y_L, y_L) = \phi(y_L) = 0$, and thus:

$$T_{\text{nl},K|L} = -G_{\sigma,K} [\Phi_\sigma(y_K, y_L) - \Phi_\sigma(y_L, y_L)] (y_K^- - y_L^-).$$

Since the function $G_{\sigma,K} \Phi_\sigma(\cdot, \cdot)$ is nondecreasing with respect to the first argument and the function $x \mapsto x^-$ is non-increasing, we obtain that $T_{\text{nl},K|L} \geq 0$. Otherwise, y_K is necessarily negative and we have:

$$T_{\text{nl},K|L} = -G_{\sigma,K} [\Phi_\sigma(y_K, y_L) - \Phi_\sigma(y_K, y_K)] (y_K^- - y_L^-),$$

which is also nonnegative, since $G_{\sigma,K} \Phi_\sigma(\cdot, \cdot)$ is nonincreasing with respect to the second argument. Let us now turn to the third term. Reordering the sum, we have:

$$T_{\text{dif}} = - \sum_{\sigma \in \mathcal{E}_{\text{int}}} \sum_{(\sigma=K|L)} D \frac{|\sigma|}{d_\sigma} (y_K^- - y_L^-) (y_K - y_L) \geq 0.$$

Finally, we have:

$$\frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* [(y_K^*)^-]^2] \leq 0,$$

and thus, if $(y_K^*)_K \geq 0 \forall K \in \mathcal{T}$, then $\frac{1}{2} \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} \rho_K (y_K^-)^2 \leq 0$. □

LEMMA 2.4 (UPPER BOUND ON y) Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Let g be a monotone numerical flux function such that $\phi(x) = g(x, x)$ vanishes for $x \geq 1$. Then, if $y_K^* \leq 1$, $\forall K \in \mathcal{T}$, the discrete solution of (2.9) also satisfies $y_K \leq 1$, $\forall K \in \mathcal{T}$.

Proof. Thanks to (2.8b) and (2.9), we get the following discrete equation for the variable $1 - y$:

$$\begin{aligned} & \frac{|K|}{\delta t} [\rho_K(1 - y_K) - \rho_K^*(1 - y_K^*)] \\ & + \sum_{\sigma=K|L} \left[F_{\sigma,K}(1 - y_\sigma) - G_{\sigma,K} \Phi_\sigma(1 - (1 - y_K), 1 - (1 - y_L)) + D \frac{|\sigma|}{d_\sigma} [(1 - y_K) - (1 - y_L)] \right] = 0. \end{aligned}$$

Let $\tilde{g}(\cdot, \cdot)$ be the function defined by $\tilde{g}(a, b) = -g(1 - a, 1 - b)$. This function is nondecreasing with respect to the first variable and nonincreasing with respect to the second variable. Moreover, $\tilde{g}(x, x) = \tilde{\phi}(x) = \phi(1 - x)$ vanishes for $x \leq 0$, as $\phi(x)$ vanishes for $x \geq 1$. Thus the assumptions of Lemma 2.3 hold and, $\forall K \in \mathcal{T}$, $(1 - y)_K$ is nonnegative, which concludes the proof. \square

REMARK 2.1 (STRICT BOUNDS ON y) Strict bounds on y may also be obtained. Indeed, if $\forall K \in \mathcal{T}$, $y_K^* > 0$, then y satisfies the strict inequality $\forall K \in \mathcal{T}$, $y_K > 0$. To prove this result, let us assume that there exists $K \in \mathcal{T}$ such that $y_K = 0$. Replacing y_K by zero in the equation (2.9) of the scheme, we get:

$$\frac{|K|}{\delta t} (-\rho_K^* y_K^*) + \sum_{\sigma=K|L} \left[-F_{\sigma,K}^- y_L + G_{\sigma,K} \Phi_\sigma(0, y_L) + D \frac{|\sigma|}{d_\sigma} (-y_L) \right] = 0.$$

The first term is, by assumption, negative, while, since $y_L \geq 0$, the second and last terms are nonpositive. Since the function $s \mapsto G_{\sigma,K} \Phi_\sigma(0, s)$ is nonincreasing and $\Phi_\sigma(0, 0) = 0$, $G_{\sigma,K} \Phi_\sigma(0, y_L) \leq 0$, and the third term also is nonpositive, which contradicts the fact that the whole sum vanishes.

By applying this result to $1 - y$, we similarly prove that, if $y_K^* < 1 \forall K \in \mathcal{T}$, then $y_K < 1 \forall K \in \mathcal{T}$. Returning to the initial physical problem, this result shows that, when using this scheme for the computation of the gas mass fraction y , monophasic zones cannot appear in the flow if they are not present at the initial time.

2.3.2 Existence for the approximate solution The existence of a solution to the scheme (2.9) is obtained through a topological degree argument. We recall this result in the following theorem and refer to [10, chapter 5] for the general theory and [12, 14] for its use in the case of other nonlinear numerical schemes).

THEOREM 2.3 (APPLICATION OF THE TOPOLOGICAL DEGREE, FINITE DIMENSIONAL CASE)

Let $(V, \|\cdot\|)$ be a normed finite-dimensional vector space on \mathbb{R} , let f be a continuous function from V to V and let $b \in V$. Let us assume that there exists a continuous function \mathcal{F} from $V \times [0, 1]$ to V , and $R > 0$ satisfying:

- (i) $\mathcal{F}(\cdot, 1) = f$;
- (ii) For all $\alpha \in [0, 1]$, if v is such that $\mathcal{F}(v, \alpha) = b$ then $v \in B(0, R)$ (that is $\|v\| < R$);
- (iii) the topological degree of $\mathcal{F}(\cdot, 0)$ with respect to b and $B(0, R)$ is equal to $d_0 \neq 0$.

Then the topological degree of $\mathcal{F}(\cdot, 1)$ with respect to b and to $B(0, R)$ is also equal to $d_0 \neq 0$; consequently, there exists at least a solution $v \in B(0, R)$ such that $f(v) = 0$.

LEMMA 2.5 (EXISTENCE OF A DISCRETE SOLUTION) Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Let g be a monotone numerical flux function such that $\phi(x) = g(x, x)$ vanishes for $x \leq 0$ and $x \geq 1$. Then, if $y_K^* \in [0, 1]$, $\forall K \in \mathcal{T}$, there exists a solution to the discrete problem (2.9).

Proof. In order to apply Theorem 2.3, we consider the space $V = X_{\mathcal{T}}$ (recall that $X_{\mathcal{T}}$ is the space of functions which are piecewise constant on \mathcal{T}) equipped with the maximum norm, which we denote by $|\cdot|_{\infty}$; we then introduce a function $\mathcal{F} : X_{\mathcal{T}} \times [0, 1] \rightarrow X_{\mathcal{T}}$, such that $\mathcal{F}(\cdot, 1) = f$ where f is a function from $X_{\mathcal{T}}$ to $X_{\mathcal{T}}$ such that any solution of the nonlinear system (2.9) is a zero of f . For $y = (y_K)_{K \in \mathcal{T}} \in X_{\mathcal{T}}$ and $\alpha \in [0, 1]$, this function \mathcal{F} is defined by $\mathcal{F}(y, \alpha) = (q_K)_{K \in \mathcal{T}} \in X_{\mathcal{T}}$, with:

$$q_K = \frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} \left[F_{\sigma,K} y_{\sigma} + \alpha G_{\sigma,K} \Phi_{\sigma}(y_K, y_L) + D \frac{|\sigma|}{d_{\sigma}} (y_K - y_L) \right].$$

The function \mathcal{F} is continuous from $X_{\mathcal{T}} \times [0, 1]$ to $X_{\mathcal{T}}$. We then remark that the lemmata 2.3 and 2.4 apply to the solution of the equation $\mathcal{F}(y, \alpha) = 0$, for $0 \leq \alpha \leq 1$. Hence, any solution to this equation belongs to the open ball:

$$B(0, 2) = \{y \in X_{\mathcal{T}}, \text{ such that } |y|_{\infty} < 2\}.$$

Moreover, thanks to the estimate $y_K \leq 1$, $\forall K \in \mathcal{T}$, the linear system $\mathcal{F}(y, 0) = 0$ has a unique solution, which belongs to $B(0, 2)$ (indeed, if it had two solutions a classical argument would allow us to conclude that the set of solutions is bounded; existence follows from uniqueness since the dimension of the space is finite). From the existence of a solution to the linear system $\mathcal{F}(y, 0) = 0$, we get that the topological degree of $\mathcal{F}(\cdot, 0)$ with respect to $B(0, 2)$ and 0 is nonzero. Applying Theorem 2.3, we then deduce that the topological degree of $\mathcal{F}(\cdot, 1)$ with respect to $B(0, 2)$ and 0 is nonzero, which in turn yields that there exists at least one solution to the nonlinear system (2.9). \square

2.3.3 Uniqueness of the approximate solution The uniqueness of the solution to the scheme (2.9) is obtained through a duality argument. First, we introduce this technique in the semidiscrete time setting. Let y and \tilde{y} be two smooth functions satisfying (2.7). Then the difference $\delta y = y - \tilde{y}$ satisfies:

$$\frac{\rho \delta y}{\delta t} + \nabla \cdot (\rho \delta y u) + \nabla \cdot \left(\rho \frac{\phi(y) - \phi(\tilde{y})}{\delta y} \delta y u_r \right) = \nabla \cdot (D \nabla \delta y).$$

Multiplying by a (smooth) test function ψ and integrating over Ω , we get:

$$\frac{1}{\delta t} \int_{\Omega} \rho \delta y \psi + \int_{\Omega} \nabla \cdot (\rho \delta y u) \psi + \int_{\Omega} \nabla \cdot \left(\rho \frac{\phi(y) - \phi(\tilde{y})}{\delta y} \delta y u_r \right) \psi + D \int_{\Omega} \nabla \delta y \cdot \nabla \psi = 0. \quad (2.18)$$

We then consider the following dual problem (with unknown \bar{y} and data y , \tilde{y} and δy , and which is consequently linear):

$$\forall \psi, \frac{1}{\delta t} \int_{\Omega} \rho \bar{y} \psi + \int_{\Omega} \nabla \cdot (\rho \psi u) \bar{y} + \int_{\Omega} \nabla \cdot \left(\rho \frac{\phi(y) - \phi(\tilde{y})}{\delta y} \psi u_r \right) \bar{y} + D \int_{\Omega} \nabla \bar{y} \cdot \nabla \psi = \int_{\Omega} \delta y \psi. \quad (2.19)$$

Under some regularity assumptions, the dual problem (2.19) is known to satisfy the maximum principle (e.g. [16, chapter 8]), and then, by an application of the Fredholm alternative, to admit a unique solution. Taking as test function $\psi = \delta y$ in the dual problem, we get:

$$\frac{1}{\delta t} \int_{\Omega} \rho \bar{y} \delta y + \int_{\Omega} \nabla \cdot (\rho \delta y u) \bar{y} + \int_{\Omega} \nabla \cdot \left(\rho \frac{\phi(y) - \phi(\tilde{y})}{\delta y} \delta y u_r \right) \bar{y} + D \int_{\Omega} \nabla \bar{y} \cdot \nabla \delta y = \int_{\Omega} (\delta y)^2.$$

However, thanks to (2.18), the left-hand side of the previous relation is equal to zero, and therefore $\int_{\Omega} (\delta y)^2 = 0$ so that $\delta y = 0$, which proves the uniqueness of the solution. The proof that we give for the discrete problem is adapted from this technique.

LEMMA 2.6 (UNIQUENESS OF THE DISCRETE SOLUTION) Let $(\rho_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$, $(\rho_K^*)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $(F_{\sigma,K})_{K \in \mathcal{T}, \sigma=K|L} \in \mathbb{R}^{2M}$ satisfy (2.8). Let g be a numerical monotone flux function. Then there exists at most one solution $y \in X_{\mathcal{T}}$ to the discrete equation (2.9).

Proof. Let $y = (y_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $\tilde{y} = (\tilde{y}_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$ be two solutions of (2.9); then $\delta y = (y_K - \tilde{y}_K)_{K \in \mathcal{T}} \in \mathbb{R}^N$ satisfies, for all $K \in \mathcal{T}$:

$$\frac{|K|}{\delta t} \rho_K \delta y_K + \sum_{\sigma=K|L} \left[F_{\sigma,K} \delta y_{\sigma} + G_{\sigma,K} \left(\Phi_{\sigma}(y_K, y_L) - \Phi_{\sigma}(\tilde{y}_K, \tilde{y}_L) \right) + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_{\sigma}} (\delta y_K - \delta y_L) \right] = 0.$$

The quotients $(\Phi_{\sigma}(a, \cdot) - \Phi_{\sigma}(b, \cdot)) / (a - b)$ and $(\Phi_{\sigma}(\cdot, a) - \Phi_{\sigma}(\cdot, b)) / (a - b)$ may be extended to the case $a = b$ thanks to the fact that Φ_{σ} is Lipschitz-continuous; we may thus write the above system as:

$$\begin{aligned} \frac{|K|}{\delta t} \rho_K \delta y_K + \sum_{\sigma=K|L} \left[F_{\sigma,K} \delta y_{\sigma} + G_{\sigma,K} \frac{\Phi_{\sigma}(y_K, y_L) - \Phi_{\sigma}(\tilde{y}_K, y_L)}{\delta y_K} \delta y_K \right. \\ \left. + G_{\sigma,K} \frac{\Phi_{\sigma}(\tilde{y}_K, y_L) - \Phi_{\sigma}(\tilde{y}_K, \tilde{y}_L)}{\delta y_L} \delta y_L + D \frac{|\sigma|}{d_{\sigma}} (\delta y_K - \delta y_L) \right] = 0. \end{aligned} \quad (2.20)$$

We now introduce the following discrete dual problem (with unknown $\bar{y} \in X_{\mathcal{T}}$ and data y, \tilde{y} and $\delta y \in X_{\mathcal{T}}$):

$$\left| \begin{array}{l} \text{Find } \bar{y} \in \mathbb{R}^N \text{ such that, } \forall \psi \in \mathbb{R}^N, \\ \sum_{K \in \mathcal{T}} \frac{|K|}{\delta t} \rho_K \bar{y}_K \psi_K + \bar{y}_K \left[F_{\sigma,K} \psi_{\sigma} + G_{\sigma,K} \left(\frac{\Phi_{\sigma}(y_K, y_L) - \Phi_{\sigma}(\tilde{y}_K, y_L)}{\delta y_K} \psi_K \right. \right. \\ \left. \left. + \frac{\Phi_{\sigma}(\tilde{y}_K, y_L) - \Phi_{\sigma}(\tilde{y}_K, \tilde{y}_L)}{\delta y_L} \psi_L \right) + D \bar{y}_K \sum_{\sigma=K|L} \frac{|\sigma|}{d_{\sigma}} (\psi_K - \psi_L) \right] = \sum_{K \in \mathcal{T}} |K| (\delta y_K) \psi_K. \end{array} \right. \quad (2.21)$$

If there exists $\bar{y} \in X_{\mathcal{T}}$ satisfying (2.21), then taking as a test function $\psi = \delta y$ we get from (2.20):

$$\sum_{K \in \mathcal{T}} |K| (\delta y_K)^2 = 0,$$

which yields $y = \tilde{y}$. In order to conclude the proof of the lemma, it only remains to show that there exists a (unique) solution to the dual problem (2.21). Let A be the matrix of the linear system (2.21), obtained with the natural ordering: the line of A associated to the control volume K is obtained by taking $\psi = 1_K$, where 1_K is the characteristic function of the element K . In this line, the diagonal entry is given by the term of the sum associated to K , and the off-diagonal entries are given by the terms of the sum corresponding to control volumes sharing an edge with K . Denoting by A_{KK} this diagonal entry, we have:

$$A_{KK} = \frac{|K|}{\delta t} \rho_K + \sum_{\sigma=K|L} \left[F_{\sigma,K}^+ + G_{\sigma,K}^+ \frac{g(y_K, y_L) - g(\tilde{y}_K, y_L)}{y_K - \tilde{y}_K} - G_{\sigma,K}^- \frac{g(\tilde{y}_L, y_K) - g(\tilde{y}_L, \tilde{y}_K)}{y_K - \tilde{y}_K} + D \frac{|\sigma|}{d_{\sigma}} \right].$$

Since the density is positive and the function g is nondecreasing with respect to its first argument and nonincreasing with respect to its second argument, A_{KK} is positive. The off-diagonal entries on the same line are given by:

$$A_{KL} = -F_{\sigma,K}^+ - G_{\sigma,K}^+ \frac{g(y_K, y_L) - g(\tilde{y}_K, y_L)}{y_K - \tilde{y}_K} + G_{\sigma,K}^- \frac{g(\tilde{y}_L, y_K) - g(\tilde{y}_L, \tilde{y}_K)}{y_K - \tilde{y}_K} - D \frac{|\sigma|}{d_{\sigma}}.$$

where, in this relation, L is a neighbouring control volume of K and $\sigma = K|L$. By a similar argument, these terms are nonpositive. Moreover, the sum of all coefficient on a line reads:

$$\sum_{L \in \mathcal{T}} A_{KL} = \frac{|K|}{\delta t} \rho_K,$$

which is positive, since $\rho_K > 0$. Thus A is a strictly diagonally dominant matrix; therefore it is invertible and problem (2.21) admits a unique solution, which completes the proof. \square

3. A fractional step algorithm for dispersed two-phase flows

In this section, we address the solution of the full system (1.1). For this purpose, we build a numerical scheme by complementing the discrete nonlinear advection-diffusion equation (2.9) with an incremental-projection-like algorithm.

Let us consider a partition $0 = t_0 < t_1 < \dots < t_N = T$ of the time interval $(0, T)$, which, for the sake of simplicity, we suppose to be uniform. Let δt be the constant time step $\delta t = t_{n+1} - t_n$ for $n = 0, 1, \dots, N-1$. In a semi-discrete time setting, the proposed algorithm consists in the following three-step scheme:

1 - Solve for y^{n+1} :

$$\begin{aligned} \frac{\rho^n y^{n+1} - \rho^{n-1} y^n}{\delta t} + \nabla \cdot (\rho^n y^{n+1} u^n) \\ + \nabla \cdot (\rho^n y^{n+1} (1 - y^{n+1}) u_r^n) - \nabla \cdot (D \nabla y^{n+1}) = 0. \end{aligned} \quad (3.1)$$

2 - Solve for \tilde{u}^{n+1} :

$$\frac{\rho^n \tilde{u}^{n+1} - \rho^{n-1} u^n}{\delta t} + \nabla \cdot (\tilde{u}^{n+1} \otimes \rho^n u^n) + \nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) = f^{n+1}. \quad (3.2)$$

3 - Solve for p^{n+1} , u^{n+1} and ρ^{n+1} :

$$\rho^n \frac{u^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - p^n) = 0, \quad (3.3a)$$

$$\frac{\rho^{n+1} - \rho^n}{\delta t} + \nabla \cdot (\rho^{n+1} u^{n+1}) = 0, \quad (3.3b)$$

$$\rho^{n+1} = \eta(p^{n+1}, y^{n+1}). \quad (3.3c)$$

After a computation of the unknown y (step 1), step 2 consists in a semi-implicit solution of the momentum balance equation to obtain a predicted velocity. Step 3 is a nonlinear pressure correction step, which boils down to the usual projection step used in incompressible flow solvers when the density is constant (e.g. [19]). Taking the divergence of (3.3a) and using (3.3b) to eliminate the unknown velocity u^{n+1} yields a nonlinear elliptic problem for the pressure. This computation is formal in the semi-discrete formulation, but, of course, is completely meaningful at the algebraic level, as described in section 3.3.

Once the pressure is computed, the first relation yields the updated velocity and the third relation gives the end-of-step density.

The main difficulty in designing such an algorithm lies in the approximation of the density. Indeed, we have to meet two requirements: first, to satisfy the compatibility condition (2.8b) when computing y (*i.e.* at Step 1 of the algorithm), second to ensure the conservativity of the scheme. The first point has been shown in Section 2.3 to be necessary for a reliable computation of the unknown y , and, in our experience, a violation of this condition may be at the origin of strong instabilities, for the estimation of y itself, but also for the whole algorithm. Still in our experience, using a nonconservative scheme for the approximation of y leads to large errors. This is especially important when the equation of state is strongly nonlinear, as for flows involving phases of very different densities. To meet both requirements, we use here a time-shift of the density: in the advection terms of both the computation of y (Equation (3.1)) and the prediction of the velocity (Equation (3.2)), the density is taken one time step before the unknown y . This technique shows remarkable stability properties, but is of course limited to first order in time; this convergence property is assessed by numerical experiments.

REMARK 3.1 (ON ANOTHER TIME DISCRETIZATION OF THE DENSITY) To satisfy condition (2.8b), another way to proceed has already been proposed [1, 14]. The idea is to use, as a preliminary stage of the time step, the mass balance equation with a known value for the velocity (for instance, the velocity at the previous time step, or any extrapolation of it) to obtain a prediction of the density. If, for stability reasons, the discretization of this equation is chosen to be implicit, this preliminary step reads:

$$\frac{\bar{\rho} - \rho^n}{\delta t} + \nabla \cdot (\bar{\rho} \bar{u}).$$

where \bar{u} and $\bar{\rho}$ are the velocity used and the density obtained in this step, respectively. The first two terms of the balance equation for y (step 3.1 of the present algorithm) are now:

$$\frac{\bar{\rho} y^{n+1} - \rho^n y^n}{\delta t} + \nabla \cdot (\bar{\rho} y^{n+1} \bar{u}) \dots$$

This approach can be easily modified to obtain a (formally) second order scheme [1]. Unfortunately, it seems difficult to consider $\bar{\rho}$ as the end-of-step value for the density, as its computation does not make use of the equation of state; this scheme thus cannot be conservative. In the present context, it may however be used at the first time step (and only at the first time step), to initialize the density by the following prediction step:

$$\frac{\rho^0 - \rho^{-1}}{\delta t} + \nabla \cdot (\rho^0 u^{-1}) = 0,$$

where ρ^{-1} and u^{-1} are suitable approximations for the initial density and the velocity, respectively.

In order to obtain the full space-time discrete algorithm, we discretize (3.1) and (3.3b) by the finite volume method as described in Section 2. There remains to discretize the steps (3.2) (Paragraph 3.2 below) and (3.3a) (Paragraph 3.3) This is performed with the mixed Crouzeix-Raviart or Rannacher-Turek finite elements, which we now describe. This choice is quite convenient here since the discrete velocities are located on the edges of the mesh while the discrete pressures are piecewise constant on the cells; it is therefore compatible with the finite volume discretization described in Section 2. Furthermore, these elements are known to satisfy the inf-sup inequality, which contributes to the stability of the scheme, especially in the incompressible limit.

3.1 The Crouzeix Raviart and Rannacher-Turek finite elements

Let us briefly describe the Crouzeix-Raviart element for simplicial meshes (see [9] for the seminal paper and, for instance, [11, p. 83–85] for a synthetic presentation), and the so-called "rotated bilinear element" introduced by Rannacher and Turek for quadrilateral or hexahedric meshes [20].

In the following, we assume that the mesh of Ω which was introduced in Section 2 consists either of simplices (triangles in 2D or tetrahedra in 3D) or, in the case where the shape of Ω allows it, of rectangles or rectangular parallelepipeds. Note that conditions (i), (ii) and (iii) of Section 2.1 hold for both types of mesh, taking for x_K the intersection of the orthogonal bisectors of the edges of K . The reference element for the Crouzeix-Raviart element is the unit d -simplex and the discrete function space is the space \mathbb{P}_1 of affine polynomials. The reference element \hat{K} for the rotated bilinear element is the unit d -cube (with edges parallel to the coordinate axes); the discrete function space on \hat{K} is $\tilde{Q}_1(\hat{K})^d$, where $\tilde{Q}_1(\hat{K})$ is defined as follows:

$$\tilde{Q}_1(\hat{K}) = \text{span} \{ 1, (x_i)_{i=1,\dots,d}, (x_i^2 - x_{i+1}^2)_{i=1,\dots,d-1} \}.$$

For both velocity elements used here, the degrees of freedom are determined by the following set of nodal functionals:

$$\{m_{\sigma,i}, \sigma \in \mathcal{E}(K), i = 1, \dots, d\}, \quad m_{\sigma,i}(v) = |\sigma|^{-1} \int_{\sigma} v_i. \quad (3.4)$$

The mapping from the reference element to the actual discretization cell is the standard affine mapping for the Crouzeix-Raviart element, and the standard Q_1 mapping for the Rannacher-Turek element. Finally, in both cases, the continuity of the average value of discrete velocities (*i.e.*, for a discrete velocity field v , $m_{\sigma,i}(v)$, $1 \leq i \leq d$) across each face of the mesh is required, thus the discrete space $V_{\mathcal{T}}$ is defined as follows:

$$\begin{aligned} V_{\mathcal{T}} = \{ & v \in L^2(\Omega)^d : v|_K \in \tilde{Q}_1(K)^d, \forall K \in \mathcal{T}; \\ & m_{\sigma,i}(v|_K) = m_{\sigma,i}(v|_L), \forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \text{ for } 1 \leq i \leq d; \\ & m_{\sigma,i}(v) = 0, \forall \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq i \leq d \}. \end{aligned} \quad (3.5)$$

For both the Crouzeix-Raviart and the Rannacher-Turek discretizations, the approximation space for the pressure is the space $X_{\mathcal{T}}$ of piecewise constant functions defined by (2.1) while the approximation space for the velocity is $V_{\mathcal{T}}$. Since only the continuity of the function average over each edge of the mesh is imposed, the velocity is generally discontinuous through each edge; the discretization is thus nonconforming in $H^1(\Omega)^d$. These pairs of approximation spaces for the velocity and the pressure are *inf-sup* stable, in the usual sense for "piecewise H^1 " discrete velocities, *i.e.* there exists $c_1 > 0$ which does not depend on the mesh such that:

$$\forall p \in X_{\mathcal{T}}, \quad \sup_{v \in V_{\mathcal{T}}} \frac{\int_{\Omega, \mathcal{T}} p \nabla \cdot v}{\|v\|_{1,b}} \geq c_1 \|p - m(p)\|_{L^2(\Omega)},$$

where $m(p)$ is the mean value of p over Ω , the symbol $\int_{\Omega, \mathcal{T}}$ stands for $\sum_{K \in \mathcal{T}} \int_K$ and $\|\cdot\|_{1,b}$ stands for the broken Sobolev H^1 semi-norm:

$$\|v\|_{1,b}^2 = \int_{\Omega, \mathcal{T}} |\nabla v|^2 = \sum_{K \in \mathcal{T}} \int_K |\nabla v|^2.$$

From the definition (3.4), each velocity degree of freedom can be uniquely associated to an element edge. Hence, the velocity degrees of freedom may be indexed by the number of the component and the associated edge, and the set of velocity degrees of freedom reads:

$$\{v_{\sigma,i}, \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d\}.$$

We define $v_{\sigma} = \sum_{i=1}^d v_{\sigma,i} e^{(i)}$ where $e^{(i)}$ is the i^{th} vector of the canonical basis of \mathbb{R}^d . We denote by $\varphi_{\sigma}^{(i)}$ the vector shape function associated to $v_{\sigma,i}$, which, by the definition of the Crouzeix-Raviart and Rannacher-Turek finite elements, reads:

$$\varphi_{\sigma}^{(i)} = \varphi_{\sigma} e^{(i)},$$

where φ_{σ} is the scalar basis function.

3.2 Space discretization of the momentum equation (1.1b)

As previously mentioned, we seek discrete approximations of u and p in the Crouzeix-Raviart or Rannacher-Turek finite element spaces, which we again denote $u \in V_{\mathcal{T}}$ and $p \in X_{\mathcal{T}}$. With the notations of the previous section, we write the discrete functions as linear combinations of the basis functions: $u = \sum_{i=1}^d \sum_{\sigma \in \mathcal{E}_{\text{int}}} u_{\sigma,i} \varphi_{\sigma}^{(i)}$ and $p = \sum p_K 1_K$, where 1_K is the characteristic function K .

A direct discretization of (3.2) with this finite element pair is in fact not quite satisfactory for stability reasons. Indeed, we slightly modify it in order to obtain a finite volume type scheme on a dual mesh, which allows us to respect a discrete counterpart of the L^2 -stability of the advection operator for the velocity, *i.e.* the discrete counterpart of the following relation:

$$\int_{\Omega} \left[\frac{\partial \rho u}{\partial t} + \nabla \cdot (\rho u) \right] \cdot u = \frac{d}{dt} \int_{\Omega} \frac{1}{2} \rho |u|^2, \quad (3.6)$$

which is satisfied for any smooth functions ρ, u satisfying (1.1a)-(1.1b). This discrete counterpart is stated in theorem 3.1 below (see [14] for the proof). It is central in the proof of *a priori* estimates (*e.g.* the kinetic energy conservation theorem for incompressible flows) for the solution of the overall system, and has been found to be essential for convection dominant flows [2]. It is also one of the ingredients used in [14] to derive a pressure correction scheme for compressible barotropic flows which conserves the entropy of the system. We state it below for the dual diamond mesh \mathcal{M} which is constructed as follows: for each internal edge $\sigma = K|L$, let $D_{K,\sigma}$ be the cone with basis σ and the center of mass of the cell K taken as the opposite vertex. The volume $D_{\sigma} = D_{K,\sigma} \cup D_{L,\sigma}$ is referred to as the "diamond cell" associated to σ and $D_{K,\sigma}$ is the half-diamond cell associated to σ and K (see Figure 1). The diamond cells D_{σ} are naturally referred to by the primal edge $\sigma \in \mathcal{E}_{\text{int}}$, and their edges by the letter ε .

THEOREM 3.1 (STABILITY OF FINITE-VOLUME ADVECTION OPERATORS) Let $\mathcal{M} = (D_{\sigma})_{\sigma \in \mathcal{E}_{\text{int}}}$. Let $(\rho_{\sigma})_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^M$, $(\rho_{\sigma}^*)_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^M$ and $(F_{\varepsilon,\sigma})_{\sigma \in \mathcal{E}_{\text{int}}, \varepsilon \in \mathcal{E}(D_{\sigma})} \in \mathbb{R}^{kM}$ (where k is the number of sides of the diamond cells) satisfying:

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad \rho_{\sigma}^* > 0, \quad \rho_{\sigma} > 0, \quad (3.7a)$$

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad \frac{|D_{\sigma}|}{\delta t} (\rho_{\sigma} - \rho_{\sigma}^*) + \sum_{\varepsilon \in \mathcal{E}(D_{\sigma})} F_{\varepsilon,\sigma} = 0 = 0, \quad (3.7b)$$

$$\forall \varepsilon = \sigma|\sigma', \quad F_{\varepsilon,\sigma} = -F_{\varepsilon,\sigma'}. \quad (3.7c)$$

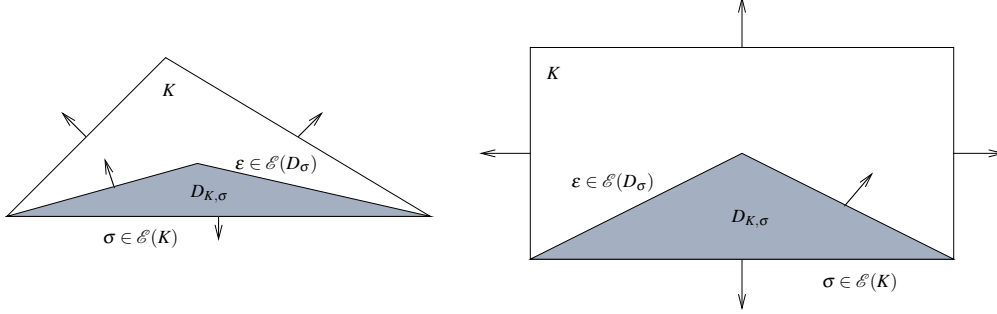


FIG. 1. Diamond-cells for the Crouzeix-Raviart and Rannacher-Turek element.

(note that this is the dual mesh equivalent of the primal mesh assumptions (2.8)). Let $(v_\sigma^*)_{\sigma \in \mathcal{E}_{\text{int}}}$ and $(v_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}$ be two families of real numbers. For any internal edge $\varepsilon = D_\sigma | D_{\sigma'}$, we define v_ε either by $v_\sigma = \frac{1}{2}(v_\sigma + v_{\sigma'})$ (centred choice), or by $v_\varepsilon = v_\sigma$ if $F_{\varepsilon,\sigma} \geq 0$ and $v_\varepsilon = v_{\sigma'}$ otherwise (upwind choice). In both cases, the following inequality holds:

$$\sum_{D_\sigma \in \mathcal{M}} v_\sigma \left[\frac{|D_\sigma|}{\delta t} (\rho_\sigma v_\sigma - \rho_\sigma^* v_\sigma^*) + \sum_{\varepsilon = D_\sigma | D_{\sigma'}} F_{\varepsilon,\sigma} v_\varepsilon \right] \geq \frac{1}{2} \sum_{D_\sigma \in \mathcal{M}} \frac{|D_\sigma|}{\delta t} [\rho_\sigma v_\sigma^2 - \rho_\sigma^* v_\sigma^{*2}]. \quad (3.8)$$

Let us then derive a discretization scheme for (3.2) for which we are able to apply Theorem 3.1, taking for v the velocity components. Let us first note that, for the Crouzeix-Raviart element and, for the rectangular (in two dimensions) or cubic (in three dimensions) Rannacher-Turek element, one has $\int_K \varphi_\sigma = |D_{K,\sigma}|$, where $|D_{K,\sigma}|$ is the measure of the half diamond cell $D_{K,\sigma}$. Thus, a mass lumping of the finite element discretization of the term $\rho^n u^{n+1}$ in Equation (3.2) associated with σ leads to an expression of the form $\rho_\sigma^n u_\sigma^{n+1}$, where ρ_σ^n is defined by:

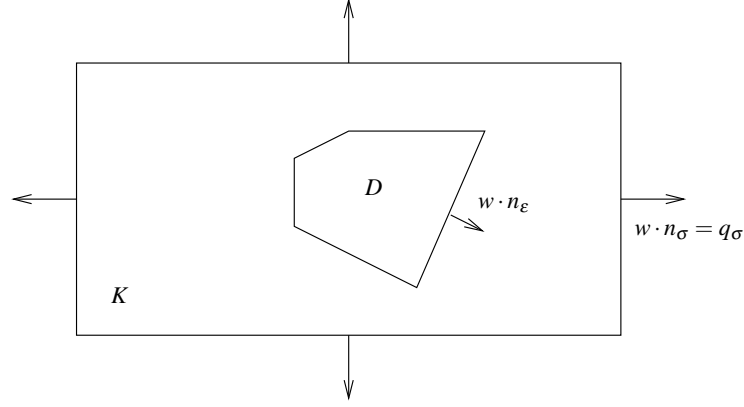
$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad |D_\sigma| \rho_\sigma^n = |D_{K,\sigma}| \rho_K^n + |D_{L,\sigma}| \rho_L^n. \quad (3.9)$$

We also note that when the diffusion term $\nabla \cdot \tau(u)$ reduces to the Laplace operator, its 2D Crouzeix-Raviart finite element discretization is identical to the finite volume discretization on the dual mesh consisting of the diamond cells D_σ [7, 6]. This property readily extends to the Rannacher-Turek element and suggests that the advection term $\nabla \cdot (\tilde{u}^{n+1} \otimes \rho^n u^n)$ in (3.2) be discretized on each edge $\sigma \in \mathcal{E}_{\text{int}}$ by the term $\sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma}^n \tilde{u}_\varepsilon^{n+1}$, where $\mathcal{E}(D_\sigma)$ is the set of the edges of D_σ , $\tilde{u}_\varepsilon^{n+1}$ is a centred approximation of \tilde{u}^{n+1} on ε and $F_{\varepsilon,\sigma}^n = |\varepsilon| q_\varepsilon^n \cdot n_\varepsilon$, where q_ε^n denotes an approximation of the momentum $\rho^n u^n$ on the edge ε , $|\varepsilon|$ is the measure of ε and n_ε is the normal to ε outward D_σ .

We then need to express $q_\varepsilon^n \cdot n_\varepsilon$ in such a way that the discrete mass balance (3.7) holds. To this goal, we use the following result [2] and give its (elementary) proof for the sake completeness.

LEMMA 3.1 (MASS BALANCE IN A SUB-VOLUME OF A MESH) Let $K \in \mathcal{T}$, let $(\rho^*, \rho) \in (\mathbb{R}_+)^2$, and consider a family $(F_{\sigma,K})_{\sigma \in \mathcal{E}(K)} \subset \mathbb{R}^k$ such that:

$$\frac{|K|}{\delta t} (\rho - \rho^*) + \sum_{\sigma \in \mathcal{E}(K)} F_{\sigma,K} = 0. \quad (3.10)$$

FIG. 2. sub-volume of K .

Let w be a momentum field on K , such that $\nabla \cdot w$ is constant over K and satisfies:

$$\int_K \nabla \cdot w = \int_{\partial K} w \cdot n_{\partial K} = \sum_{\sigma \in \mathcal{E}(K)} F_{\sigma, K}, \quad (3.11)$$

where ∂K and $n_{\partial K}$ stand for the boundary of K and the normal vector to ∂K outward to K , respectively. Let D be a subset of K with boundary ∂D (see Figure 2). Then the following property holds:

$$\frac{|D|}{\delta t} (\rho - \rho^*) + \int_{\partial D} w \cdot n_{\partial D} = 0,$$

where $n_{\partial D}$ stands for the normal vector to ∂D outward to D .

Proof. Multiplying (3.10) by $|D|/|K|$ and using (3.11) yields that:

$$\frac{|D|}{|\delta t|} (\rho - \rho^*) + \frac{|D|}{|K|} \int_K \nabla \cdot w = 0,$$

which concludes the proof thanks to the fact that $\nabla \cdot w$ is constant over K . \square

Now let us apply this lemma to obtain the mass conservation property on the diamond mesh \mathcal{M} . At step n , the discrete mass balance for ρ^n is obtained from the solution of (3.3b) at step $n-1$; its (finite volume) discretization is (2.3) with $\rho = \rho^n$, $\rho^* = \rho^{n-1}$ and $u = u^n$:

$$\frac{K}{\delta t} (\rho_K^n - \rho_K^{n-1}) + \sum_{\sigma \in \mathcal{E}(K)} F_{\sigma, K}^{up}(\rho^n, u^n) = 0, \quad \forall K \in \mathcal{T}. \quad (3.12)$$

We construct the momenta on the edges $\varepsilon \in \mathcal{E}(D_\sigma)$ by building on each cell K a field w^n with constant divergence and such that:

$$\forall \sigma \in \mathcal{E}(K), \quad \int_\sigma w^n \cdot n_\sigma = F_{\sigma, K}^{up}(\rho^n, u^n). \quad (3.13)$$

Assuming the field w^n to be constructed, the value $F_{\varepsilon,\sigma}^n$ of the mass flux on the edge or face ε of the diamond cell is then computed by integrating $w^n \cdot n_\varepsilon$ over ε , *i.e.*:

$$F_{\varepsilon,\sigma}^n = \int_{\varepsilon} w^n \cdot n_\varepsilon, \quad (3.14)$$

which yields a discrete mass balance over both half-diamond cells; thanks to the fact that ρ_σ^n is defined by (3.9), summing the two discrete balance equations on the half diamonds $D_{K,\sigma}$ and $D_{L,\sigma}$ gives the discrete mass balance over $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$, that is:

$$\frac{|D_\sigma|}{\delta t} (\rho_\sigma^n - \rho_\sigma^{n-1}) + \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma}^n = 0;$$

Condition (3.7b) of Theorem 3.1 is therefore satisfied.

Let us then turn to the construction of the field w^n ; such a field is derived for the Crouzeix-Raviart element by direct interpolation of quantities $((\rho u)_\sigma^n)_{\sigma \in \mathcal{E}(K)}$ (*i.e.* using the standard expansion of the Crouzeix-Raviart elements):

$$w^n(x) = \sum_{\sigma \in \mathcal{E}(K)} \varphi_\sigma(x) (\rho u)_\sigma^n$$

where $(\rho u)_\sigma^n$ is such that $\int_\sigma (\rho u)_\sigma^n \cdot n_\sigma = F_{\sigma,K}^{up}(\rho^n, u^n)$. With the chosen discretization for the mass balance (see (2.4)), the natural choice for $(\rho u)_\sigma^n$ reads $(\rho u)_\sigma^n = \rho_{up,\sigma}^n u_\sigma^n$ where $\rho_{up,\sigma}^n$ is the upwind density on σ with respect to u^n . Indeed, thanks to the fact that the Crouzeix Raviart basis functions satisfy $\int_\sigma \varphi_\sigma = |\sigma|$, one has $\int_\sigma w^n \cdot n_\sigma = F_{\sigma,K}^{up}(\rho^n, u^n)$.

For the Rannacher-Turek element, the divergence of discrete functions is no longer constant, but when the mesh is rectangular or cubic, we may use the following interpolation formula:

$$w^n(x) = \sum_{\sigma \in \mathcal{E}(K)} \alpha_\sigma(x \cdot n_\sigma) F_{\sigma,K}^{up}(\rho^n, u^n) n_\sigma$$

where the α_σ are affine interpolation functions which are determined in such a way that the relations (3.13) hold. The extension to more general grids is underway.

Finally, the standard (Crouzeix-Raviart) finite element expansion is used to discretize the terms $\nabla p^n - \nabla \cdot \tau(u^{n+1})$ of (3.2), and we thus obtain the following discrete momentum balance equation:

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d,$$

$$\begin{aligned} \frac{|D_\sigma|}{\delta t} (\rho_\sigma^n \tilde{u}_{\sigma,i}^{n+1} - \rho_\sigma^{n-1} u_{\sigma,i}^n) + \sum_{\substack{\varepsilon \in \mathcal{E}(D_\sigma), \\ \varepsilon = D_\sigma | D_{\sigma'} \\ \varepsilon = D_{\sigma'} | D_\sigma}} \frac{1}{2} F_{\varepsilon,\sigma}^n (\tilde{u}_{\sigma,i}^{n+1} + \tilde{u}_{\sigma',i}^{n+1}) \\ - \int_{\Omega, \mathcal{T}} p^n \nabla \cdot \varphi_\sigma^{(i)} + a_d(\tilde{u}^{n+1}, \varphi_\sigma^{(i)}) = \int_{\Omega} f^{n+1} \cdot \varphi_\sigma^{(i)} \end{aligned} \quad (3.15)$$

where

$$a_d(v, w) = \begin{cases} \mu \int_{\Omega, \mathcal{T}} \left[\nabla v : \nabla w + \frac{1}{3} \nabla \cdot v \nabla \cdot w \right] & \text{if (1.3) holds (case of constant viscosity),} \\ \int_{\Omega, \mathcal{T}} \tau(v) : \nabla w & \text{with } \tau \text{ given by (1.2) otherwise.} \end{cases}$$

3.3 Spatial discretization of the projection step

The discretization of the first projection step (3.3a) is consistent with that of the momentum equation (3.2); a mass lumping is performed for the unsteady term and a standard finite element formulation is used for the gradient of the pressure increment, yielding:

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d, \quad \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (u_{\sigma,i}^{n+1} - \tilde{u}_{\sigma,i}^{n+1}) - \int_{\Omega, \mathcal{T}} (p^{n+1} - p^n) \nabla \cdot \varphi_\sigma^{(i)} dx = 0,$$

where ρ_σ^n is defined by (3.9). Since the pressure is piecewise constant, the discrete gradient operator takes the form of the transposed of the standard finite volume discretization of the divergence (based on the primal mesh \mathcal{T} , and not on the diamond mesh \mathcal{M}) and can be rewritten as follows:

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \quad \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (u_\sigma^{n+1} - \tilde{u}_\sigma^{n+1}) + |\sigma| [(p_L^{n+1} - p_L^n) - (p_K^{n+1} - p_K^n)] n_{KL} = 0. \quad (3.16)$$

The mass balance equation (3.3a) is discretized as described in Section 2, which ensures that the density stays positive at all times; taking $\rho = \rho^{n+1}$, $\rho^* = \rho^n$ and $u = u^{n+1}$ in the finite volume scheme (2.3) yields:

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} [\rho_K^{n+1} - \rho_K^n] + \sum_{\sigma=K|L} F_{\sigma,K}^{up}(\rho^{n+1}, u^{n+1}) = 0, \quad (3.17a)$$

$$\rho_K^{n+1} = \eta(p_K^{n+1}, y_K^{n+1}). \quad (3.17b)$$

The projection step therefore consists of finding $(u_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$ and $(p_K^{n+1})_{K \in \mathcal{T}}$ that solve of the nonlinear system (3.16)–(3.17). Under some assumptions on the function η , we now prove that this projection step admits one solution.

LEMMA 3.2 Let us suppose that the equation of state η is such that for any $y \in [0, 1]$, the function $p \mapsto \eta(p, y)$ is defined and increasing on $[0, +\infty)$, $\eta(0, y) = 0$ and $\lim_{p \rightarrow +\infty} \eta(p, y) = +\infty$. Then the nonlinear algebraic system (3.16)–(3.17) admits at least one solution.

Proof. The proof of this lemma consists of an application of the Brouwer fixed point theorem (see *e.g.* [10, chapter 5]). Let $\tilde{u}^{n+1} = (\tilde{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}} \in \mathbb{R}^{M \times d}$, $y^{n+1} = (y_K^{n+1})_{K \in \mathcal{T}} \in \mathbb{R}^N$, $\rho^n = (\rho_K^n)_{K \in \mathcal{T}} \in \mathbb{R}^N$ and $p^n = (p_K^n)_{K \in \mathcal{T}} \in \mathbb{R}^N$ be known families of velocities, mass fractions, densities and pressures (calculated at the previous iteration or at the prediction step), such that $\rho_K^n > 0$ for all $K \in \mathcal{T}$. Let H be the mapping from $V_{\mathcal{T}} \times X_{\mathcal{T}}$ into itself defined by $(u, p) = ((u_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}, (p_K)_{K \in \mathcal{T}}) \mapsto H(u, p) = (v, q) = ((v_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}, (q_K)_{K \in \mathcal{T}})$, with (v, q) satisfying:

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} [\rho_K - \rho_K^n] + \sum_{\sigma=K|L} F_{\sigma,K}^{up}(\rho, u) = 0, \quad (3.18a)$$

$$\forall K \in \mathcal{T}, \quad \rho_K = \eta(q_K, y_K^{n+1}), \quad (3.18b)$$

$$\forall \sigma = K|L \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d, \quad \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (v_{\sigma,i} - \tilde{u}_{\sigma,i}^{n+1}) + \int_{\Omega, \mathcal{T}} (q - p^n) \nabla \cdot \varphi_\sigma^{(i)} dx = 0, \quad (3.18c)$$

where $F_{\sigma,K}^{up}$ is defined by (2.4) and ρ_σ^n by (3.9); note that $\rho_\sigma^n > 0$ for any $\sigma \in \mathcal{E}_{\text{int}}$. Equation (3.18a) is linear in ρ , and by Lemma 2.1, it has a unique solution $(\rho_K)_{K \in \mathcal{T}}$ satisfying $\rho_K > 0$ for all $K \in \mathcal{T}$.

Thanks to the assumption on η , for each $K \in \mathcal{T}$, there exists a unique q_K such that (3.18b) is satisfied. Finally, when q is computed, Equation (3.18c) yields a diagonal system for v . Hence the system (3.18) has a unique solution, so that the function H is well defined. We then note that any fixed point of the function H is a solution to the system (3.16)-(3.17). By conservativity of the finite volume scheme (3.17a), we easily see that $\sum_{K \in \mathcal{T}} |K| \rho_K = \sum_{K \in \mathcal{T}} |K| \rho_K^n$, and therefore, there exists $c_\rho \in \mathbb{R}$ such that $0 < \max_{K \in \mathcal{T}} \rho_K < c_\rho$ and again from the assumption on η , also for the pressure: $\max_{K \in \mathcal{T}} q_K \leq c_p$. Multiplying (3.18c) by $v_{\sigma,i}$ and summing over $\sigma \in \mathcal{E}_{\text{int}}$ and $i = 1, \dots, d$, we get that:

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma^n |v_\sigma|^2 \leq \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma^n v_\sigma \tilde{u}_\sigma^{n+1} + \delta t \int_{\Omega, \mathcal{T}} (q - p^n) \nabla \cdot v.$$

Using the positivity of ρ_σ^n and the fact that all norms are equivalent on a finite dimensional space, we get that there exists $c_u \geq 0$ such that $\|v\|_2^2 = \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| |v_\sigma|^2 \leq c_u^2$. By Brouwer's theorem, the mapping H therefore admits a fixed point in the convex set $\mathcal{C} = \{(v, q) \in X_{\mathcal{T}} \times V_{\mathcal{T}} \text{ such that } \|v\|_2 \leq c_u \text{ and } \max_{K \in \mathcal{T}} q_K \leq c_p\}$. \square

REMARK 3.2 We know from Remark 2.1 that, if $y^n > 0$ ($y^n < 1$), then $y^{n+1} > 0$ (resp. $y^{n+1} < 1$), and thus, by induction, pure monophasic zones do not appear if not already present at initial time. Hence, if no monophasic zone exists at initial time, it is sufficient for the existence of the solution that assumptions on the equation of state of Lemma 3.2 hold for $y \in (0, 1)$; this allows one to deal with cases where one phase is supposed to be incompressible.

Let us now combine the two algebraic relations (3.16) and (3.17) to build a discrete elliptic problem for the pressure. To this end, let us introduce the algebraic formulation of this system:

$$\frac{1}{\delta t} M_{\rho^n} (u^{n+1} - \tilde{u}^{n+1}) + B^t (p^{n+1} - p^n) = 0, \quad (3.19a)$$

$$\frac{1}{\delta t} R (\eta(p^{n+1}, y^{n+1}) - \rho^n) - B Q_{\rho^{n+1}, u^{n+1}}^{\text{up}} u^{n+1} = 0. \quad (3.19b)$$

where:

- M_{ρ^n} stands for the diagonal mass matrix weighted by the density ρ_σ^n at t^n (at edges or faces center) (*i.e.* the d diagonal entries of M_{ρ^n} associated to the edge σ are given by $|D_\sigma| \rho_\sigma^n$);
- B^t is the $((dM) \times N)$ - gradient operator matrix (recall that $M = \text{card}(\mathcal{E}_{\text{int}})$) and $N = \text{card}(\mathcal{T})$ and B is the opposite of the divergence operator matrix.
- $Q_{\rho^{n+1}, u^{n+1}}^{\text{up}}$ is a diagonal matrix; its entry corresponding to an edge $\sigma = K|L \in \mathcal{E}_{\text{int}}$ is $\rho_{\text{up}, \sigma}^{n+1}$, *i.e.* the upwind density with respect to u^{n+1} (see (2.5)). The matrix R is diagonal and, for any $K \in \mathcal{T}$, $R_{K,K} = |K|$.

The elliptic problem for the pressure is obtained by multiplying (3.19a) by $B Q_{\rho^{n+1}, u^{n+1}}^{\text{up}} (M_{\rho^n})^{-1}$ and using (3.19b). The resulting equation reads:

$$L_{\rho^{n+1}, u^{n+1}}^{\text{up}} p^{n+1} + \frac{1}{\delta t^2} R (\eta(p^{n+1}, y^{n+1})) = L_{\rho^{n+1}, u^{n+1}}^{\text{up}} p^n + \frac{1}{\delta t^2} R p^n + \frac{1}{\delta t} B Q_{\rho^{n+1}, u^{n+1}}^{\text{up}} \tilde{u}^{n+1}, \quad (3.20)$$

where $L_{\rho^{n+1}, u^{n+1}}^{\text{up}} = B Q_{\rho^{n+1}, u^{n+1}}^{\text{up}} (M_{\rho^n})^{-1} B^t$ can be evaluated through a finite volume type expression of the Laplace operator [14]:

$$(L_{\rho^{n+1}, u^{n+1}}^{\text{up}} p^{n+1})_K = \sum_{\sigma=K|L} \frac{\rho_{\text{up}, \sigma}^{n+1}}{\rho_{\sigma}^n} \frac{|\sigma|^2}{|D_{\sigma}|} (p_K^{n+1} - p_L^{n+1}), \forall K \in \mathcal{T}.$$

Note that, even if $\forall \sigma \in \mathcal{E}_{\text{int}}, \rho_{\text{up}, \sigma}^{n+1} = \rho_{\sigma}^n$, the operator $L_{\rho^{n+1}, u^{n+1}}^{\text{up}}$ differs from the usual finite volume Laplace operator (in the case of rectangles or cubes, by a factor $d = 2$ or 3 respectively [14]); this is linked to the fact that the nonconforming $(V_{\mathcal{T}}, X_{\mathcal{T}})$ finite element approximation is not consistent for the mixed form approximation of the Laplace equation, neither in the Crouzeix-Raviart case nor if in the Rannacher-Turek case. Provided that p^{n+1} is known, (3.19a) gives the updated value of the velocity:

$$u^{n+1} = \tilde{u}^{n+1} - \delta t (M_{\rho^n})^{-1} B^t (p^{n+1} - p^n). \quad (3.21)$$

Note that contrary to usual projection methods, equations (3.20) and (3.21) are not decoupled, because of the upwinding of the density with respect to u^{n+1} . We thus implement an iterative algorithm, which reads as follows when the equation of state is linear with respect to the pressure:

Initialization: $p_0^{n+1} = p^n$ and $u_0^{n+1} = \tilde{u}^{n+1}$.

Step 4.1 – Solve for p_{k+1}^{n+1} :

$$L_{\rho_k^{n+1}, u_k^{n+1}}^{\text{up}} p_{k+1}^{n+1} + \frac{1}{\delta t^2} R p(p_{k+1}^{n+1}, y^{n+1}) = L p^n + \frac{1}{\delta t^2} R p^n + \frac{1}{\delta t} B Q_{\rho^{n+1}, u_k^{n+1}}^{\text{up}} \tilde{u}^{n+1},$$

with $\rho_k^{n+1} = \eta(p_k^{n+1}, y^{n+1})$,

Step 4.2 – Compute u_{k+1}^{n+1} as :

$$u_{k+1}^{n+1} = \tilde{u}^{n+1} - \delta t (M_{\rho^n})^{-1} B^t (p_{k+1}^{n+1} - p^n).$$

Convergence criterion: $\max [\|p_{k+1}^{n+1} - p_k^{n+1}\|, \|u_{k+1}^{n+1} - u_k^{n+1}\|] < \varepsilon$.

When the equation of state is nonlinear with respect to p , which is in general the case, step 4.1 is replaced by one iteration of a quasi-Newton algorithm where only the diagonal term $\rho(p_{k+1}^{n+1}, y^{n+1})$ is differentiated with respect to p_{k+1}^{n+1} .

3.4 The fully discrete algorithm

To sum up, we consider the following algorithm:

Initialization – Let $(y_K^0)_{K \in \mathcal{T}} \in [0, 1]$, $(u_K^0)_{K \in \mathcal{T}} \in \mathbb{R}^d$ and $(\rho_K^{-1})_{K \in \mathcal{T}} \subset \mathbb{R}_+$, compute $(\rho_K^0)_{K \in \mathcal{T}}$ that solves:

$$\forall K \in \mathcal{T}, \quad \frac{|K|}{\delta t} [\rho_K^0 - \rho_K^{-1}] + \sum_{\sigma=K|L} F_{\sigma, K}^{up}(\rho^0, u^0) = 0, \quad (3.22)$$

with $F_{\sigma, K}^{up}$ defined by (2.4), and let $(p_K^0)_{K \in \mathcal{T}}$ be given by $\rho_K^0 = \eta(p_K^0, y_K^0)$.

Then, for $n = 0, 1, 2, \dots$:

1. **Computation of y** – Compute $(y_K^{n+1})_{K \in \mathcal{T}}$ by an upwind finite volume discretization, as explained in Section 2:

$$\begin{aligned} \forall K \in \mathcal{T}, \quad & \frac{|K|}{\delta t} (\rho_K^n y_K^{n+1} - \rho_K^{n-1} y_K^n) + \sum_{\sigma=K|L} F_{\sigma,K} y_\sigma^{n+1} \\ & + \sum_{\sigma=K|L} G_{\sigma,K} \Phi_\sigma(y_K^{n+1}, y_L^{n+1}) + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (y_K^{n+1} - y_L^{n+1}) = 0. \end{aligned} \quad (3.23)$$

2. **Prediction of the velocity** – Compute $(\tilde{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$ by equation (3.15), with $(\rho_\sigma^n)_{\sigma \in \mathcal{E}_{\text{int}}}$ defined by (3.9) and $F_{\epsilon,\sigma}^n$ defined by (3.14).
3. **Projection step** – Compute $(u_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$ and $(p_K^{n+1})_{K \in \mathcal{T}}$ from equations (3.16) and (3.17).

The following theorem gathers the properties of the proposed numerical scheme.

THEOREM 3.2 (PROPERTIES OF THE NUMERICAL SCHEME) Under the assumptions for the equation of state of Lemma 3.2, there exists a set of families $(u^n)_{n \geq 0}$, $(p^n)_{n \geq 0}$, $(\rho^n)_{n \geq -1}$ and $(y^n)_{n \geq 1}$ given by the proposed algorithm. Moreover, the following properties hold, for all $n \leq N$:

- (i) positivity of the density:

$$\rho_K^n > 0, \quad \forall K \in \mathcal{T};$$

- (ii) L^∞ stability property:

$$y_K^n \in [0, 1], \quad \forall K \in \mathcal{T};$$

- (iv) conservativity property:

$$\sum_{K \in \mathcal{T}} |K| \rho_K^{n-1} y_K^n = M_y^0 \text{ and } \sum_{K \in \mathcal{T}} |K| \rho_K^n = M^0, \quad \forall n \geq 1,$$

where M_y^0 (resp. M^0) denotes the initial gas mass (resp. total mass) in Ω .

Proof. The existence, uniqueness and positivity of ρ^0 satisfying (3.22) follows from Lemma 2.1, which also gives that $\rho_K^0 > 0$ for any $K \in \mathcal{T}$. Thanks to the assumptions on η and since $\rho^0 > 0$ and $y^0 \in [0, 1]$, ρ^0 is uniquely determined by the equation of state $\rho_K^0 = \eta(p_K^0, y_K^0)$.

Then, we get from Theorem 2.2 that there exists a unique family $(y_K^1)_{K \in \mathcal{T}}$ satisfying (3.23), and that $y_K^1 \in [0, 1]$. Summing Equation (3.23) over $K \in \mathcal{T}$ and using the conservativity of the fluxes (2.6) and (2.10) yields that $\sum_{K \in \mathcal{T}} |K| \rho_K^0 y_K^1 = M_y^0$.

The predicted velocity \tilde{u}^1 satisfies the set of linear equations (3.15). Multiplying (3.15) by \tilde{u}_σ^1 , summing over the edges $\sigma \in \mathcal{E}_{\text{int}}$ and using Theorem 3.1 (which we may use thanks to the careful discretization of the convection term) yields that there exists $C \in \mathbb{R}_+$ such that $\|\tilde{u}^1\| \leq C_1 \|f\|_{L^2(\Omega)}$, which in turn yields the existence and uniqueness of \tilde{u}^1 .

We then obtain from Lemma 3.2 that there exists u^1 , p^1 and ρ^1 satisfying (3.16)-(3.17), and from Lemma 2.1, that $\rho^1 > 0$. Summing Equation (3.17a) over $K \in \mathcal{T}$ and using the conservativity of the fluxes (2.6) yields that $\sum_{K \in \mathcal{T}} |K| \rho_K^1 = M^0$.

The proof of theorem is then completed by an easy induction. \square

4. Numerical results

In this section, we present two numerical tests to assess the behaviour of the fractional step scheme described above. In these tests, we compute the flow of an isothermal two-phase mixture of immiscible liquid and gas; the model (1.1) is then the so-called drift-flux model, that is a mixture model which takes into account the relative velocity u_r between the liquid and the gas phase (the so-called drift velocity), for which a phenomenological relation must be supplied. Then, ρ stands for the mixture density and takes the general form $\rho = (1 - \alpha_g)\rho_\ell + \alpha_g\rho_g(p)$ where α_g stands for the void fraction and $\rho_g(p)$ yields the gas density as a function of the pressure. Introducing the mass gas fraction y and using the relation $\alpha_g\rho_g = \rho y$ leads to the equation of state:

$$\eta(p, y) = \frac{1}{y/\rho_g(p) + (1-y)/\rho_\ell}. \quad (4.1)$$

In the perfect gas approximation and for a constant temperature, $\rho_g(p)$ is simply proportional to the pressure:

$$\rho_g(p) = \frac{p}{RT}, \quad (4.2)$$

where R is the gas constant and T is the absolute temperature. We further assume that the liquid phase can be considered as incompressible so that the liquid density ρ_ℓ is constant. In this section, the nonlinear function ϕ is given by $\phi(y) = y(1-y)$ and the flux function by $g(a, b) = \underline{a} - \underline{b}^2$ with $\underline{x} = \max(0, \min(x, 1))$ for any $x \in \mathbb{R}$.

4.1 Assessing the convergence of the scheme against an analytic solution

We first assess the convergence rate of the proposed scheme with respect to space and time discretizations through a model problem whose analytical solution is known.

We choose for the computational domain $\Omega = (0, 1) \times (-0.5, 0.5)$, and for the momentum and density the following expressions:

$$\rho(x, t) \quad u(x, t) = -\frac{1}{4} \cos(\pi t) \begin{bmatrix} \sin(\pi x_1) \\ \cos(\pi x_2) \end{bmatrix} \quad \rho(x, t) = 1 + \frac{1}{4} \sin(\pi t) [\cos(\pi x_1) - \sin(\pi x_2)]$$

The pressure and the gas mass fraction are linked to the density by the equation of state (4.1), where the liquid density ρ_ℓ is set at $\rho_\ell = 5$ and the product RT in the equation of state of the gas (4.2) is given by $RT = 1$ (so $\rho_g = p$). We choose the following expression for the unknown y :

$$y(x, t) = \frac{2.5 - 0.5 \rho(x, t)}{4.5 \rho(x, t)}$$

The relative velocity is constant and is given by $u_r = (0, 1)^t$ and the diffusive coefficient D is equal to 0.1. The analytical expression for the pressure is obtained from the equation of state. These functions satisfy the mass balance equation (1.1a); for the gas mass fraction equation (1.1c) and momentum balance equation (1.1b), we add the corresponding right-hand side. In this latter equation, we also assume that the divergence of the stress tensor is given by (1.3) with $\mu = 10^{-2}$ and we use the corresponding expression for the bilinear form $a_d(\cdot, \cdot)$.

For the Rannacher-Turek element, computations are made with 20×20 , 40×40 and 80×80 uniform meshes. For the Crouzeix-Raviart element, the meshes are built as follows: the computational domain

is first split in square subdomains, then each subdomain is split in 26 simplices, all having angles of at most 80° , according to the pattern given in [3, Figure 5 – bbbb], see Figure 3. The first splitting of the domain yields 20×20 , 40×40 and 80×80 uniform grids.

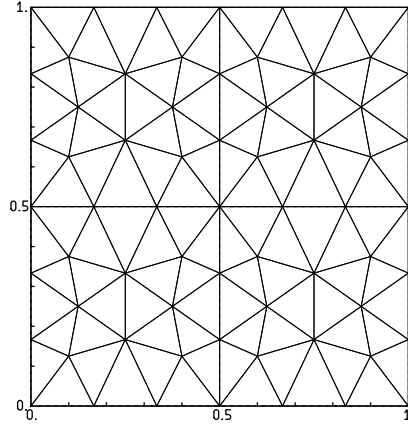


FIG. 3. Pattern used to split square cells into acute-angle triangles, for the Crouzeix-Raviart element.

Velocity, pressure and gas mass fraction errors obtained at $t = 0.5$ as a function of the time step are drawn on Figure 4, Figure 5 and Figure 6, respectively. These errors are evaluated in the L^2 norm for the velocity and in the discrete L^2 norms for the pressure and the gas mass fraction. For large time steps, these curves show a decrease corresponding to approximately a first order convergence in time, until a plateau is reached, due to the fact that errors are bounded by below by the residual spatial discretization error. The value of the errors on this plateau then shows a spatial convergence order close to one, which is consistent with the choice of an upwind discretization for the advection terms in the gas mass fraction and mass balance equations. Finally, the numerical results seem to be significantly more accurate with the Rannacher-Turek element.

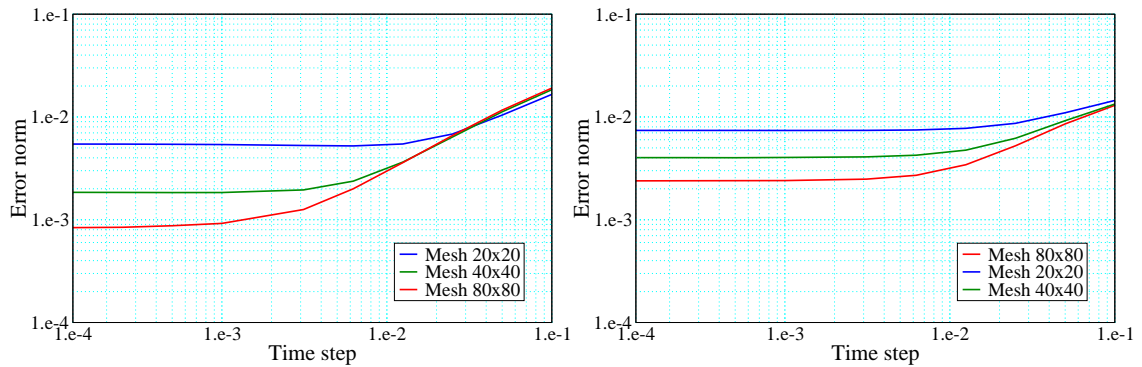


FIG. 4. Velocity error as a function of the time step. Left: Rannacher-Turek element, right: Crouzeix-Raviart element

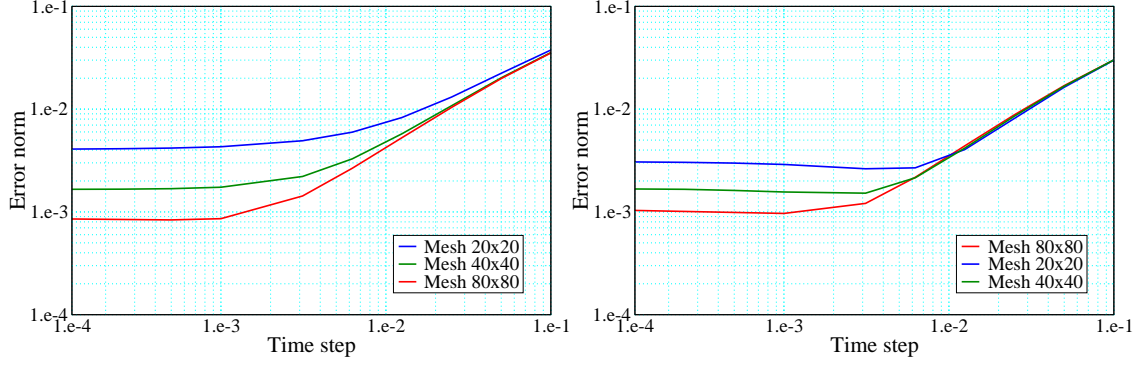


FIG. 5. Pressure error as a function of the time step. left: Rannacher-Turek element, right: Crouzeix-Raviart element

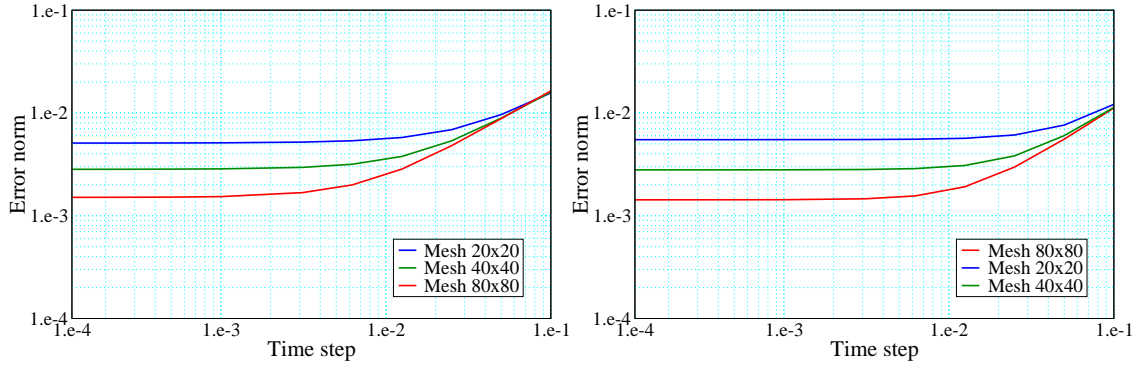


FIG. 6. Gas mass fraction error as a function of the time step. left: Rannacher-Turek element, right: Crouzeix-Raviart element

4.2 A phase separation problem

We now present numerical results obtained for a phase separation problem, with data inspired by a classical benchmark test for the simulation of two-phase flows [8, 15, 21] with two-fields models (*i.e.* models considering separate balance equations for each phase). The physical domain considered is a vertical tube of length $L = 7.5m$, filled at initial time with a two-phase mixture of air and water with $\alpha = 0.5$, $u = 0$ and $p = p_0$ where $p_0 = 10^5 Pa$ is the ambient pressure. Under the action of gravity (with $g = 9.81 m.s^{-2}$), phases separate and the solution at $t = +\infty$ is the superposition of a zone of pure water and a zone of pure air, both at rest. In the original problem, the interactions between the two phases are neglected; instead, we assume here that the relative velocity is constant and given by $u_r = 1 m.s^{-1}$, which is clearly nonphysical (at small times, water droplets just fall with a constant acceleration g).

However, even under this assumption, the solution of the problem qualitatively reproduces the original phase separation phenomenon.

The equation of state for the mixture is the same as in the previous test case and the densities, for water and air respectively, are $\rho_\ell = 1000 \text{ kg.m}^{-3}$ and $\rho_g = p/RT$ where RT is such that $\rho_g = 1.2 \text{ kg.m}^{-3}$ at $p = 10^5 \text{ Pa}$. The diffusion coefficient D and the viscosity μ are set to zero. At the top and the bottom boundaries, both the velocity and the relative velocity are prescribed to be zero.

For this test case, we use a regular mesh composed of rectangular cells (with the Rannacher-Turek element); since this problem is one-dimensional, only one cell is used in the horizontal direction, and 200 in the vertical direction. Calculations with time steps up to $\delta t = 10^{-1} \text{ s}$ have been performed without observing any instability. With respect to the time discretization, the convergence for the void fraction and the density is readily achieved, and profiles obtained with $\delta t \leq 10^{-2} \text{ s}$ are all similar; the results with $\delta t = 10^{-2} \text{ s}$ are reported on Figure 7. However, probably because of the large variations of the density near $y = 1$, themselves due to the large difference of the densities of the two phases, convergence for y is more difficult to reach, and variations in the profiles obtained are observed when decreasing the time step down to $\delta t = 5 \cdot 10^{-4} \text{ s}$.

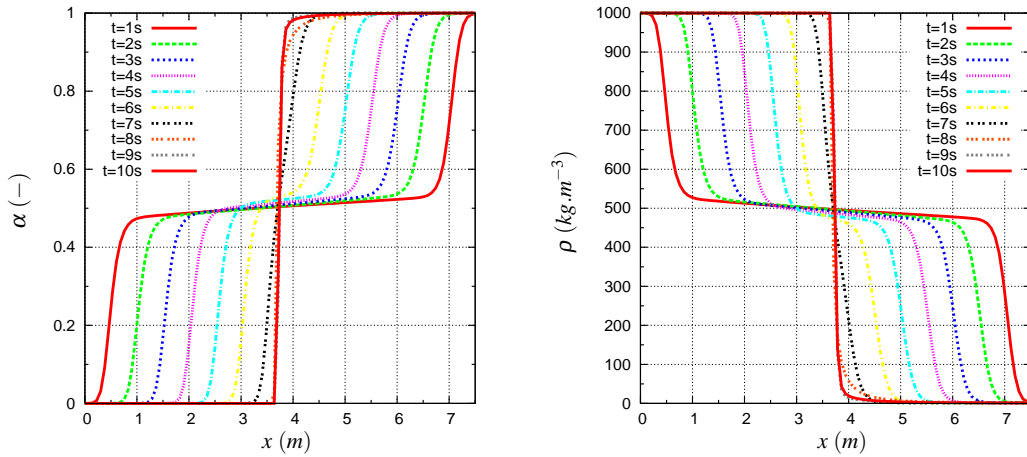


FIG. 7. void fraction and density profiles for the phase separation problem.

5. Conclusion

In this paper, we first addressed a parabolic equation that models the phase mass balance in two-phase flows. It differs from the mass balance for chemical species in compressible multi-component flows studied in [18] by the addition of a nonlinear term of the form $\nabla \cdot \rho \phi(y) u_r$, where y is the unknown, $\phi(\cdot)$ is a regular function such that $\phi(0) = \phi(1) = 0$ and u_r is a general (in particular not necessarily divergence free) velocity field. We proved the existence and uniqueness of the finite volume approximation together with the fact that it remains within physical bounds. As in [18], the necessary condition for this L^∞ stability result is that the discretization of the convection operator is such that it vanishes

for constant y , which amounts to demanding that some discrete mass balance equation be satisfied. The second ingredient of this scheme is a discretization of the nonlinear term based on the notion of monotone flux functions [13]. This work extends the theory developed in [18] in two directions: it copes with a new nonlinear term, and introduces different techniques which appear to be well-suited for nonlinear problems: proof of an L^∞ *a priori* estimate for the solution and proof of its existence (by a topological degree argument) and its uniqueness.

In a second part of the paper, we proposed a discretization by a fractional step method for the set of equations (1.1). This algorithm decouples the resolution of the phase mass balance from the resolution of the Navier-Stokes equations and meets two essential requirements: it is conservative, and the discrete mass balance needed for a stable computation of y is satisfied. To achieve this goal, the key ingredient is a particular time-discretization of the density terms, which unfortunately limits the time accuracy of the scheme to first order. This technique is now routinely used in the ISIS computer code [1] developed at IRSN for the modelling of reacting flows; it demonstrates very satisfactory stability properties, even in cases of large density variations for which numerical difficulties are often reported in the literature. Let us also mention that the proposed numerical scheme degenerates to a classical projection method in the incompressible limit.

As far as extensions of this work are concerned, first, (formally) second order in space discretizations (typically using MUSCL-like techniques) should be developed. Second, the proposed algorithm is based on the simplest and computationally cheapest fractional step approach (decoupling all equations) so that it should be retained as far as it works. Unfortunately, in the specific case of two-phase compressible flows involving phases of very different densities, instabilities are observed, the cure to which seems to be a drastic time step reduction. These instabilities appear to be linked to the fact that the present algorithm does not preserve a constant pressure through moving interfaces between phases (*i.e.* contact discontinuities of the underlying hyperbolic system); a solution to this problem, still based on the same essential ingredients for the evaluation of the density terms and the discretization of the phase mass balance but coupling this latter equation to the projection step, is now under development and shows promising results.

Acknowledgment

The authors would like to thank the referees for their useful suggestions, which led to a significant improvement of the paper.

REFERENCES

- [1] F. Babik, T. Gallouët, J.-C. Latché, S. Suard, and D. Vola. On some fractional step schemes for combustion problems. In *Finite Volumes for Complex Applications IV (FVCA IV)*, pages 505–514. Editions Hermès, Paris, 2005.
- [2] F. Babik, J.-C. Latché, and D. Vola. An L^2 -stable approximation of the Navier-Stokes advective operator for non conforming finite elements. In *Mini-Workshop on Variational Multiscale Methods and Stabilized Finite Elements, Lausanne*, 2007.
- [3] M. Bern, D. Eppstein, and J. Gilbert. Provably good mesh generation. *Journal of Computer and System Sciences*, 48:384–409, 1994.
- [4] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, 1991.
- [5] C. Calgaro, E. Creusé, and T. Goudon. An hybrid finite-volume-finite element method for variable density incompressible flows. *Journal of Computational Physics*, 227:4671–4696, 2008.
- [6] P. Chatzipantelidis. A finite volume method based on the Crouzeix-Raviart element for elliptic PDE's in two

- dimensions. *Numer. Math.*, 82(3):409–432, 1999.
- [7] S. H. Chou. Analysis and convergence of a covolume method for the generalized Stokes problem. *Math. Comp.*, 66(217):85–104, 1997.
- [8] F. Coquel, K. El Amine, E. Godlewski, B. Perthame, and P. Rascle. A numerical method using upwind schemes for the resolution of two-phase flows. *Journal of Computational Physics*, 136:272–288, 1997.
- [9] M. Crouzeix and P.-A. Raviart. Conforming and nonconforming finite element methods for solving the stationary Stokes equations I. *Revue Française d'Automatique, Informatique et Recherche Opérationnelle (R.A.I.R.O.)*, R-3:33–75, 1973.
- [10] P. Drábek and J. Milota. *Methods of nonlinear analysis*. Birkhäuser Advanced Texts. 2007.
- [11] A. Ern. *Aide Mémoire Éléments finis*. Dunod, Paris, 2005.
- [12] R. Eymard, T. Gallouët, M. Ghilani, and R. Herbin. Error estimates for the approximate solutions of a nonlinear hyperbolic equation given by finite volume schemes. *IMA Journal of Numerical Analysis*, 18(4):563–594, 1998.
- [13] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. In P. Ciarlet and J.L. Lions, editors, *Handbook of Numerical Analysis, Volume VII*, pages 713–1020. North Holland, 2000.
- [14] 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.
- [15] T. Gallouët, J.-M. Hérard, and N. Seguin. Numerical modeling of two-phase flows using the two-fluid two-pressure approach. *Mathematical Models and Methods in Applied Sciences*, 14(5):663–700, 2004.
- [16] D. Gilbarg and N. Trudinger. *Elliptic Partial Differential Equations of Second Order*. Classics in Mathematics. Springer, third edition, 2001.
- [17] 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.
- [18] B. Larrouturou. How to preserve the mass fractions positivity when computing compressible multi-component flows. *Journal of Computational Physics*, 95:59–84, 1991.
- [19] M. Marion and R. Temam. Navier-Stokes equations: Theory and approximation. In P. Ciarlet and J.L. Lions, editors, *Handbook of Numerical Analysis, Volume VI*. North Holland, 1998.
- [20] R. Rannacher and S. Turek. Simple nonconforming quadrilateral Stokes element. *Numerical Methods for Partial Differential Equations*, 8:97–111, 1992.
- [21] H. Staedtke, G. Franchello, B. Worth, U. Graf, P. Romstedt, A. Kumbaro, J. Garca-Cascales, H. Paillère, H. Deconinck, M. Ricchiuto, B. Smith, F. De Cachard, E.F. Toro, E. Romenski, and S. Mimouni. Advanced three-dimensional two-phase flow simulation tools for application to reactor safety (ASTAR). *Nuclear Engineering and Design*, 235(2-4):379–400, 2005.