

UNIVERSITÉ DE PROVENCE  
U.F.R. de Mathématiques, Informatique et Mécanique  
ÉCOLE DOCTORALE DE MATHÉMATIQUES ET INFORMATIQUE  
E.D. numéro 184

**THÈSE**

présentée pour obtenir le grade de  
DOCTEUR DE L'UNIVERSITÉ DE PROVENCE

*Spécialité : Mathématiques*

par

**Laura GASTALDO**

sous la direction du Pr. Raphaèle HERBIN

*Titre :*

**MÉTHODES DE CORRECTION DE PRESSION POUR LES  
ÉCOULEMENTS COMPRESSIBLES :  
APPLICATION AUX ÉQUATIONS DE NAVIER-STOKES BAROTROPES  
ET AU MODÈLE DE DÉRIVE**

soutenue publiquement le 29 novembre 2007

JURY

M. Franck BOYER	Université Paul Cézanne	<i>Examineur</i>
M. Miloslav FEISTAUER	Université Charles de Prague	<i>Rapporteur</i>
Mme Vivette GIRAULT	Université Paris VI	<i>Examinatrice</i>
M. Hervé GUILLARD	INRIA, Sophia Antipolis	<i>Rapporteur</i>
M. Philippe HELLUY	Université Louis Pasteur	<i>Examineur</i>
M. Jean-Marc HÉRARD	EDF, Chatou	<i>Examineur</i>
Mme Raphaèle HERBIN	Université de Provence	<i>Directeur de thèse</i>
M. Jean-Claude LATCHÉ	IRSN, Cadarache	<i>Encadrant</i>



# Table des matières

Le manuscrit est organisé de la manière suivante. Le premier chapitre est un résumé de l'ensemble des travaux effectués. Le deuxième chapitre est un article soumis à *Mathematical Modelling and Numerical Analysis* et détaille la construction ainsi que l'analyse d'un schéma de correction de pression pour la résolution des équations de Navier-Stokes compressibles et barotropes. Le troisième chapitre est un article soumis à *IMA Journal of Numerical Analysis* et décrit l'approximation par volumes finis de l'équation de transport de l'une des deux phases dans le modèle de dérive; elle est ensuite intégrée à un schéma entièrement découplé pour la résolution de ce modèle. Le dernier chapitre est un article soumis, dans une version raccourcie, à *Mathematical Modelling and Numerical Analysis* et décrit la construction et l'analyse, à partir des résultats fournis dans les deux chapitres précédents, d'un schéma de correction de pression pour la résolution du modèle de dérive.

<b>I Synthèse générale</b>	<b>9</b>
I.1 Introduction . . . . .	9
I.2 Contexte . . . . .	11
I.3 Modèle mathématique . . . . .	12
I.3.1 Equations locales instantanées . . . . .	13
I.3.2 Modèle diphasique à l'échelle macroscopique . . . . .	13
I.3.3 Modèle à vitesse de dérive ou modèle drift-flux . . . . .	14
I.4 Quelques résultats fondamentaux en volumes finis . . . . .	16
I.4.1 Discrétisation spatiale et temporelle . . . . .	16
I.4.2 Une condition de compatibilité centrale : la satisfaction du bilan de masse . . . . .	17
I.4.3 Stabilité de l'opérateur d'advection . . . . .	17
I.4.4 Stabilité induite par le travail des forces de pression . . . . .	18
I.4.5 Stabilité $L^\infty$ pour une équation d'advection-diffusion non-linéaire . . . . .	21
I.5 Schémas de correction de pression pour les écoulements compressibles . . . . .	23
I.5.1 Discrétisation spatiale . . . . .	24
I.5.2 Schéma de correction de pression pour les équations de Navier-Stokes barotropes . . . . .	25
I.5.3 Schéma de correction de pression pour le modèle à vitesse de dérive . . . . .	30
I.6 Résultats numériques . . . . .	34
I.6.1 Un cas de ballottement . . . . .	35
I.6.2 Colonne à bulles . . . . .	36
I.7 Conclusion . . . . .	37
<b>II An unconditionally stable pressure correction scheme for compressible barotropic Navier-Stokes equations</b>	<b>41</b>
II.1 Introduction . . . . .	41
II.2 Analysis of a class of discrete problems . . . . .	43

II.2.1	The discrete problem . . . . .	44
II.2.2	On the pressure control induced by the pressure forces work . . . . .	46
II.2.3	Existence of a solution . . . . .	48
II.2.4	Some cases of application . . . . .	53
II.3	A pressure correction scheme . . . . .	54
II.3.1	Time semi-discrete formulation . . . . .	55
II.3.2	Stability of the advection operator : a finite-volume result . . . . .	56
II.3.3	Space discretization of the density prediction and the momentum balance equation . . . . .	58
II.3.4	Space discretization of the projection step . . . . .	60
II.3.5	Renormalization steps . . . . .	62
II.3.6	An overview of the algorithm . . . . .	63
II.3.7	Stability analysis . . . . .	63
II.3.8	Implementation . . . . .	66
II.3.9	Numerical experiments . . . . .	68
II.4	Conclusion . . . . .	69
<b>III On a discretization of phases mass balance in segregated algorithms for the drift-flux model</b>		<b>71</b>
III.1	Introduction . . . . .	71
III.2	Discretization for the nonlinear advection-diffusion equation . . . . .	73
III.2.1	An $L^\infty$ stability property . . . . .	75
III.2.2	Existence for the approximate solution . . . . .	80
III.2.3	Uniqueness of the approximate solution . . . . .	81
III.3	A fractional step algorithm for dispersed two-phase flows . . . . .	84
III.3.1	The semi-discretized in time algorithm . . . . .	84
III.3.2	Spatial discretization of the momentum balance equation . . . . .	86
III.3.3	Spatial discretization of the projection step . . . . .	89
III.3.4	Spatial discretization for the nonlinear advection diffusion equation . . . . .	92
III.3.5	An overview of the algorithm . . . . .	93
III.4	Numerical results . . . . .	94
III.4.1	Assessing the convergence against an analytic solution . . . . .	94
III.4.2	A phase separation problem . . . . .	96
III.4.3	Flow of a sedimenting dilute suspension over a rectangular bump . . . . .	98
III.5	Conclusion . . . . .	100
<b>IV An entropy preserving finite-element/finite-volume pressure correction scheme for the drift-flux model</b>		<b>101</b>
IV.1	Introduction . . . . .	101
IV.2	The numerical algorithm . . . . .	104
IV.2.1	Time semi-discrete formulation . . . . .	104
IV.2.2	Spatial discretization . . . . .	105
IV.2.3	Spatial discretization of the momentum balance equation . . . . .	107
IV.2.4	Spatial discretization of the pressure correction step . . . . .	110
IV.2.5	Spatial discretization of the correction step for $y$ . . . . .	112
IV.2.6	Some properties of the scheme . . . . .	113
IV.3	The stability induced by the pressure forces work . . . . .	116
IV.3.1	Abstract estimates . . . . .	116
IV.3.2	The case of a constant density liquid and an ideal gas . . . . .	120

IV.4	Stability analysis . . . . .	121
IV.4.1	First case : $u_r = 0, D = 0$ . . . . .	122
IV.4.2	Dissipativity of the drift term . . . . .	125
IV.5	Numerical results . . . . .	129
IV.5.1	Assessing the convergence against an analytic solution . . . . .	129
IV.5.2	The glass of water problem . . . . .	130
IV.5.3	Transport of interfaces . . . . .	131
IV.5.4	Two-dimensional sloshing in cavity . . . . .	132
IV.5.5	Bubble column . . . . .	137
IV.6	Conclusion . . . . .	138
IV.7	Existence of a solution to a class of discrete diphasic problems . . . . .	140

<b>Bibliographie</b>	<b>151</b>
----------------------	------------



# Chapitre I

## Synthèse générale

### I.1 Introduction

L'Institut de Radioprotection et de Sûreté Nucléaire (IRSN) réalise des expertises et mène des recherches dans les domaines de la sûreté nucléaire, de la protection contre les rayonnements ionisants, du contrôle et de la protection des matières nucléaires et de la protection contre les actes de malveillance. Une part essentielle de l'analyse de sûreté est constituée par l'étude des différentes situations auxquelles un réacteur nucléaire peut se trouver confronté depuis les conditions normales de fonctionnement jusqu'aux accidents graves qui sont le cadre général de cette thèse.

Nous développons ici un outil de simulation pour les bains à bulles, tels que rencontrés dans certaines phases tardives des scénarios de fusion de cœur pour les réacteurs nucléaires à eau pressurisée, lorsque les matériaux fondus issus de la cuve viennent interagir avec le béton du radier. La modélisation est basée sur le modèle dit "à vitesse de dérive", constitué des équations de bilan de masse et de quantité de mouvement pour le mélange (équations de Navier-Stokes) et d'une équation de conservation de masse de la phase gazeuse.

La définition de méthodes numériques pour la résolution de ces équations aux dérivées partielles doit tenir compte de certains aspects du modèle. Le nombre de Mach caractéristique de l'écoulement est faible, avec potentiellement des zones incompressibles. Le schéma numérique devra ainsi répondre en premier lieu à l'impératif de conserver ses propriétés de stabilité et de convergence dans une large gamme de nombre de Mach, et notamment jusqu'à l'incompressible. Par ailleurs, la satisfaction du principe de maximum (garder la fraction massique de la phase gazeuse entre 0 et 1, par exemple) est indispensable.

Pour la construction d'un algorithme inconditionnellement stable, *i.e.* quelque soit le nombre de Mach, il est possible de s'inspirer de schémas initialement développés dans le contexte des écoulements incompressibles. Parmi ceux-ci, les méthodes de projection ont, depuis les travaux originels de Chorin et Témam (voir, par exemple, [14, 68] pour les papiers d'origine, [53] pour une introduction ou [38] pour une revue de la plupart des variantes), acquis une popularité croissante. Ce succès tient dans le fait qu'elles découplent à chaque pas de temps les équations de bilan de quantité de mouvement et de bilan de masse, substituant ainsi à un problème mixte, de résolution difficile, une succession de problèmes elliptiques plus aisés à résoudre. Le principe de ces méthodes est le suivant. Dans une première étape, on obtient une prédiction de la vitesse pour la résolution de l'équation de bilan de quantité de mouvement, dans laquelle la pression est ignorée (méthode originelle) ou approchée par une formule explicite (méthode dite incrémentale). La seconde étape s'apparente à un problème de Darcy, qui est classiquement réécrit comme un problème elliptique pour la pression (méthode originelle) ou l'incrément de pression (méthode incrémentale).

Le développement des méthodes de correction de pression pour les équations de Navier-Stokes

compressibles remonte aux années soixante avec le travail de Harlow et Amsden [41, 42], lesquels développèrent un algorithme itératif (nommé méthode ICE) comprenant une étape de correction elliptique pour la pression. Par la suite, des équations de correction de pression apparaissent dans les schémas numériques proposés par différents chercheurs, essentiellement dans le cadre des volumes finis, en utilisant soit une disposition des inconnues co-localisée [59, 20, 49, 60, 46, 54] soit décalée [11, 44, 45, 48, 7, 16, 69, 73, 74, 70, 72]; dans le premier cas, certaines actions correctives doivent être menées afin d'éviter le découplage pair-impair de la pression à faible nombre de Mach. Certains de ces algorithmes sont essentiellement implicites; la solution de fin de pas est alors obtenue avec un processus itératif de type SIMPLE [71, 48, 20, 49, 60, 46, 54]. Les autres schémas [44, 45, 59, 7, 16, 75, 69, 74, 70, 72] sont des méthodes de prédiction-corrrection, avec essentiellement deux étapes en séquence : une prédiction semi-explicite de la quantité de mouvement ou de la vitesse (et éventuellement de l'énergie, pour les écoulements anisothermes) suivie d'une étape de correction au cours de laquelle la pression de fin de pas est calculée et la quantité de mouvement ou la vitesse sont corrigées, comme dans les méthodes de projection pour les écoulements incompressibles. Le schéma Characteristic-Based Split (CBS) (voir [56] pour un article récent ou [76] pour le papier d'origine), développé dans le contexte des éléments finis, appartient également à cette dernière classe de méthodes.

Les extensions des algorithmes de correction de pression aux écoulements multiphasiques sont plus rares, et semblent se restreindre aux algorithmes itératifs, similaires à l'algorithme SIMPLE [65, 55, 50].

Dans ce mémoire, nous proposons, dans un premier temps, un schéma de correction de pression non-itératif pour la résolution des équations de Navier-Stokes compressibles et barotropes basé sur une approximation par éléments finis mixtes non conformes. Son développement a été guidé par le souci de respecter, au niveau discret, les deux égalités de stabilité suivantes :

$$\int_{\Omega} \frac{\partial \rho u}{\partial t} u + \int_{\Omega} \nabla \cdot (\rho u \otimes u) u = \frac{d}{dt} \int_{\Omega} \frac{1}{2} \rho u^2 \quad (\text{I.1})$$

$$- \int_{\Omega} p \nabla \cdot u = \frac{d}{dt} \int_{\Omega} \rho f(\rho) \quad (\text{I.2})$$

où  $u$  désigne la vitesse,  $\rho$  la masse volumique,  $p$  la pression et  $f$  l'énergie libre dans l'écoulement. La première relation est le théorème de l'énergie cinétique et assure la stabilité de l'opérateur d'advection appliqué à chacune des composantes de la vitesse, la deuxième fournit un contrôle, par le travail des forces de pression, de l'énergie potentielle  $\rho f(\rho)$ . Ces deux résultats nous permettent d'une part d'établir des estimations *a priori* pour la vitesse et la pression qui assurent l'existence d'une solution à chaque étape de l'algorithme et d'autre part de démontrer la stabilité du schéma au sens de la décroissance de l'entropie discrète.

L'extension de cet algorithme aux écoulements diphasiques comporte une autre difficulté : conserver la fraction massique de gaz dans ses bornes physiques, à savoir entre 0 et 1. L'équation de transport de la phase gazeuse est une équation d'advection-diffusion qui diffère du bilan de masse usuel pour les espèces chimiques dans les écoulements réactifs compressibles étudiés par Larrouturou [51] par l'ajout d'un terme non-linéaire  $\nabla \cdot \rho \varphi(y) u_r$ , où  $\varphi(\cdot)$  est une fonction régulière telle que  $\varphi(0) = \varphi(1) = 0$  (dans le cas présent,  $\varphi(y) = y(1 - y)$ ). Nous construisons pour cette équation de conservation une approximation par volumes finis ayant les propriétés souhaitées; ce résultat de stabilité garantit l'existence de la solution à ce problème, dont on démontre ensuite l'unicité par une technique de problème dual. Ce développement étend à l'équation spécifique du modèle de dérive les résultats théoriques connus dans le cas linéaire.

Enfin, sur la base de ces travaux, un schéma conservatif, monotone et stable indépendamment du nombre de Mach est construit pour le modèle de dérive. Le développement de cet algorithme



nécessite la transposition au cas diphasique de la relation (I.2), soit l'analogue discret de l'identité suivante :

$$-\int_{\Omega} p \nabla \cdot u = \frac{d}{dt} \int_{\Omega} z f(\rho, z)$$

où  $z$  est la masse volumique partielle du gaz, et  $f$  désigne, dans ce cas, l'énergie libre de la phase gazeuse qu'on suppose compressible. Comme dans le cas précédent, cette propriété jointe à la stabilité de l'opérateur d'advection et au fait que l'algorithme conserve les variables dans leurs bornes physiques, permet de démontrer l'existence d'une solution discrète ainsi que la stabilité du schéma numérique, soit la décroissance de l'entropie discrète. Ce modèle présente en outre une caractéristique importante pour garantir la robustesse du schéma : la capacité de transporter une interface entre phases à vitesse et pression constantes, si telle est la solution du problème continu. Cette dernière propriété, comme la stabilité du schéma, nécessite l'emploi d'une étape de correction de pression couplant les bilans de masse du mélange et de la phase gazeuse.

Après une rapide présentation du contexte de la thèse ainsi que du modèle physique utilisé, nous décrivons les propriétés de stabilité qui sont à la base des deux schémas développés au cours de cette thèse. Ensuite nous construisons les algorithmes dédiés à la résolution des équations de Navier-Stokes barotropes compressibles ainsi que du modèle de dérive à partir des résultats précédents. Enfin, nous présentons deux cas tests de la littérature qui démontrent le bon fonctionnement de ce dernier.

## I.2 Contexte

Comme dit précédemment, cette thèse s'inscrit dans un programme de recherche visant à la simulation des accidents graves dans un réacteur à eau pressurisée (REP). Nous décrivons ici brièvement le déroulement d'un tel accident.

L'évènement indicateur est la perte par une brèche de réfrigérant du circuit primaire, qui entraîne le réchauffement du cœur. Si l'on suppose alors la défaillance simultanée d'un certain nombre de dispositifs de sûreté, l'accident peut se poursuivre jusqu'à la fusion du cœur. Ce dernier se relocalise alors en fond de cuve, ce qui peut conduire à la rupture de cette dernière (figure I.1). Dans cette éventualité, un bain de corium (mélange de matériaux fondus du cœur et de la cuve) se forme dans le puits de cuve composé de béton. Le corium, encore chauffé par le dégagement de puissance résiduelle dû à la désintégration des produits de fission, interagit avec les structures en béton qui le contiennent, et le bain pénètre peu à peu dans le radier ainsi que dans les parois latérales. Cette interaction s'accompagne de relâchements importants de gaz : vaporisation de l'eau contenue dans le béton et formation de dioxyde de carbone par décomposition du calcaire, principalement. Le bain est alors traversé par un flux de bulles. Cette étape de l'accident est appelée interaction corium-béton. Le puits de cuve est une des barrières de confinement du corium, dont il est primordial de connaître le temps de percée. Ce dernier dépend essentiellement de la vitesse et la direction (horizontale ou verticale) de l'ablation du béton.

Le bain de corium contient des oxydes lourds en provenance du cœur, des oxydes légers en provenance du béton et des métaux, le tout étant soumis au brassage induit par les gaz de décomposition du béton. De plus la mise en contact du corium à haute température avec le béton plus froid peut entraîner sa solidification, des débris étant susceptibles d'être mis en suspension dans le liquide. Le bain de corium est donc un milieu multiphasique (liquide, solide, gaz) dont la composition et les propriétés physiques évoluent constamment au cours de l'interaction corium-béton du fait de la décomposition du béton et des réactions chimiques entre composés du mélange liquide, et entre ces derniers et les gaz dégagés. Dans les premières heures de l'interaction, suivant le débit du gaz relâché par le béton et les masses volumiques respectives des phases oxydes et

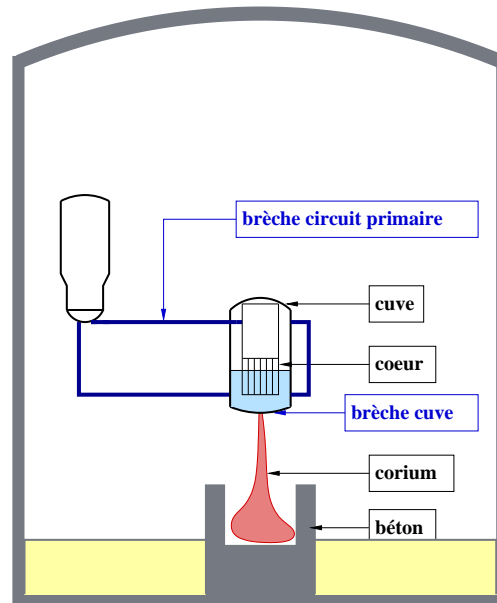


FIG. I.1 – Schéma simplifié d'un réacteur en situation accidentelle.

métalliques, on observe des configurations avec une seule couche (oxyde et métal mélangés) ou avec plusieurs couches (oxyde et métal stratifiés). Par la suite, la phase métallique est entièrement oxydée, notamment par la vapeur d'eau relâchée par le béton ; le bain est alors dit homogène. Cette thèse s'inscrit dans un programme de recherche visant à décrire cette phase tardive des scénarios de fusion du cœur, dans laquelle le mélange corium-béton est composé d'une seule phase oxyde traversée par un flux de bulles (figure I.2). Une approche physique relativement simplifiée de la thermohydraulique du bain est adoptée ; elle est basée sur un modèle de dérive, dont on trouvera la formulation en régime isotherme dans la prochaine section de ce document.

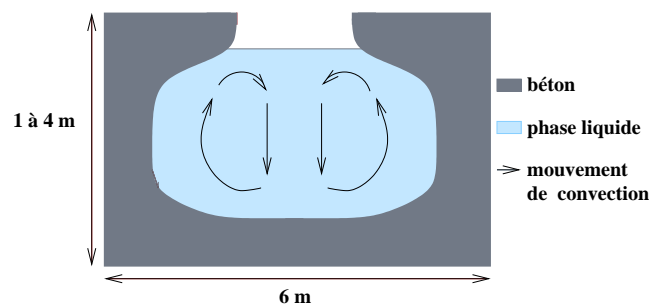


FIG. I.2 – Interaction corium-béton dans une situation homogène.

### I.3 Modèle mathématique

Il existe une grande variété d'écoulements diphasiques. On s'intéresse ici aux écoulements dispersés, caractérisés par le fait que l'une des deux phases est présente sous forme d'inclusions (phase dispersée, ici la phase gazeuse) et l'autre est continue (dans le cas qui nous intéresse, la phase liquide).

A l'échelle microscopique, les écoulements diphasiques peuvent être vus comme un ensemble de régions purement monophasiques, séparées par des interfaces au travers desquelles s'opèrent des transferts de masse, de quantité de mouvement et d'énergie. Pour développer un modèle macroscopique, on écrit d'abord les équations de bilan valables à l'intérieur de chaque région purement monophasique, puis, par des procédés de prise de moyenne, on étend ces équations à tout le domaine et on modélise les transferts aux interfaces [23, 43]. Nous décrivons ici la mise en œuvre de cette démarche pour le problème qui nous intéresse.

### I.3.1 Equations locales instantanées

On considère un mélange non miscible diphasique de liquide (indiqué ici par l'indice  $k = \ell$ ) et de gaz ( $k = g$ ). On adopte le point de vue de la mécanique des milieux continus et on suppose que en tout point  $(x, t)$  de l'espace et du temps, il est possible de définir de la manière habituelle une masse volumique  $\rho$  et une vitesse  $u$  qui sont la masse volumique et la vitesse de la particule fluide présente en  $x$  au temps  $t$  quelle que soit la nature (liquide ou gaz) de cette particule. Cependant, dans le liquide et dans le gaz, les équations instantanées locales régissant l'évolution des variables  $\rho$ ,  $u$  sont différentes. A l'échelle microscopique, on suppose que l'évolution de la phase liquide est gouvernée par les équations de Navier-Stokes pour les écoulements incompressibles tandis que la phase gazeuse est décrite par les équations de Navier-Stokes pour les écoulements compressibles :

Phase liquide :

$$\left\{ \begin{array}{l} \nabla \cdot u_\ell = 0 \\ \frac{\partial u_\ell \rho_\ell}{\partial t} + \nabla \cdot (\rho_\ell u_\ell \otimes u_\ell) + \nabla p_\ell - \nabla \cdot \tau_\ell = \rho_\ell g \\ \rho_\ell = \text{constante} \end{array} \right.$$

Phase gazeuse :

$$\left\{ \begin{array}{l} \frac{\partial \rho_g}{\partial t} + \nabla \cdot (\rho_g u_g) = 0 \\ \frac{\partial u_g \rho_g}{\partial t} + \nabla \cdot (\rho_g u_g \otimes u_g) + \nabla p_g - \nabla \cdot \tau_g = \rho_g g \\ p_g = \rho_g R T \end{array} \right.$$

où  $g$  est la force de gravité et, pour  $k = g, \ell$ ,  $p_k$  est la pression et  $\tau_k$  le tenseur des contraintes visqueuses de la phase  $k$ . La pression et la masse volumique de la phase gazeuse sont reliées par la loi d'état des gaz parfaits,  $R$  étant la constante des gaz parfaits et  $T$  la température, qu'on suppose constante. Ces systèmes d'équations aux dérivées partielles doivent être complétés par des conditions initiales et aux limites.

### I.3.2 Modèle diphasique à l'échelle macroscopique

Après prise de moyenne, le système d'équations de transport du modèle à deux fluides s'écrit sous la forme générale suivante :

$$\left\{ \begin{array}{l} \frac{\partial}{\partial t}(\alpha_k \rho_k) + \nabla \cdot (\alpha_k \rho_k u_k) = \Gamma_k \\ \frac{\partial}{\partial t}(\alpha_k \rho_k u_k) + \nabla \cdot (\alpha_k \rho_k u_k \otimes u_k) + \nabla(\alpha_k p_k) = \nabla \cdot (\alpha_k (\tau_k + \tau_k^T)) + \alpha_k \rho_k g + M_k \end{array} \right.$$

Les grandeurs qui apparaissent dans ce système correspondent maintenant aux grandeurs moyennes associées à la phase  $k$  et  $\alpha_k = V_k/V$  désigne la fraction volumique de la phase  $k$ ,  $V_k \subset V$  étant le volume occupé par celle-ci. Dans la seconde équation les termes  $\tau_k$  et  $\tau_k^T$  représentent respectivement le tenseur des contraintes visqueuses et sa partie turbulente ou tenseur des contraintes de Reynolds.

Les termes  $\Gamma_k$  et  $M_k$  apparaissant aux seconds membres des équations du modèle diphasique représentent les termes de transfert de masse et de quantité de mouvement entre phases [23, 43].

Ces termes vérifient les conditions d'interface :

$$\sum_k \Gamma_k = 0, \quad \sum_k M_k = 0$$

Le terme d'échange de quantité de mouvement  $M_k$  contient plusieurs contributions distinctes :

$$M_k = \Gamma_k u_{kI} + p_{kI} \nabla \alpha_k + F_k^d + \dots$$

La première traduit les effets dus aux transferts de masse aux interfaces, la deuxième exprime les forces de pression aux interfaces et la dernière correspond aux forces de frottement s'exerçant sur la phase  $k$ . Les grandeurs indicées par  $kI$  désignent les grandeurs d'interface moyennées associées à la phase  $k$ .

Il convient maintenant de modéliser les termes d'interaction entre les deux phases, c'est-à-dire les grandeurs moyennes aux interfaces et les termes d'échange. Commençons par les termes de pression. L'hypothèse la plus simple consiste à supposer l'égalité entre pression moyenne et pression d'interface ; on fera ici en outre l'hypothèse que la pression en un même point dans les deux phases est identique :

$$p_\ell = p_g = p_{\ell I} = p_{gI} = p.$$

De plus, on suppose que cette pression commune est toujours reliée à la masse volumique de la phase gazeuse par la loi d'état de cette dernière, donnée dans le paragraphe précédent.

En second lieu, on fait l'hypothèse que le terme de force de frottement interfacial  $F_k^d$  est proportionnel à la différence des vitesses entre les deux phases :  $F_k^d = \lambda (u_{k'} - u_k)$  avec  $\lambda$  paramètre positif. De plus, on négligera le terme de turbulence  $\tau_k^T$  ainsi que le transfert de masse  $\Gamma_k$ . Ceci nous conduit au modèle diphasique suivant :

$$\left\{ \begin{array}{l} \frac{\partial}{\partial t} (\alpha_k \rho_k) + \nabla \cdot (\alpha_k \rho_k u_k) = 0 \\ \frac{\partial}{\partial t} (\alpha_k \rho_k u_k) + \nabla \cdot (\alpha_k \rho_k u_k \otimes u_k) + \alpha_k \nabla p = \nabla \cdot (\alpha_k \tau_k) + \lambda (u_{k'} - u_k) + \alpha_k \rho_k g \end{array} \right. \quad (\text{I.3})$$

### I.3.3 Modèle à vitesse de dérive ou modèle drift-flux

Dans le cadre de nos travaux, le modèle à deux fluides (I.3) apparaît trop complexe et un modèle de mélange peut être suffisant. Ce type de modèle se caractérise par un système d'équations de bilan pour le mélange complété par une équation de transport de la masse volumique partielle de l'une des deux phases [23, 43].

Soient  $\rho$  et  $u$ , respectivement, la masse volumique et la vitesse barycentrique du mélange définies par :

$$\rho u = (\alpha_\ell \rho_\ell + \alpha_g \rho_g) u = \alpha_\ell \rho_\ell u_\ell + \alpha_g \rho_g u_g$$

L'équation de conservation de masse du mélange s'obtient en sommant les équations de conservation de masse de chaque phase :

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \quad (\text{I.4})$$

L'équation de diffusion du modèle s'obtient à partir de l'équation de conservation de la masse de l'une des deux phases en faisant apparaître la vitesse de mélange :

$$\frac{\partial \rho y_g}{\partial t} + \nabla \cdot (\rho y_g u) = -\nabla \cdot (\rho y_g V_{gm}) + \nabla \cdot (D \nabla y_g) \quad (\text{I.5})$$

avec  $y_k$  fraction massique de la phase  $k$  définie par  $\rho y_k = \alpha_k \rho_k$  et  $V_{km} = u_k - u$  vitesse de diffusion de la phase  $k$ . Le terme  $\nabla \cdot (D \nabla y_g)$  a été rajouté au second membre de l'équation précédente afin de tenir compte de l'effet de dispersion de la phase gazeuse [47, 64].

Comme pour l'équation de conservation de la masse, les équations de conservation de la quantité de mouvement du mélange s'obtiennent en sommant les équations de transport relatives à chaque phase du système :

$$\frac{\partial}{\partial t}(\rho u) + \nabla \cdot (\rho u \otimes u) + \nabla p = -\nabla \cdot \left[ \sum_k \alpha_k \rho_k V_{km} \otimes V_{km} \right] + \nabla \cdot \left( \sum_k \alpha_k \tau_k \right) + \rho g \quad (\text{I.6})$$

Pour fermer le modèle, il faut alors établir des lois de fermeture pour la vitesse de diffusion  $V_{km}$  et pour la somme des tenseurs de contraintes visqueuses effectives  $\sum_k \alpha_k \tau_k$ . Afin de garantir l'existence d'une entropie compatible avec les termes du deuxième ordre [58], on posera :

$$\sum_k \alpha_k \tau_k = \tau(u) \quad (\text{I.7})$$

où  $\tau(u)$  est un tenseur effectif de la quantité de mouvement que nous supposerons donné par une loi de comportement de fluide Newtonien pour le mélange :

$$\tau(u) = \mu (\nabla u + \nabla^t u - \frac{2}{3}(\nabla \cdot u) I) \quad (\text{I.8})$$

où  $\mu > 0$  est le coefficient effectif de viscosité dynamique.

Il reste à définir une expression pour la vitesse de diffusion  $V_{km}$ . Soit  $u_r$  la vitesse relative entre les phases, définie par :

$$u_r = u_g - u_\ell$$

En faisant apparaître la vitesse relative dans l'expression de la vitesse de diffusion  $V_{km}$ , on obtient :

$$V_{km} = y_{k'} (u_k - u_{k'})$$

soit, en posant  $y_g = y$  (notation que nous garderons tout au long du mémoire) :

$$V_{\ell m} = -y u_r \quad \text{et} \quad V_{gm} = (1 - y) u_r$$

De même, le tenseur des vitesses de diffusion s'écrit :

$$\sum_k \alpha_k \rho_k V_{km} \otimes V_{km} = \rho (1 - y) y u_r \otimes u_r$$

Pour préciser une loi de fermeture pour la vitesse relative, nous proposons la démarche suivante. En multipliant l'équation de conservation de la quantité de mouvement de la phase  $k = \ell$  par  $\alpha_g \rho_g$ , celle de la phase  $k = g$  par  $\alpha_\ell \rho_\ell$  et en prenant la différence, on obtient :

$$\alpha_\ell \alpha_g \rho_\ell \rho_g \left( \frac{d_\ell u_\ell}{dt} - \frac{d_g u_g}{dt} \right) + \alpha_\ell \rho_\ell \nabla \cdot (\alpha_g \tau_g) - \alpha_g \rho_g \nabla \cdot (\alpha_\ell \tau_\ell) + \alpha_\ell \alpha_g (\rho_g - \rho_\ell) \nabla p = \lambda \rho u_r$$

où on a adopté la notation  $\frac{d_k}{dt} = \frac{\partial}{\partial t} + u_k \cdot \nabla$ . On fait ainsi apparaître une différence d'accélération entre les deux phases que, compte-tenu de la faible masse volumique de la phase gazeuse, nous choisissons de négliger. De plus, dans nombre de cas, la pression, la gravité et le terme de traînée

sont les termes dominants dans les équations de quantité de mouvement des deux phases, et nous négligeons donc les tenseurs des contraintes  $\tau_k$ . La vitesse relative est donnée alors par la relation :

$$\lambda u_r = \frac{\alpha_\ell \alpha_g (\rho_g - \rho_\ell)}{\rho} \nabla p = [\alpha_\ell y - \alpha_g (1 - y)] \nabla p \quad (\text{I.9})$$

qui est l'analogie d'une loi de Darcy pour le déséquilibre des vitesses [67]. L'expression précédente peut être retrouvée par un développement asymptotique de Chapman-Enskog (voir [13] pour une description complète de cette méthode) du système (I.3), dans la limite des temps de relaxation des vitesses tendant vers 0. Pour une preuve de cette affirmation, voir [58, 39]. On prouve ainsi que la vitesse relative est une correction de premier ordre par rapport au temps de relaxation des vitesses. Il est donc légitime de l'introduire dans les équations de conservation de la masse mais il n'est pas justifié du point de vue de l'analyse asymptotique de l'introduire dans les équations de conservation de la quantité de mouvement où elle apparaît sous la forme d'un terme d'ordre deux. On négligera donc ce terme dans l'équation (I.6).

On peut remarquer que, dans ce genre de modèles, la relation de fermeture pour la vitesse relative n'est pas en général obtenue par une analyse asymptotique du système, mais est simplement fournie sous forme algébrique sur la base de corrélations expérimentales [58].

Finalement, on obtient un modèle de mélange fermé, composé de deux équations de continuité (I.4) et (I.5) ainsi que d'une équation de quantité de mouvement pour le mélange (I.6) et complété par la relation (I.9) et la définition (I.7) des flux visqueux. Il est nommé modèle de dérive ou modèle drift-flux et s'écrit sous la forme :

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \quad (\text{I.10})$$

$$\frac{\partial \rho y}{\partial t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho (1 - y) y u_r) = \nabla \cdot (D \nabla y) \quad (\text{I.11})$$

$$\frac{\partial}{\partial t}(\rho u) + \nabla \cdot (\rho u \otimes u) + \nabla p - \nabla \cdot \tau(u) = \rho g \quad (\text{I.12})$$

Le modèle de dérive présente deux avantages : c'est un modèle simple et les équations sont sous la forme conservative, et de ce fait bien adaptées à des discrétisations par des schémas volumes finis. Cependant, les modèles de dérive ont, en général, une structure complexe étant donné que les sous-modèles impliqués, par exemple, les relations donnant la vitesse de dérive ou le terme de friction, ont une forme compliquée donnée très souvent par des tabulations de données empiriques.

## I.4 Quelques résultats fondamentaux en volumes finis

Après avoir défini une discrétisation admissible en volumes finis du domaine, nous présentons, dans ce paragraphe, trois résultats fondamentaux pour la construction des schémas numériques développés dans cette thèse : la stabilité de l'opérateur d'advection (*i.e.* un analogue discret du théorème de l'énergie cinétique), la stabilité induite par le travail des forces de pression et un résultat de stabilité  $L^\infty$  pour une équation d'advection-diffusion non linéaire.

### I.4.1 Discrétisation spatiale et temporelle

On suppose que le problème est posé sur un domaine connexe  $\Omega$ , ouvert et borné de  $\mathbb{R}^d$ ,  $d \leq 3$ , et sur un intervalle de temps  $(0, T)$  discrétisé de manière uniforme avec un pas de temps fixe  $\delta t$ .

Une discrétisation admissible de  $\Omega$  en volumes finis [26] est définie par :

- (i) une famille  $\mathcal{M}$  de polygones ( $d = 2$ ) ou polyèdres ( $d = 3$ ) convexes et disjoints, appelés volumes de contrôle, inclus dans  $\Omega$  et tels que  $\bar{\Omega} = \bigcup_{K \in \mathcal{M}} \bar{K}$ .
- (ii) une famille  $\mathcal{E}$  de sous-espaces bornés d'hyperplans de  $\mathbb{R}^d$  inclus dans  $\bar{\Omega}$ , qui représentent les côtés ( $d = 2$ ) ou les faces ( $d = 3$ ) des volumes de contrôle. On note  $\mathcal{E}_{\text{ext}}$  l'ensemble des faces appartenant au bord de  $\Omega$  et  $\mathcal{E}_{\text{int}}$  l'ensemble des faces internes (*i.e.*  $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ ). Pour tout  $K, L \in \mathcal{M}$ , on suppose que soit  $\bar{K} \cap \bar{L}$  est réduit à l'ensemble vide, à un point ou, pour  $d = 3$ , à un segment, soit  $\bar{K} \cap \bar{L} \in \mathcal{E}_{\text{int}}$ . Dans ce dernier cas, la face commune à  $K$  et  $L$  est notée  $K|L$ .
- (iii) une famille de points de  $\Omega$ ,  $\mathcal{P} = (x_K)_{K \in \mathcal{M}}$ , telle que  $x_K \in K$  pour tout  $K \in \mathcal{M}$  et telle que pour chaque face  $\sigma = K|L$ , la droite passant par  $x_K$  et  $x_L$  soit orthogonale à  $\sigma$ .

Tout au long de ce mémoire, on utilisera les notations suivantes. On note  $\mathcal{E}(K)$  l'ensemble des faces d'une maille  $K \in \mathcal{M}$ . Pour chaque face interne  $\sigma = K|L$ ,  $n_{KL}$  est le vecteur normal à  $\sigma$ , orienté de  $K$  à  $L$  (ainsi  $n_{KL} = -n_{LK}$ ). Les notations  $|K|$  et  $|\sigma|$  désignent, respectivement, la mesure du volume de contrôle  $K$  et de la face  $\sigma$ . Pour tout  $K \in \mathcal{M}$  et tout  $\sigma \in \mathcal{E}(K)$ ,  $d_{K,\sigma}$  représente la distance euclidienne entre  $x_K$  et  $\sigma$ . Pour tout  $\sigma \in \mathcal{E}$ , on définit  $d_\sigma = d_{K,\sigma} + d_{L,\sigma}$ , si  $\sigma \in \mathcal{E}_{\text{int}}$  (dans ce cas  $d_\sigma$  est la distance euclidienne entre  $x_K$  et  $x_L$ ) et  $d_\sigma = d_{K,\sigma}$  si  $\sigma \in \mathcal{E}_{\text{ext}}$ .

Notons que la condition d'orthogonalité (iii) est requise pour permettre la consistance de l'approximation des flux de diffusion, issus par exemple du second membre de l'équation (I.11). Elle n'est pas utile pour la discrétisation d'équations d'advection telles que l'équation (I.10) que nous allons examiner en premier lieu.

On note  $X_{\mathcal{M}}$  l'espace des fonctions constantes par morceaux sur chaque volume de contrôle  $K \in \mathcal{M}$ .

Enfin, pour tout réel  $a$ , on définit  $a^+ = \max(a, 0)$  et  $a^- = -\min(a, 0)$ ; de cette façon  $a = a^+ - a^-$  et  $a^+$  et  $a^-$  sont tous deux positifs.

### I.4.2 Une condition de compatibilité centrale : la satisfaction du bilan de masse

Dans l'ensemble des schémas proposés dans ce mémoire, la masse volumique sera considérée constante par maille. Soient donc  $\rho^*$  et  $\rho$  deux fonctions de  $X_{\mathcal{M}}$  représentant la masse volumique à deux instants consécutifs. Les trois résultats qui font l'objet de la suite de cette section sont tous obtenus sous la même condition de compatibilité entre les flux massiques et les masses volumiques de début et fin de pas, qui s'écrit, avec les notations définies ci-dessus :

$$\forall K \in \mathcal{M}, \quad \rho_K > 0, \quad \rho_K^* > 0 \quad \text{et} \quad |K| \frac{\rho_K - \rho_K^*}{\delta t} + \sum_{\sigma=K|L} F_{\sigma,K} = 0 \quad (\text{I.13})$$

où  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  est une quantité conservative associée à la face  $\sigma$  et au volume de contrôle  $K$ , *i.e.* telle que  $F_{\sigma,K} = -F_{\sigma,L}$ ,  $\forall \sigma = K|L$ , et qui peut être vue comme le flux massique issu de  $K$  traversant la face  $\sigma$ . La relation précédente représente ainsi une approximation du bilan de masse (I.10) sur la maille  $K$  en volumes finis.

### I.4.3 Stabilité de l'opérateur d'advection

Ce premier résultat est un équivalent discret d'une identité de stabilité pour l'opérateur d'advection, qu'on peut écrire pour des fonctions suffisamment régulières  $\rho$ ,  $s$  et  $u$  de la manière suivante :

$$\int_{\Omega} \left[ \frac{\partial \rho s}{\partial t} + \nabla \cdot (\rho s u) \right] s \, dx = \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho s^2 \, dx \quad (\text{I.14})$$

et qui est satisfaite si  $\rho$  et  $u$  vérifient le bilan de masse (I.10). Afin de mieux comprendre le rôle de cette dernière dans l'obtention de l'identité (I.14), nous en donnons ici la preuve (toujours sous hypothèse de régularité de  $\rho$ ,  $s$  et  $u$ ). En développant les dérivées, le premier membre de la relation (I.14) s'écrit :

$$\int_{\Omega} \left[ \frac{\partial \rho s}{\partial t} + \nabla \cdot (\rho s u) \right] s \, dx = \int_{\Omega} \left[ \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) \right] s^2 \, dx + \int_{\Omega} \left[ \rho s \frac{\partial s}{\partial t} + (\rho s u) \cdot \nabla s \right] \, dx$$

La première intégrale de la relation précédente s'annule en raison de (I.10), et la seconde s'écrit en intégrant par parties :

$$\int_{\Omega} \left[ \rho s \frac{\partial s}{\partial t} + (\rho s u) \cdot \nabla s \right] \, dx = \frac{1}{2} \int_{\Omega} \left[ \rho \frac{\partial s^2}{\partial t} - s^2 \nabla \cdot (\rho u) \right] \, dx$$

Soit en utilisant une deuxième fois le bilan de masse :

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

On obtient ainsi la relation (I.14).

Un analogue de cette identité dans le cadre discret est donné par le théorème suivant, que nous démontrons au chapitre II.

#### **Théorème I.4.1**

*Soient  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  trois familles de réels qui vérifient la condition (I.13). Soient  $(s_K^*)_{K \in \mathcal{M}}$  et  $(s_K)_{K \in \mathcal{M}}$  deux familles de réels. Pour toute face interne  $\sigma = K|L$ , on définit  $s_{\sigma}$  soit par  $s_{\sigma} = \frac{1}{2}(s_K + s_L)$ , soit par  $s_{\sigma} = s_K$  si  $F_{\sigma,K} \geq 0$  et  $s_{\sigma} = s_L$  autrement. Le premier choix est nommé "choix centré", le second "choix décentré vers l'amont" ou "choix upwind". Dans les deux cas, on a l'inégalité de stabilité suivante :*

$$\sum_{K \in \mathcal{M}} s_K \left[ \frac{|K|}{\delta t} (\rho_K s_K - \rho_K^* s_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} s_{\sigma} \right] \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K s_K^2 - \rho_K^* s_K^{*2}]$$

Si on pose  $s_K = u_{i,K}$ ,  $\forall K \in \mathcal{M}$ , avec  $u_{i,K}$   $i$ -ème composante de la vitesse relative à maille  $K$ , l'inégalité précédente s'écrit :

$$\sum_K u_{i,K} \left[ \frac{|K|}{\delta t} (\rho_K u_{i,K} - \rho_K^* u_{i,K}^*) + \sum_{\sigma=K|L} F_{\sigma,K} u_{i,K} \right] \geq \frac{1}{2} \sum_K \frac{|K|}{\delta t} [\rho_K (u_{i,K})^2 - \rho_K^* (u_{i,K}^*)^2]$$

On obtient alors, en sommant sur les composantes, un contrôle en norme  $L^2$  de la dérivée en temps discrète de l'énergie cinétique. Cette relation sera fondamentale pour garantir la stabilité de l'équation de quantité de mouvement.

#### **I.4.4 Stabilité induite par le travail des forces de pression**

Ce deuxième résultat consiste en un équivalent discret de la relation

$$\int_{\Omega} -p \nabla \cdot u \, dx = \frac{d}{dt} \int_{\Omega} \rho f(\rho) \, dx \quad (\text{I.15})$$

pour un fluide monophasique compressible, où  $f$  est l'énergie libre associée à l'écoulement ; nous étendrons ensuite ce résultat au cas d'un écoulement diphasique à phases dispersées avec une phase continue incompressible et une phase dispersée compressible.



### Écoulement compressible et barotrope

On considère un écoulement instationnaire compressible et barotrope, c'est à dire au sein duquel la masse volumique est donnée par  $\rho = \varrho(p)$ . On définit une énergie libre pour cet écoulement comme une fonction  $f(\cdot)$  telle que :

$$f'(s) = \frac{\varphi(s)}{s^2} \quad (\text{I.16})$$

où  $\varphi(\cdot)$  est la réciproque de la fonction  $\varrho(\cdot)$ , *i.e.* la fonction qui donne la pression en fonction de la masse volumique du fluide. Ici encore, nous insistons sur l'importance du rôle du bilan de masse dans l'obtention de l'identité (I.15), en donnant d'abord la preuve formelle de cette identité dans le cas continu. Le point de départ de la démonstration est le bilan de masse (I.10) que l'on multiplie par la dérivée de  $\rho f(\rho)$  par rapport à  $\rho$  :

$$[\rho f(\rho)]' \left[ \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) \right] = 0$$

Cette relation implique :

$$\frac{\partial [\rho f(\rho)]}{\partial t} + [\rho f(\rho)]' [u \cdot \nabla \rho + \rho \nabla \cdot u] = 0 \quad (\text{I.17})$$

Ainsi :

$$\frac{\partial [\rho f(\rho)]}{\partial t} + u \cdot \nabla [\rho f(\rho)] + [\rho f(\rho)]' \rho \nabla \cdot u = 0$$

En développant les dérivées, on obtient :

$$\frac{\partial [\rho f(\rho)]}{\partial t} + \underbrace{u \cdot \nabla [\rho f(\rho)] + \rho f(\rho) \nabla \cdot u}_{\nabla \cdot (\rho f(\rho) u)} + \underbrace{\rho^2 [f(\rho)]' \nabla \cdot u}_{\rho \nabla \cdot u} = 0 \quad (\text{I.18})$$

Si l'on suppose que la vitesse obéit à une condition de Dirichlet homogène, chose que nous ferons lors de l'ensemble des développements théoriques présentés dans ce mémoire, on obtient le résultat (I.15) en intégrant simplement la relation précédente sur le domaine  $\Omega$ . Pour obtenir un résultat analogue dans le cadre discret, outre la vérification de la condition de compatibilité (I.13), nous devons faire deux hypothèses supplémentaires :

- (i) la fonction  $s \mapsto s f(s)$  est convexe,
- (ii) L'approximation de la masse volumique à l'interface de deux volumes de contrôle est effectuée avec un choix décentré vers l'amont par rapport au signe de la vitesse normale à la face que, par souci de cohérence avec la discrétisation ultérieure de la vitesse, nous notons  $u_\sigma \cdot n_{KL}$ . Il est à noter que, outre le fait qu'il est nécessaire pour le résultat énoncé ici, ce choix garantit la positivité de la masse volumique.

On a donc  $F_{\sigma,K} = v_{\sigma,K}^+ \rho_K - v_{\sigma,K}^- \rho_L$ , avec  $v_{\sigma,K} = |\sigma| u_\sigma \cdot n_{KL} = v_{\sigma,K}^+ - v_{\sigma,K}^-$ . La relation (I.13) s'écrit alors :

$$\forall K \in \mathcal{M}, \quad \rho_K > 0, \quad \rho_K^* > 0 \quad \text{et} \quad |K| \frac{\rho_K - \rho_K^*}{\delta t} + \sum_{\sigma=K|L} v_{\sigma,K}^+ \rho_K - v_{\sigma,K}^- \rho_L = 0 \quad (\text{I.19})$$

Avec ces notations, on obtient le théorème suivant, qui est démontré au chapitre II :

**Théorème I.4.2 (Stabilité induite par le travail des forces de pression - cas monophasique)**

Soit  $f$  une fonction "énergie libre" définie par (I.16) telle que la fonction  $s \mapsto s f(s)$  soit continuellement différentiable et strictement convexe. Soient  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(v_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  trois familles de réels qui vérifient la condition (I.19), et soit  $(p_K)_{K \in \mathcal{M}}$  une famille de réels telle que  $\varrho(p_K) = \rho_K$ . Alors :

$$\sum_{K \in \mathcal{M}} -p_K \sum_{\sigma=K|L} v_{\sigma,K} \geq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [\rho_K f(\rho_K) - \rho_K^* f(\rho_K^*)] \quad (\text{I.20})$$

**Ecoulement diphasique à phases dispersées**

On considère maintenant un écoulement diphasique à phases dispersées avec une phase continue (phase liquide) incompressible et une phase dispersée (phase gazeuse) compressible, dont la masse volumique est donnée par :

$$\rho = (1 - \alpha_g) \rho_l + \alpha_g \varrho_g(p) \quad (\text{I.21})$$

où on rappelle que  $\alpha_g$  désigne la fraction volumique du gaz et  $\varrho_g(p)$  exprime la masse volumique du gaz en fonction de la pression ; on suppose que la pression est donnée par la loi d'état des gaz parfaits :

$$\varrho_g(p) = \frac{p}{RT} \quad (\text{I.22})$$

avec  $R$  constante des gaz parfaits et  $T$  température absolue. On note  $\rho_l$  la masse volumique (constante) du liquide.

Dans ce cas, on peut choisir pour l'énergie libre de la phase dispersée la quantité  $f(\rho_g)$ , où la fonction  $f(\cdot)$  est donnée par :

$$f(s) = RT \log s \quad (\text{I.23})$$

On vérifie alors que  $f'(\rho_g) = p/\rho_g^2$ . La masse volumique du gaz peut aussi s'écrire comme une fonction de la masse volumique du mélange  $\rho$  et de la masse volumique partielle de la phase gazeuse  $z = \rho y$  de la façon suivante, avec un léger abus de notation ( $\varrho_g$  étant maintenant une fonction de  $\rho, z$ ) :

$$\varrho_g(\rho, z) = \frac{z \rho_l}{z + \rho_l - \rho} \quad (\text{I.24})$$

ce qui permet d'exprimer l'énergie libre de la phase gazeuse comme une fonction de  $\rho$  et de  $z$ .

Comme il est démontré formellement ci-après, pour un écoulement diphasique à phases dispersées, le travail des forces de pression s'écrit :

$$\int_{\Omega} -p \nabla \cdot u \, dx = \frac{d}{dt} \int_{\Omega} z f(\rho, z) \, dx \quad (\text{I.25})$$

Comme pour le cas monophasique, le bilan de masse est nécessaire à la validité de ce résultat ainsi que à celle de son analogue discret. Du fait que  $f$  dépend maintenant également de  $z$ , cette dernière variable doit aussi vérifier la même équation de conservation :

$$\frac{\partial z}{\partial t} + \nabla \cdot (zu) = 0$$

qui correspond au bilan de masse de la phase gazeuse dans le cas où la vitesse de drift et le coefficient de diffusion sont nuls.

Le point de départ de la preuve de la relation (I.25) est la somme du bilan de masse multiplié par la dérivée partielle de  $zf$  par rapport à  $\rho$  et de l'équation précédente multipliée par la dérivée partielle de  $zf$  par rapport à  $z$  :

$$\left(\frac{\partial zf}{\partial \rho}\right)_z \left[\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u)\right] + \left(\frac{\partial zf}{\partial z}\right)_\rho \left[\frac{\partial z}{\partial t} + \nabla \cdot (zu)\right] = 0$$

Par des calculs analogues à ceux effectués dans le cas monophasique on obtient :

$$\frac{\partial zf}{\partial t} + \nabla \cdot (zf u) + \left[\rho \left(\frac{\partial zf}{\partial \rho}\right)_z + z \left(\frac{\partial zf}{\partial z}\right)_\rho - zf\right] \nabla \cdot u = 0 \quad (\text{I.26})$$

On conclut alors en vérifiant que le terme multipliant  $\nabla \cdot u$  est exactement égal à la pression  $p$  et en intégrant sur  $\Omega$ .

Pour obtenir un résultat analogue dans un cadre discret, on démontre tout d'abord que la fonction  $(\rho, z) \mapsto zf(\rho, z)$  est convexe. On doit ensuite choisir, comme dans le cas monophasique, une discrétisation par volumes finis de l'équation de bilan de masse de la phase gazeuse ainsi qu'une discrétisation décentrée vers l'amont des termes d'advection des bilans de masse du mélange et de la phase gazeuse. Ce choix garantit la positivité de la masse volumique.

On montre alors au chapitre III un équivalent de la relation I.25.

#### **Théorème I.4.3 (Stabilité induite par le travail des forces de pression - cas diphasique)**

Soient  $f(\rho, z)$  l'énergie libre de la phase dispersée définie par (I.23) et  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(v_{\sigma, K})_{K \in \mathcal{M}, \sigma=K|L}$  trois familles de réels qui vérifient la condition (I.19). Soient  $(z_K)_{K \in \mathcal{M}}$  et  $(z_K^*)_{K \in \mathcal{M}}$  deux familles de réels positifs telles que :

$$|K| \frac{z_K - z_K^*}{\delta t} + \sum_{\sigma=K|L} v_{\sigma, K}^+ z_K - v_{\sigma, K}^- z_L = 0.$$

Alors pour  $(p_K)_{K \in \mathcal{M}}$  telle que  $(p_K)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(z_K)_{K \in \mathcal{M}}$  vérifient (I.24), on a l'estimation suivante :

$$\sum_{K \in \mathcal{M}} -p_K \sum_{\sigma=K|L} v_{\sigma, K} \geq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [z_K f(\rho_K, z_K) - z_K^* f(\rho_K^*, z_K^*)]$$

#### **I.4.5 Stabilité $L^\infty$ pour une équation d'advection-diffusion non-linéaire**

Soit une discrétisation par volumes finis de l'équation de conservation discrète de la quantité  $\rho s$ , qui s'écrit :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K s_K - \rho_K^* s_K^*) + \sum_{\sigma=K|L} F_{\sigma, K} s_\sigma + \dots [\text{éventuels termes de diffusion}] \dots = 0$$

Larroutou a démontré que la relation (I.13), qui peut être vue comme un bilan de masse discret, garantit une stabilité  $L^\infty$  pour l'inconnue  $s$  de cette équation de transport, pourvu d'utiliser un choix upwind (ou tout autre choix monotone) pour exprimer  $s_\sigma$  [51]. Dans le troisième chapitre de ce mémoire, nous donnons de ce résultat une démonstration différente, mieux adaptée à des

généralisations à des équations non-linéaires, en le présentant comme une conséquence directe du lemme suivant, dont on notera la similitude avec le théorème I.4.1 :

**Lemme I.4.4**

Soient  $(\rho_K)_{K \in \mathcal{M}}$ ,  $(\rho_K^*)_{K \in \mathcal{M}}$  et  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  tels que la condition (I.13) soit vérifiée. On a alors l'inégalité de stabilité suivante :

$$\begin{aligned} - \sum_{K \in \mathcal{M}} s_K^- \left[ \frac{|K|}{\delta t} (\rho_K s_K - \rho_K^* s_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} s_\sigma \right] \\ \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K (s_K^-)^2 - \rho_K^* (s_K^*)^2] \end{aligned} \quad (\text{I.27})$$

avec  $s_\sigma$  défini de la façon suivante :  $s_\sigma = s_K$  si  $F_{\sigma,K} \geq 0$ ,  $s_\sigma = s_L$  autrement.

Considérons maintenant l'équation d'advection-diffusion non-linéaire suivante :

$$\frac{\rho s - \rho^* s^*}{\delta t} + \nabla \cdot (\rho s u) + \nabla \cdot (\rho \varphi(s) u_r^*) = \nabla \cdot (D \nabla s) \quad \text{sur } \Omega \times (0, T) \quad (\text{I.28})$$

où le coefficient de diffusion  $D$  et le vecteur  $u_r^*$  sont des quantités connues et la fonction non-linéaire  $\varphi \in C^1([0, 1], \mathbb{R})$  est telle que  $\varphi(0) = \varphi(1) = 0$ . Nous en proposons la discrétisation par volumes finis suivante :

$$\begin{aligned} \forall K \in \mathcal{M}, \\ |K| \frac{\rho_K s_K - \rho_K^* s_K^*}{\delta t} + \sum_{\sigma=K|L} F_{\sigma,K} s_\sigma \\ + \sum_{\sigma=K|L} [G_{\sigma,K}^+ g(s_K, s_L) - G_{\sigma,K}^- g(s_L, s_K)] + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (s_K - s_L) = 0 \end{aligned} \quad (\text{I.29})$$

où  $G_{\sigma,K}$  est définie par :  $G_{\sigma,K} = \rho_\sigma \int_\sigma u_r^* \cdot n_{KL}$  et  $\rho_\sigma$  est une approximation (positive) de la masse volumique à la face  $\sigma$ . La fonction  $g(\cdot, \cdot)$  est un flux numérique monotone, c'est à dire satisfait la définition suivante (voir [26] pour la théorie complète ainsi que pour quelques exemples) :

**Définition I.4.5 (Flux numérique monotone)**

Soit  $g(\cdot, \cdot)$  une fonction de  $C(\mathbb{R}^2, \mathbb{R})$  vérifiant les hypothèses suivantes :

1.  $g(a_1, a_2)$  est non-décroissante par rapport à  $a_1$  et non-croissante par rapport à  $a_2$ , pour tous réels  $a_1$  et  $a_2$ ,
2.  $g(\cdot, \cdot)$  est Lipschitzienne sur  $\mathbb{R}$  par rapport à ses deux variables,
3.  $g(a, a) = \varphi(a)$ , pour tout  $a \in \mathbb{R}$ .

Alors  $g(\cdot, \cdot)$  est appelé flux numérique monotone.

En multipliant chacune des relations (I.29) par  $s_K^-$ , en sommant sur les volumes de contrôle et en utilisant pour les termes advectifs le résultat précédent (lemme I.4.4) et pour le terme de drift les propriétés des flux numériques monotones, on obtient les deux résultats suivants :

**Lemme I.4.6 (Non-négativité de  $s$ )**

Soient  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  tels que la condition (I.13) soit satisfaite et soit  $g(\cdot, \cdot)$  un flux numérique monotone tel que la fonction  $x \mapsto g(x, x)$  s'annule pour  $x \leq 0$ . Alors, si  $s_K^* \geq 0$ ,  $\forall K \in \mathcal{M}$ , la solution discrète de (I.29) vérifie  $s_K \geq 0$ ,  $\forall K \in \mathcal{M}$ .

**Lemme I.4.7 ( $s$  borné par 1)**

Soient  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  et  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  tels que la condition (I.13) soit satisfaite et soit  $g(\cdot, \cdot)$  un flux numérique monotone tel que  $\varphi(x) = g(x, x)$  s'annule pour  $x \geq 1$ . Alors, si  $s_K^* \leq 1$ ,  $\forall K \in \mathcal{M}$ , la solution discrète de (I.29) vérifie  $s_K \leq 1$ ,  $\forall K \in \mathcal{M}$ .

En utilisant ces deux résultats, on peut par ailleurs prouver, avec un argument de degré topologique, l'existence d'une solution discrète de (I.29), solution dont on démontrera l'unicité par une technique de dualité (chapitre III).

## I.5 Schémas de correction de pression pour les écoulements compressibles

Dans le développement d'un schéma numérique dédié à la résolution du système (I.10)-(I.12), nous sommes confrontés à deux principales difficultés. Tout d'abord, dans les écoulements considérés peuvent coexister des zones liquides (incompressibles) et à fort taux de vide (fortement compressibles); le nombre de Mach est ainsi fortement variable. Le schéma proposé devra donc répondre en premier lieu à l'impératif de conserver ses propriétés de stabilité et de convergence dans une large gamme de nombres de Mach, et notamment, jusqu'à l'incompressible. En second lieu, il devra garder certaines inconnues (ici essentiellement la fraction massique de la phase gazeuse), dans leurs bornes physiques. Un schéma répondant à cette dernière contrainte a été proposé dans le paragraphe précédent (section I.4.5). Pour résoudre la première difficulté, on peut soit étendre des schémas dédiés au calcul des écoulements compressibles, en développant des techniques de préconditionnement de solveurs de Riemann [40], soit choisir comme point de départ des schémas pour les écoulements incompressibles, en extrapolant les méthodes de correction de pression largement utilisées dans ce contexte [14, 68, 7]; c'est cette dernière approche que nous adoptons ici.

La discrétisation spatiale choisie utilise des éléments finis mixtes non conformes (éléments finis de Crouzeix-Raviart ou Rannacher-Turek). Elle permet une approximation naturelle des termes visqueux et elle est intrinsèquement stable (*i.e.* sans ajout de termes de stabilisation pour obtenir la condition dite *inf-sup* ou condition de Babuska-Brezzi, voir par exemple [9]) dans la limite de l'incompressible. Cette discrétisation est décrite dans la première partie de cette section.

En s'inspirant d'un certain nombre de travaux dont ceux menés par Wesseling *et al.*, nous présentons ensuite un schéma de correction de pression pour les équations de Navier-Stokes compressibles et barotropes, sous-problème que l'on obtient en posant  $y = 1$  dans le système (I.10)-(I.12). Enfin, sur la base de ces travaux, un schéma conservatif, monotone et stable indépendamment du nombre de Mach est construit pour le modèle de dérive. Nous présentons ce schéma dans la dernière partie de cette section.

### I.5.1 Discrétisation spatiale

Soit  $\mathcal{M}$  une partition du domaine  $\Omega$  en quadrilatères convexes ( $d = 2$ ), en hexahédres ( $d = 3$ ) ou en simplexes. On note  $\mathcal{E}$  et  $\mathcal{E}(K)$  l'ensemble de toutes les  $(d - 1)$ -faces  $\sigma$  du maillage et de l'élément  $K \in \mathcal{M}$  respectivement. On nomme  $\mathcal{E}_{\text{ext}}$  l'ensemble des faces appartenant au bord de  $\Omega$  et  $\mathcal{E}_{\text{int}}$  l'ensemble des faces internes (*i.e.*  $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ ). La partition de  $\mathcal{M}$  est supposée régulière dans le sens usuel de la littérature éléments finis (e.g. [15]), et, en particulier,  $\mathcal{M}$  vérifie les propriétés suivantes :  $\Omega = \bigcup_{K \in \mathcal{M}} \bar{K}$  ; si  $K, L \in \mathcal{M}$ , alors soit  $\bar{K} \cap \bar{L}$  est réduit à l'ensemble vide, à un point ou, pour  $d = 3$ , à un segment, soit  $\bar{K} \cap \bar{L}$  est la face commune de  $K$  et  $L$ , qu'on note  $K|L$ . Pour chaque face interne du maillage  $\sigma = K|L$ ,  $n_{KL}$  désigne le vecteur normal à  $\sigma$ , orienté de  $K$  vers  $L$ . De plus,  $|K|$  et  $|\sigma|$  représentent, respectivement, la mesure de  $K$  et de la face  $\sigma$ .

La discrétisation spatiale utilise une technique d'éléments finis mixtes non conformes : éléments finis de Crouzeix-Raviart (voir [18] pour le papier originel et, par exemple, [24, p. 83–85] pour une présentation synthétique) pour les maillages triangulaires, ou de Rannacher-Turek [61] pour les maillages quadrilatéraux ou hexahédraux. L'élément de référence  $\hat{K}$  pour l'élément fini de Rannacher-Turek est le  $d$ -cube unitaire (avec des faces parallèles aux axes de coordonnées) ; l'espace des fonctions de forme pour l'élément  $\hat{K}$  est  $\tilde{Q}_1(\hat{K})^d$ , où  $\tilde{Q}_1(\hat{K})$  est défini de la façon suivante :

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

L'élément de référence pour l'élément fini de Crouzeix-Raviart est le  $d$ -simplexe unitaire et l'espace des fonctions de forme est l'espace  $P_1$  des polynômes affines. Pour les deux types d'éléments utilisés ici, les degrés de liberté sont déterminés par l'ensemble des fonctions de forme globales suivant :

$$\{F_{\sigma,i}, \sigma \in \mathcal{E}(K), i = 1, \dots, d\}, \quad F_{\sigma,i}(v) = |\sigma|^{-1} \int_{\sigma} v_i \, d\gamma \quad (\text{I.30})$$

La transformation de l'élément de référence en un élément quelconque du maillage est, pour l'élément de Rannacher-Turek, la transformation  $Q_1$  standard et, pour l'élément de Crouzeix-Raviart, la transformation affine standard. Enfin, dans les deux cas, on demande que la valeur moyenne des vitesses discrètes (*i.e.*, pour tout champ de vitesse discret  $v$ ,  $F_{\sigma,i}(v)$ ) soit continue sur chaque face du maillage. En prenant en compte les conditions aux limites de Dirichlet homogènes, l'espace discret  $W_h$  est ainsi défini de la façon suivante :

$$\begin{aligned} W_h = \{ & v_h \in L^2(\Omega) : v_h|_K \in W(K)^d, \forall K \in \mathcal{M} ; \\ & F_{\sigma,i}(v_h) \text{ continu sur chaque face } \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d ; \\ & F_{\sigma,i}(v_h) = 0, \forall \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq i \leq d \} \end{aligned}$$

où  $W(K)$  est l'espace des fonctions sur  $K$ , générées à partir de  $\tilde{Q}_1(\hat{K})$  par la transformation  $Q_1$  standard pour l'élément de Rannacher-Turek et l'espace des fonctions affines sur  $K$  pour l'élément de Crouzeix-Raviart.

Pour la discrétisation de Rannacher-Turek comme pour celle de Crouzeix-Raviart, la pression est approchée par l'espace des fonctions continues par morceaux  $L_h$  :

$$L_h = \{q_h \in L^2(\Omega) : q_h|_K = \text{constante}, \forall K \in \mathcal{M}\}$$

Etant donné que seule la continuité de l'intégrale aux faces du maillage est requise, les vitesses peuvent être discontinues aux faces ; la discrétisation est ainsi non conforme dans  $H^1(\Omega)^d$ . Ces paires d'espaces d'approximation pour la vitesse et la pression sont *inf-sup* stables, au sens usuel des vitesses discrètes "H<sup>1</sup> par morceaux", *i.e.* il existe  $c_i > 0$  indépendant du maillage tel que :

$$\forall p \in L_h, \quad \sup_{v \in W_h} \frac{\int_{\Omega, h} p \nabla \cdot v \, dx}{\|v\|_{1, b}} \geq c_i \|p - \bar{p}\|_{L^2(\Omega)}$$

où  $\bar{p}$  est la valeur moyenne de  $p$  sur  $\Omega$ , le symbole  $\int_{\Omega,h}$  désigne  $\sum_{K \in \mathcal{M}} \int_K$  et  $\|\cdot\|_{1,b}$  est la semi-norme de l'espace de Sobolev  $H^1$  "brisée" :

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

Avec la définition (I.30), chaque degré de liberté de la vitesse peut être associé de façon univoque à une face de l'élément. Les degrés de liberté de la vitesse seront alors indexés par le numéro de la composante et de la face auxquelles ils sont associés, si bien que l'ensemble des degrés de liberté de la vitesse s'écrit :

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

On définit  $v_\sigma = \sum_{i=1}^d v_{\sigma,i} e_i$  où  $e_i$  est le  $i$ -ème vecteur de la base canonique de  $\mathbb{R}^d$ . On note  $\varphi_\sigma^{(i)}$  la fonction de base associée à  $v_{\sigma,i}$ , qui, par définition des éléments finis considérés, s'écrit :

$$\varphi_\sigma^{(i)} = \varphi_\sigma e_i$$

où  $\varphi_\sigma$  est une fonction scalaire. De manière analogue, chaque degré de liberté de la pression est associé à la maille  $K$ , et on désigne l'ensemble des degrés de liberté de la pression par  $\{p_K, K \in \mathcal{M}\}$ .

### I.5.2 Schéma de correction de pression pour les équations de Navier-Stokes barotropes

Si l'on suppose  $y \equiv 1$  dans le système (I.10)-(I.12), on obtient les équations de Navier-Stokes pour un écoulement compressible, barotrope et instationnaire. La solution de ce problème dans le cadre continu [52, 28, 57] satisfait les trois estimations *a priori* suivantes :

(i) stricte positivité de la masse volumique :

$$\rho(x,t) > 0, \quad \forall x \in \Omega, \forall t \in (0, T)$$

si cette propriété est respectée par la condition initiale ;

(ii) conservation de la masse :

$$\int_{\Omega} \rho(x,t) dx = \int_{\Omega} \rho(x,0) dx, \quad \forall t \in (0, T) \quad (\text{I.31})$$

(iii) identité d'énergie qui garantit le contrôle de la quantité :

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho(x,t) |u(x,t)|^2 dx + \frac{d}{dt} \int_{\Omega} \rho(x,t) f(\rho(x,t)) dx \\ + \int_{\Omega} \tau(u(x,t)) : \nabla u(x,t) dx, \quad \forall t \in (0, T) \end{aligned}$$

en fonction du second membre et des conditions initiales et aux limites du problème,

où  $u$  désigne la vitesse,  $\tau$  le tenseur de cisaillement et  $f$  l'énergie libre de l'écoulement définie par (I.16).

La quantité  $\frac{1}{2} \int_{\Omega} \rho(x,t) |u(x,t)|^2 dx + \int_{\Omega} \rho(x,t) f(\rho(x,t)) dx$  est appelée entropie du système.

La démonstration de (I.31-(iii)) repose d'une part sur une propriété de stabilité de l'opérateur d'advection qui s'écrit :

$$\int_{\Omega} \left[ \frac{\partial}{\partial t}(\rho u) + \nabla \cdot (\rho u \otimes u) \right] \cdot u \, dx = \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho u^2 \, dx \quad (\text{I.32})$$

Notons que cette propriété s'obtient directement en prenant  $s = u_i$  dans I.14 et en sommant pour  $i = 1, \dots, d$ ; un équivalent discret s'obtiendrait donc en utilisant de même le théorème I.4.1. On utilise d'autre part dans cette démonstration la relation I.15 (stabilité induite par le travail des forces de pression), dont on a introduit un équivalent discret au paragraphe I.4.4 (théorème I.4.2). Nous allons maintenant présenter un schéma à pas fractionnaires, de type correction de pression, qui s'appuie sur ces résultats discrets pour satisfaire les mêmes propriétés de stabilité.

### Algorithme semi-discret

Nous proposons l'algorithme suivant pour la résolution des équations de Navier-Stokes compressibles :

$$1 - \text{Trouver } \tilde{\rho}^{n+1} \text{ tel que : } \frac{\tilde{\rho}^{n+1} - \rho^n}{\delta t} + \nabla \cdot (\tilde{\rho}^{n+1} u^n) = 0 \quad (\text{I.33})$$

$$2 - \text{Trouver } \tilde{p}^{n+1} \text{ tel que : } -\nabla \cdot \left( \frac{1}{\tilde{\rho}^{n+1}} \nabla \tilde{p}^{n+1} \right) = -\nabla \cdot \left( \frac{1}{\sqrt{\tilde{\rho}^{n+1}} \sqrt{\rho^n}} \nabla p^n \right) \quad (\text{I.34})$$

3 - Trouver  $\tilde{u}^{n+1}$  tel que :

$$\frac{\tilde{\rho}^{n+1} \tilde{u}^{n+1} - \rho^n u^n}{\delta t} + \nabla \cdot (\tilde{\rho}^{n+1} u^n \otimes \tilde{u}^{n+1}) + \nabla \tilde{p}^{n+1} - \nabla \cdot \tau(\tilde{u}^{n+1}) = \tilde{\rho}^{n+1} g^{n+1} \quad (\text{I.35})$$

4 - Trouver  $\bar{u}^{n+1}, p^{n+1}, \rho^{n+1}$  tel que :

$$\left\{ \begin{array}{l} \tilde{\rho}^{n+1} \frac{\bar{u}^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - \tilde{p}^{n+1}) = 0 \\ \frac{\varrho(p^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho(p^{n+1}) \bar{u}^{n+1}) = 0 \\ \rho^{n+1} = \varrho(p^{n+1}) \end{array} \right. \quad (\text{I.36})$$

$$5 - \text{Trouver } u^{n+1} \text{ tel que : } \sqrt{\rho^{n+1}} u^{n+1} = \sqrt{\tilde{\rho}^{n+1}} \bar{u}^{n+1} \quad (\text{I.37})$$

La première étape de l'algorithme est une prédiction de la masse volumique, utilisée pour la discrétisation du terme instationnaire dans le bilan de quantité de mouvement. Comme observé par Wesseling *et al* [7, 75], cette étape peut être évitée dans la résolution des équations d'Euler : en effet, dans ce cas, le débit massique peut être choisi comme inconnue, en utilisant la vitesse explicite comme un champ advectif dans la discrétisation du terme de convection dans le bilan de quantité de mouvement ; la vitesse est alors mise à jour en divisant le débit massique par la masse volumique de fin de pas. En revanche, en présence d'un flux visqueux, si la discrétisation du terme de diffusion est choisie implicite, le débit massique et la vitesse apparaissent comme inconnues dans le bilan de quantité de mouvement ; l'ajout d'une étape de prédiction de la masse volumique devient alors nécessaire.

La deuxième étape, introduite afin de garantir la stabilité du schéma [31], correspond à une renormalisation de la pression. Une technique similaire a déjà été utilisée par Guermond et Quartapelle pour des écoulements incompressibles à masse volumique variable [37].

L'étape 3 consiste en une résolution semi-implicite classique de l'équation de conservation de la quantité de mouvement, obtenant ainsi une vitesse prédite.

L'étape 4 est une étape de correction de pression non linéaire. En prenant la divergence de la première équation de (I.36) et en utilisant la seconde pour éliminer l'inconnue  $\bar{u}^{n+1}$ , on obtient un



problème elliptique non linéaire pour la pression. Une fois la pression calculée, la première et la troisième relation fournissent, respectivement, la vitesse et la masse volumique de fin de pas.

Enfin la dernière étape est une renormalisation de la vitesse introduite, comme la seconde étape, pour garantir la stabilité du schéma.

Le schéma dégénère en une méthode de projection incrémentale standard pour un écoulement incompressible (e.g. [53]).

### Algorithme discret

La vérification de l'équation (I.32) nécessite le développement d'une discrétisation *ad hoc* de l'opérateur d'advection. Pour ce faire, nous utilisons une technique de lumping pour la discrétisation du terme instationnaire du bilan de quantité de mouvement. Nous pouvons par ailleurs remarquer que dans le cas bidimensionnel la matrice de masse est déjà diagonale avec l'élément fini de Crouzeix-Raviart. Soit  $|D_\sigma|$  tel que :

$$|D_\sigma| \stackrel{\text{def}}{=} \int_{\Omega} \varphi_\sigma \, dx > 0 \quad (\text{I.38})$$

Pour toute face interne  $\sigma = K|L$  et tout volume de contrôle  $K$ , soit  $D_{K,\sigma}$  le cône de base  $\sigma$  et de sommet le centre de masse de la maille. Le volume  $D_{K,\sigma}$  désigne la demi-maille diamant relative à  $\sigma$  et à  $K$  (voir figure I.3). On nomme  $|D_{K,\sigma}|$  la mesure de  $D_{K,\sigma}$ . Pour chaque face  $\sigma \in \mathcal{E}$ , on définit la "maille diamant"  $D_\sigma$  relative à  $\sigma$  par  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$  si  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , et par  $D_\sigma = D_{K,\sigma}$  si  $\sigma \in \mathcal{E}_{\text{ext}}$ ,  $K$  désignant, dans ce cas, le seul volume de contrôle adjacent à  $\sigma$ . Ainsi, pour l'élément fini de Crouzeix-Raviart,  $|D_\sigma|$  représente la mesure de la maille diamant relative à  $\sigma$ . La même propriété est valable pour l'élément fini de Rannacher-Turek pour des rectangles ( $d = 2$ ) ou des parallélépipèdes rectangles ( $d = 3$ ), seuls cas considérés ici.

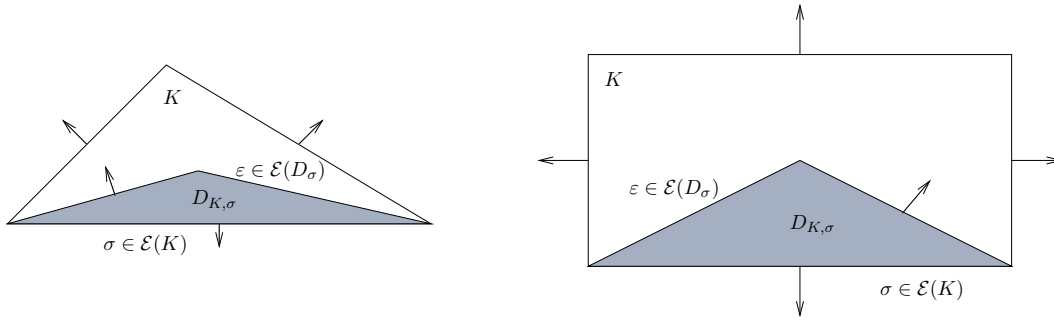


FIG. I.3 – Demi-mailles diamant pour les éléments finis de Crouzeix-Raviart et de Rannacher-Turek.

La discrétisation du terme instationnaire dans les équations associées à la vitesse relative à la face  $\sigma$  aboutit donc à une expression de la forme suivante :

$$\sum_{K \in \mathcal{T}} \int_K \frac{\rho_K^{n+1} u^{n+1} - \rho_K^n u^n}{\delta t} \varphi_\sigma \, dx \rightarrow \sum_{D_\sigma} |D_\sigma| \frac{\rho_\sigma^{n+1} u_\sigma^{n+1} - \rho_\sigma^n u_\sigma^n}{\delta t}, \quad \sigma \in \mathcal{E}_{\text{int}}$$

où  $\rho_\sigma^n$  est la moyenne des masses volumiques dans les mailles adjacentes à  $\sigma$ , pondérées par la mesure des demi-diamants :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad |D_\sigma| \rho_\sigma^n = |D_{K,\sigma}| \rho_K^n + |D_{L,\sigma}| \rho_L^n \quad (\text{I.39})$$

On retrouve ainsi une approximation de type volumes finis basée sur un maillage dual (mailles diamant).

Afin de vérifier la condition de stabilité (I.13), une prédiction de la masse volumique au temps  $t^{n+1}$  est nécessaire; elle est obtenue en résolvant le bilan de masse sur ces mêmes volumes de contrôle :

$$\forall \sigma \in \mathcal{E}, \quad \frac{|D_\sigma|}{\delta t} (\tilde{\rho}_\sigma^{n+1} - \rho_\sigma^n) + \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon, \sigma}^{n+1} = 0 \quad (\text{I.40})$$

où  $\mathcal{E}(D_\sigma)$  est l'ensemble des faces de  $D_\sigma$  et  $F_\varepsilon^{n+1}$  le flux massique associé à  $\varepsilon \in \mathcal{E}(D_\sigma)$ , défini par :

$$F_\varepsilon^{n+1} = |\varepsilon| u_\varepsilon^n \cdot n_{\varepsilon, \sigma} \tilde{\rho}_\varepsilon^{n+1}$$

où  $|\varepsilon|$  est la mesure de  $\varepsilon$ ,  $n_{\varepsilon, \sigma}$  est le vecteur normal à  $\varepsilon$  sortant de  $D_\sigma$ , la vitesse  $u_\varepsilon^n$  est obtenue par interpolation au centre de  $\varepsilon$  du champ de vitesse  $u^n$  et  $\tilde{\rho}_\varepsilon^{n+1}$  est la masse volumique à la face, calculée avec une technique upwind standard (*i.e.* soit  $\tilde{\rho}_\sigma^{n+1}$  si  $u_\varepsilon^n \cdot n_{\varepsilon, \sigma} \geq 0$ , soit  $\tilde{\rho}_{\sigma'}^{n+1}$ , avec  $\sigma'$  tel que  $\varepsilon = D_\sigma | D_{\sigma'}$ ). Cette discrétisation du flux est également utilisée dans le terme de convection non linéaire  $\nabla \cdot (\tilde{\rho}^{n+1} u^n \otimes \tilde{u}^{n+1})$  de l'étape 3 de prédiction de vitesse, avec un choix centré pour la vitesse convectée; pour  $i = 1, \dots, d$ , la  $i$ -ème composante du terme de convection non linéaire discret associé à la face  $\sigma$  s'écrit donc :

$$\sum_{\substack{\varepsilon \in \mathcal{E}(D_\sigma), \\ \varepsilon = D_\sigma | D_{\sigma'}}} \frac{1}{2} F_{\varepsilon, \sigma}^{n+1} (\tilde{u}_{\sigma, i}^{n+1} + \tilde{u}_{\sigma', i}^{n+1})$$

On retrouve alors les hypothèses permettant d'appliquer le résultat de stabilité I.4.1, avec  $s = u_i$ , ce qui assure la stabilité de l'opérateur d'advection du bilan de quantité de mouvement.

Une discrétisation standard par éléments finis est utilisée pour le gradient de pression ainsi que pour le terme de viscosité :

$$\nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) \quad \rightarrow \quad - \sum_{K \in \mathcal{M}} \int_K p^n \nabla \cdot \varphi_\sigma^{(i)} + a_d(\tilde{u}^{n+1}, \varphi_\sigma^{(i)})$$

La forme bilinéaire  $a_d(\cdot, \cdot)$  représente le terme de viscosité et est définie de la façon suivante, pour tout  $v$  et  $w$  appartenant à  $W_h$  :

$$a_d(v, w) = \begin{cases} \mu \int_{\Omega, h} \left[ \nabla v : \nabla w + \frac{1}{3} \nabla \cdot v \nabla \cdot w \right] & \text{si } \mu \text{ est constant,} \\ \int_{\Omega, h} \tau(v) : \nabla w & \text{avec } \tau \text{ donné par (I.8) autrement.} \end{cases}$$

Du fait que les pressions sont constantes par maille, la discrétisation spatiale du bilan de masse s'écrit sous la forme d'un bilan de type volume fini. En s'appuyant sur cette caractéristique du schéma, nous démontrons qu'une discrétisation de cette équation décentrée dans le sens de la vitesse, outre le fait qu'elle garantit la conservation de la masse et la positivité de la masse volumique, permet également de vérifier un analogue discret de l'identité (I.15), par application du théorème I.4.2. Ce décentrement permet donc de profiter de la stabilité induite par le travail des forces de pression, qui fournit *in fine* un ensemble d'estimations *a priori* pour la vitesse et la pression assurant l'existence de la solution pour l'étape de correction de pression.

De plus, les résultats I.4.1 et I.4.2 permettent d'obtenir l'analogue discret de la propriété I.31-(iii), démontrant ainsi que le schéma est stable dans le sens où l'entropie discrète est décroissante :

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^{n+1}}^2 + \int_{\Omega} \rho^{n+1} f(\rho^{n+1}) dx + \delta t \sum_{k=1}^{n+1} a_d(\tilde{u}^k, \tilde{u}^k) + \frac{\delta t^2}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 \\ \leq \frac{1}{2} \|u^0\|_{h, \rho^0}^2 + \int_{\Omega} \rho^0 f(\rho^0) dx + \frac{\delta t^2}{2} |p^0|_{h, \tilde{\rho}^0}^2 \end{aligned}$$

où, pour tout  $n$ ,  $\|\cdot\|_{h,\rho}^2$  définit une norme  $L^2$  discrète sur  $W_h$  pondérée par  $\rho$  et  $|\cdot|_{h,\rho}^2$  peut être interprétée comme une semi-norme  $H^1$  discrète classique dans le contexte des volumes finis, pondérée par  $1/\rho$ .

### Propriétés du schéma numérique

Les différentes propriétés énoncées ci-dessus sont récapitulées dans le théorème suivant :

#### **Théorème I.5.8 (Propriétés du schéma numérique)**

Soient  $(u^n)_{1 \leq n \leq N}$ ,  $(p^n)_{1 \leq n \leq N}$  et  $(\rho^n)_{1 \leq n \leq N}$  les solutions du schéma numérique décrit ci-dessus, avec une condition initiale positive pour  $\rho$ . On suppose que le terme visqueux est dissipatif (i.e.  $\forall v \in W_h$ ,  $a_d(v, v) \geq 0$ , ce qui est vérifié avec l'expression proposée dans le cas où la viscosité est constante). Si la loi d'état  $\varrho(\cdot)$  satisfait les hypothèses suivantes :

1.  $\varrho(\cdot)$  est croissante,  $\varrho(0) = 0$  et  $\lim_{z \rightarrow +\infty} \varrho(z) = +\infty$ ,
2. il existe une énergie libre  $f$  telle que (I.16) soit vérifiée et la fonction  $s \mapsto s f(s)$  soit continue, différentiable, strictement convexe et telle que  $s f(s) \geq -C_P$ ,  $\forall s \in (0, +\infty)$ , avec  $C_P$  constante non-négative,

alors le schéma (I.33)-(I.37) vérifie les propriétés suivantes, pour tout  $n \leq N$  :

- (i) il existe une solution à chaque étape du schéma ;
- (ii) les inconnues restent dans leurs bornes physiques :

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

- (iii) le schéma est conservatif :

$$\sum_{K \in \mathcal{M}} |K| \rho_K^n = \sum_{K \in \mathcal{M}} |K| \rho_K^0;$$

- (iv) le schéma est stable (i.e. l'entropie discrète décroît).

**Remarque I.5.9** Pour démontrer l'existence d'une solution à l'étape de correction de pression nous utilisons une technique de degré topologique qui peut s'appliquer à la discrétisation considérée ici de tout problème du type :

$$\left\{ \begin{array}{ll} A u + \nabla p = f & \text{in } \Omega \\ \frac{\varrho(p) - \rho^*}{\delta t} + \nabla \cdot (\varrho(p) u) = 0 & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{array} \right. \quad (\text{I.41})$$

où  $A$  désigne un opérateur elliptique quelconque et la loi d'état est donnée par :

$$\varrho(p) = p^\gamma \quad \text{avec } \gamma \geq 1$$

On obtient, en particulier, l'existence d'une solution discrète pour l'algorithme proposé dans le cas où  $\rho = p$ , alors qu'en continu elle n'est démontrée pour les équations de Navier-Stokes que dans le cas  $\gamma > d/2$  [52, 28, 57].

Cette méthode numérique est testée sur un cas possédant par construction une solution analytique, afin d'en vérifier la convergence (chapitre II). Les erreurs de la vitesse et de la pression,

respectivement en norme  $L^2$  et en norme  $L^2$  discrète en fonction du pas de temps, pour différents maillages, montrent d'abord une décroissance qui correspond à une convergence en temps d'ordre environ un pour la vitesse et la pression, avant d'atteindre un plateau, dû au fait que les erreurs sont bornées inférieurement par l'erreur résiduelle de discrétisation en espace. Pour la vitesse ainsi que pour la pression, la valeur des erreurs sur le plateau montre une convergence en espace (en norme  $L^2$ ) d'ordre compris entre un et deux.

### I.5.3 Schéma de correction de pression pour le modèle à vitesse de dérive

Nous allons maintenant nous intéresser au modèle de dérive pour un écoulement diphasique à phases dispersées. Nous rappelons que la phase continue (dans notre cas, la phase liquide) est supposée incompressible, tandis que la phase dispersée (ici, la phase gazeuse) a un comportement barotrope.

#### Algorithme semi-discret

Le schéma développé entre dans la classe des méthodes de correction de pression et s'écrit dans un formalisme semi-discret sous la forme suivante :

1 - Trouver  $\tilde{u}^{n+1}$  tel que :

$$\frac{\rho^n \tilde{u}^{n+1} - \rho^{n-1} u^n}{\delta t} + \nabla \cdot (\rho^n u^n \otimes \tilde{u}^{n+1}) + \nabla p^n + \nabla \cdot \tau(\tilde{u}^{n+1}) = f^{n+1} \quad (\text{I.42})$$

2 - Trouver  $p^{n+1}$ ,  $u^{n+1}$ ,  $\rho^{n+1}$  et  $z^{n+1}$  tels que :

$$\left\{ \begin{array}{l} \rho^n \frac{u^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - p^n) = 0 \\ \frac{\varrho(p^{n+1}, z^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho(p^{n+1}, z^{n+1}) u^{n+1}) = 0 \\ \frac{z^{n+1} - \rho^n y^n}{\delta t} + \nabla \cdot (z^{n+1} u^{n+1}) = 0 \\ \rho^{n+1} = \varrho(p^{n+1}, z^{n+1}) \end{array} \right. \quad (\text{I.43})$$

3 - Trouver  $y^{n+1}$  tel que :

$$\frac{\rho^{n+1} y^{n+1} - z^{n+1}}{\delta t} + \nabla \cdot (\rho^{n+1} y^{n+1} (1 - y^{n+1}) u_r^{n+1}) = \nabla \cdot (D \nabla y^{n+1}) \quad (\text{I.44})$$

La première étape est une résolution semi-implicite du bilan de quantité de mouvement qui permet d'obtenir la vitesse prédite  $\tilde{u}^{n+1}$ .

L'étape de correction de pression (étape 2) consiste en un problème elliptique non linéaire pour la pression (deux premières équations) couplé avec les termes de transport du bilan de masse de la phase gazeuse (troisième équation). Afin de faciliter la convergence de l'algorithme de Newton utilisé pour résoudre cette étape, nous introduisons la variable  $z = \rho y$ . En effet, la loi d'état  $\varrho(\cdot, \cdot)$  s'écrit alors comme une fonction de la pression et de  $z$  linéaire par rapport à  $z$  :

$$\varrho(p, z) = z \left( 1 - \frac{\rho_l R T}{p} \right) + \rho_l \quad (\text{I.45})$$

La dernière étape, enfin, prend en compte les termes restants du bilan de masse de la phase gazeuse et permet de calculer la fraction massique  $y$  de fin de pas.

### Algorithme discret

Dans l'étape de prédiction, comme dans le cas monophasique, l'opérateur d'advection est discrétisé par une technique de volumes finis s'appuyant sur un maillage dual (mailles diamant). On obtient ainsi l'approximation suivante pour le terme de convection :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d \quad \text{terme de convection} \rightarrow \sum_{\substack{\varepsilon \in \mathcal{E}(D_\sigma), \\ \varepsilon = D_\sigma | D_{\sigma'}}} \frac{1}{2} F_{\varepsilon, \sigma}^{n+1} (\tilde{u}_{\sigma, i}^{n+1} + \tilde{u}_{\sigma', i}^{n+1}) \quad (\text{I.46})$$

où le flux massique  $F_{\varepsilon, \sigma}^n$  est défini de la façon suivante :

$$F_{\varepsilon, \sigma}^n = |\varepsilon| (\tilde{\rho} u)_{|\varepsilon} \cdot n_{\varepsilon, \sigma}$$

Pour garantir sa stabilité (*i.e.* vérifier le théorème I.4.1), il est nécessaire de satisfaire à une condition de compatibilité entre les flux massiques et les masses volumiques de début et de fin de pas (I.13) analogue au bilan de masse, alors que ce dernier n'est pas résolu à ce stade. Nous pouvons donc soit utiliser une étape de prédiction de la masse volumique, comme dans le cas monophasique, soit effectuer un simple décalage en temps des masses volumiques ; c'est ce dernier choix que nous effectuons, parce qu'il permet d'obtenir un schéma conservatif vis à vis de la quantité de mouvement. La condition de compatibilité devra être, dans ce dernier cas, vérifiée grâce au bilan de masse résolu dans l'étape de correction de pression du pas de temps précédent. Or, cette dernière est écrite sur les mailles primales alors que la condition de compatibilité (I.13) traduit un bilan de masse sur les mailles diamants. Pour venir à bout de cette difficulté, nous utiliserons le résultat suivant [3].

#### Lemme I.5.10 (Bilan de masse sur le sous-volume d'une maille)

Soit  $K \in \mathcal{M}$ . On suppose qu'il existe deux nombres réels  $\rho^*$  et  $\rho$ , constants sur  $K$  et une famille de flux  $(|\sigma| (\rho u)_\sigma \cdot n_\sigma)_{\sigma \in \mathcal{E}(K)}$  tels que :

$$\frac{|K|}{\delta t} (\rho - \rho^*) + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma \cdot n_\sigma = 0 \quad (\text{I.47})$$

On suppose qu'il existe un débit massique  $w$ , régulier sur  $K$ , tel que  $\nabla \cdot w$  soit constant sur  $K$  et tel que :

$$\int_K \nabla \cdot w = \int_{\partial K} w \cdot n_K = \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma \cdot n_\sigma$$

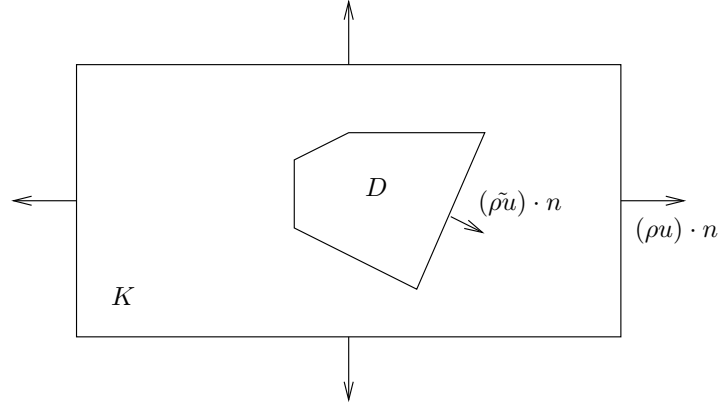
où  $\partial K$  et  $n_K$  désignent, respectivement, le bord de  $K$  et le vecteur normal à  $\partial K$  sortant de  $K$ . Soit  $D$  un sous-volume de  $K$ , de bord  $\partial D$  (voir figure I.4). Alors, la propriété suivante est satisfaite :

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

avec  $n_D$  vecteur normal à  $\partial D$  sortant de  $D$ .

La pression étant constante par maille, l'équation de conservation de la masse au pas de temps précédent s'écrit sous la forme d'un bilan de type volumes finis :

$$\frac{|K|}{\delta t} (\rho^n - \rho^{n-1}) + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma^n \cdot n_\sigma = 0$$


 FIG. I.4 – Sous-volume de  $K$ .

où  $|\sigma| (\rho u)_\sigma^n \cdot n_\sigma$  est le flux de masse relatif à  $\sigma$ . En utilisant le lemme précédent, pour obtenir le débit massique que nous cherchons, nous allons construire sur chaque maille  $K$  un champ  $w$  de divergence constante et tel que :

$$\forall \sigma \in \mathcal{E}(K), \quad \int_\sigma w \cdot n_\sigma = |\sigma| (\rho u)_\sigma^n \cdot n_\sigma \quad (\text{I.48})$$

La vitesse associée à la face  $\varepsilon$  de la maille diamant,  $(\rho \tilde{u})|_\varepsilon$ , sera alors calculée en intégrant  $w$  sur  $\varepsilon$ , et un bilan de masse discret sera satisfait sur les demi-maillages diamant ; en sommant ces relations, on obtiendra la condition de compatibilité souhaitée. Pour l'élément fini de Crouzeix-Raviart, un tel champ  $w$  est donné par une interpolation directe des quantités  $((\rho u)_\sigma^n)_{\sigma \in \mathcal{E}(K)}$  :

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

où  $\varphi_\sigma$  est la fonction de forme de Crouzeix-Raviart associée au nœud de vitesse de  $\sigma$ . Pour l'élément fini de Rannacher-Turek, quand la maille est un rectangle ou un parallélépipède rectangle, il est obtenu par la formule d'interpolation suivante :

$$\tilde{\rho u}(x) = \sum_{\sigma \in \mathcal{E}(K)} \alpha_\sigma(x \cdot n_\sigma) [(\rho u)_\sigma^n \cdot n_\sigma] n_\sigma$$

où  $\alpha_\sigma(\cdot)$  sont des fonctions d'interpolation affines déterminées de façon à vérifier (I.48). Une extension à des maillages plus généraux est en cours.

Une discrétisation standard par éléments finis est utilisée pour le gradient de pression ainsi que pour le terme de viscosité :

$$\nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) \quad \rightarrow \quad - \sum_{K \in \mathcal{M}} \int_K p^n \nabla \cdot \varphi_\sigma^{(i)} + a_d(\tilde{u}^{n+1}, \varphi_\sigma^{(i)})$$

où la forme bilinéaire  $a_d(\cdot, \cdot)$  est définie comme précédemment.

Les pressions étant constantes par maille, la discrétisation par éléments finis du bilan de masse coïncide avec une formulation par volume finis ; comme précédemment, nous choisissons une discrétisation décentrée du terme advectif. Les mêmes techniques sont appliquées à la troisième équation

de cette étape. La positivité des masses volumiques est ainsi assurée et les lemmes I.4.6 et I.4.7 s'appliquent, si bien que la fraction massique de gaz reste bornée entre 0 et 1. Ce décentrement permet en outre de profiter de la stabilité induite par le travail des forces de pression (théorème I.4.3), qui fournit une estimation *a priori* essentielle pour prouver l'existence de la solution.

Dans la dernière étape, basée sur une discrétisation par volumes finis (avec en particulier une discrétisation monotone du terme  $\nabla \cdot (\rho y(1-y) u_r)$ ), sont traités les termes restants du bilan de masse de la phase gazeuse. La théorie présentée en section I.4.5 garantit l'existence et l'unicité d'une solution à cette étape.

Le schéma ainsi construit est donc conservatif, monotone et stable indépendamment du nombre de Mach (pour un écoulement incompressible, il dégénère en une méthode de projection incrémentale standard). Il présente, en outre, la propriété suivante : si les conditions initiales et aux limites le permettent, une interface est transportée à vitesse et à pression constantes (propriété de transport de la discontinuité de contact pour le système hyperbolique sous-jacent). Celle-ci semble essentielle, au vu des expériences numériques, pour la robustesse du schéma, en particulier dans le cas d'écoulements diphasiques avec de fortes différences de masse volumiques entre les deux phases. Sa preuve est basée sur la linéarité de la loi d'état par rapport à la masse volumique partielle de la phase gazeuse à pression constante.

Enfin, les théorèmes I.4.1 et I.4.3 nous permettent de démontrer la stabilité du schéma (moyennant l'ajout d'une étape de renormalisation de la pression que nous avons choisi de ne pas implémenter en pratique), au sens de la décroissance de l'entropie discrète, c'est à dire de l'analogue discret de la quantité suivante :

$$\frac{1}{2} \int_{\Omega} \rho(x, t) |u(x, t)|^2 dx + \int_{\Omega} z(x, t) f(\rho(x, t), z(x, t)) dx$$

Cette propriété est donnée par la relation suivante :

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^n}^2 + \int_{\Omega} z^{n+1} f_g(\rho^{n+1}, z^{n+1}) dx + \delta t \sum_{k=1}^{n+1} a_d(\tilde{u}^k, \tilde{u}^k) + \frac{\delta t^2}{2} |p^{n+1}|_{h, \rho^n}^2 \\ \leq \frac{1}{2} \|u^0\|_{h, \rho^0}^2 + \int_{\Omega} z^0 f_g(\rho^0, z^0) dx + \frac{\delta t^2}{2} |p^0|_{h, \rho^0}^2 \end{aligned}$$

Cette égalité est vérifiée en l'absence de drift ou lorsque la vitesse relative est donnée par la relation (I.9).

### Propriétés du schéma numérique

Les différentes propriétés énoncées ci-dessus sont récapitulées dans le théorème suivant :

**Théorème I.5.11 (Propriétés du schéma numérique)**

Soient  $(u^n)_{1 \leq n \leq N}$ ,  $(p^n)_{1 \leq n \leq N}$ ,  $(\rho^n)_{1 \leq n \leq N}$  et  $(y^n)_{1 \leq n \leq N}$  les solutions du schéma numérique considéré, avec une condition initiale strictement positive pour  $\rho$  et une condition initiale pour  $y$  à valeurs dans  $(0, 1]$ . On suppose que le terme visqueux est dissipatif (i.e.  $\forall v \in W_h$ ,  $a_d(v, v) \geq 0$ , ce qui est vérifié pour l'expression choisie lorsque la viscosité est constante). Soit  $f$  l'énergie libre définie par (I.23). Alors le schéma (I.42)-(I.44) vérifie les propriétés suivantes, pour tout  $n \leq N$  :

- (i) il existe une solution à chaque étape du schéma ;
- (ii) les inconnues restent dans leurs bornes physiques :

$$\rho_K^n > 0, \quad z_K^n > 0, \quad p_K^n > 0, \quad y_K^n \in (0, 1], \quad \forall K \in \mathcal{M}$$

- (iii) le schéma est conservatif :

$$\sum_{K \in \mathcal{M}} |K| \rho_K^n = \sum_{K \in \mathcal{M}} |K| \rho_K^0, \quad \sum_{K \in \mathcal{M}} |K| z_K^n = \sum_{K \in \mathcal{M}} |K| z_K^0,$$

- (iv) le schéma est stable (i.e. l'entropie discrète décroît) lorsque  $u_r = 0$  ou  $u_r$  est donnée par la relation (I.9) ;
- (v) une interface entre phases est transportée à vitesse et pression constantes : si  $u_K^n = u_0$  et  $p_K^n = p_0$  pour tout  $K \in \mathcal{M}$  et si les conditions aux limites le permettent, alors  $p_K^{n+1} = p_0$  et  $u_K^{n+1} = u_0$ , pour tout  $K \in \mathcal{M}$ .

Cette méthode numérique est testée sur un cas possédant par construction une solution analytique, afin d'en vérifier la convergence (chapitre IV). Les erreurs de la vitesse, de la pression et de la masse volumique de la phase gazeuse, respectivement en norme  $L^2$  et en norme  $L^2$  discrète en fonction du pas de temps, pour différents maillages, montrent d'abord une décroissance qui correspond à une convergence en temps d'ordre environ un pour ces trois variables, avant d'atteindre un plateau, dû au fait que les erreurs sont bornées inférieurement par l'erreur résiduelle de discrétisation en espace. Pour la vitesse, la pression et la masse volumique de la phase gazeuse, la valeur des erreurs sur le plateau montre une convergence en espace (en norme  $L^2$ ) d'ordre environ un, ce qui est cohérent avec le choix d'une discrétisation décentrée vers l'amont pour les termes de convection des bilans de masse du mélange et de la phase gazeuse.

**Remarque I.5.12** *En résolvant en séquence le bilan de masse de la phase gazeuse et les équations de Navier-Stokes, on pourrait obtenir un schéma plus facile à implémenter et moins coûteux en temps de calcul, tout en conservant les propriétés de monotonie énoncées ci-dessus. Mais un tel schéma ne peut vérifier la propriété de transport des interfaces à vitesse et pression constantes, ce qui engendre de possibles instabilités, en particulier en cas de fortes différences de masse volumiques entre les deux phases.*

## I.6 Résultats numériques

Le schéma numérique pour la résolution du modèle de dérive a fait l'objet de nombreuses expériences numériques, permettant de vérifier les propriétés de stabilité et monotonie démontrées sur le plan théorique, ainsi que sa consistance. Nous présentons ici un cas test de ballotement, pour lequel une solution analytique est donnée [12], et la simulation d'un écoulement dans une colonne à bulles couramment calculé dans la littérature [5].



### I.6.1 Un cas de ballotement

On étudie l'écoulement bidimensionnel de deux fluides non visqueux (ici l'eau et l'air) de masse volumique  $\rho_\ell$  et  $\rho_g$  dans un domaine dont la géométrie est décrite par la figure I.5. Les deux fluides sont soumis à la gravité  $g$ , le fluide le plus lourd étant placé en dessous. A l'état initial, les deux fluides sont au repos. A  $t = 0$ , on soumet le système à une accélération horizontale donnée par  $a_0 = 0.1 \text{ m.s}^{-2}$ . Le domaine de calcul a une largeur de  $L = 1 \text{ m}$  et une hauteur de  $2.25 \text{ m}$ . Les hauteurs d'eau et d'air sont respectivement  $h_\ell = 1 \text{ m}$  et  $h_g = 1.25 \text{ m}$ . Initialement, le système est au repos et la pression, uniforme, vaut  $p = 10^5 \text{ Pa}$ . La masse volumique de l'eau est fixée à  $\rho_\ell = 1000 \text{ kg.m}^{-3}$  et celle de l'air est donnée par  $\rho_g = p/RT$ , où  $RT$  est tel que  $\rho_g = 1.2 \text{ kg.m}^{-3}$  à la pression initiale. Le coefficient de diffusion  $D$  ainsi que la vitesse de dérive sont nuls.

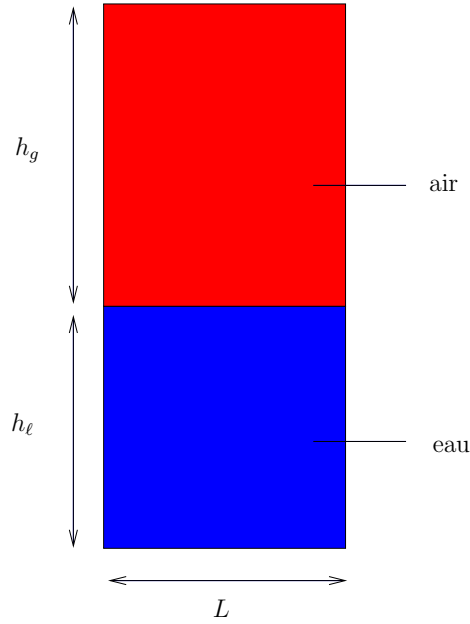


FIG. I.5 – Définition du problème de ballotement.

Il est démontré dans [12] que la solution analytique décrivant la surface libre entre les deux phases est donnée par :

$$\xi = \frac{a_0}{g} \left[ x - \frac{L}{2} + \sum_{n \geq 0} \frac{4}{L k_{2n+1}^2} \cos(\omega_{2n+1} t) \cos(k_{2n+1} t) \right]$$

où le nombre  $k_n$  est défini par :

$$k_n = \frac{2 \pi n}{L}$$

et  $\omega_n$  est donné par :

$$\omega_n^2 = \frac{g k_n (\rho_\ell - \rho_g)}{\rho_g \coth(k_n h_g) + \rho_\ell \coth(k_n h_\ell)}$$

Pour ce cas test, on utilise un maillage régulier composé de rectangles (éléments finis de Rannacher-Turek) avec environ 41 000 mailles. Dans la direction verticale, on raffine autour de l'interface entre les deux phases, obtenant ainsi des mailles ayant une hauteur pouvant varier de  $0.0005 \text{ m}$  près de l'interface, jusqu'à  $0.05 \text{ m}$  près des parois supérieure et inférieure du domaine.

Dans la direction horizontale, on choisit un maillage uniforme composé de 70 mailles. Les résultats numériques obtenus avec un pas de temps  $\delta t = 10^{-2} s$  et une viscosité  $\mu = \rho/1000$  sont reportés sur la figure I.6 et montrent que la solution numérique reste proche de la solution analytique.

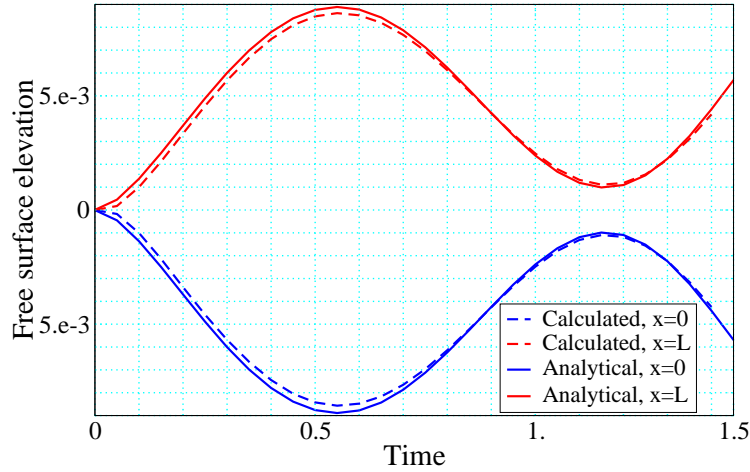


FIG. I.6 – Cas de ballottement : comparaison entre les solutions analytique et numérique.

### I.6.2 Colonne à bulles

La colonne a une section transversale rectangulaire de largeur  $L = 50 cm$ , de profondeur  $8 cm$  et de hauteur  $H = 200 cm$ . Elle est partiellement remplie d'eau jusqu'à une hauteur de  $h = 150 cm$ . Le gaz est introduit dans le système avec un débit  $q = 8 l/mn$  par un diffuseur, placé à  $15 cm$  de la paroi gauche. Ce dernier est circulaire, avec un diamètre de  $40 mm$ .

En entrée, on impose le taux de vide  $\alpha_{g,imp} = 1$  et la vitesse :

$$u_{imp} = \frac{q}{S \alpha_{g,imp}}$$

où  $S$  est la taille de la portion du maillage de la surface associée au diffuseur. Le long des autres parois, on utilise des conditions de Dirichlet homogènes pour la vitesse. Les conditions initiales sont  $u = 0 m.s^{-1}$  pour la vitesse et  $p = p_0$  pour la pression avec  $p_0 = 10^5 Pa$ , soit la pression atmosphérique. La masse volumique du liquide est fixée à  $\rho_\ell = 1000 kg.m^{-3}$  et celle du gaz est donnée par  $\rho_g = p/RT$  où  $RT$  est tel que  $\rho_g = 1.2 kg.m^{-3}$  à la pression initiale. Le coefficient de diffusion  $D$  est nul, la vitesse de dérive est constante et a pour valeur  $u_r = (0, 0.2)^t m.s^{-1}$ .

Pour ce cas test, on utilise un maillage régulier composé de rectangles (éléments finis de Rannacher-Turek), avec 76 mailles dans la direction horizontale parmi lesquelles 4 pour le diffuseur de gaz et 300 dans la direction verticale.

Les résultats numériques obtenus avec un pas de temps  $\delta t = 10^{-2} s$  et une viscosité  $\mu = 1 Pa.s$  sont reportés sur la figure I.7. On constate le bon comportement du schéma (stabilité de la surface libre, possibilité de traiter des zones purement monophasiques liquides), ainsi qu'un accord qualitatif avec l'expérience satisfaisant compte tenu de la simplicité de la modélisation physique (viscosité constante pour représenter un écoulement turbulent notamment).

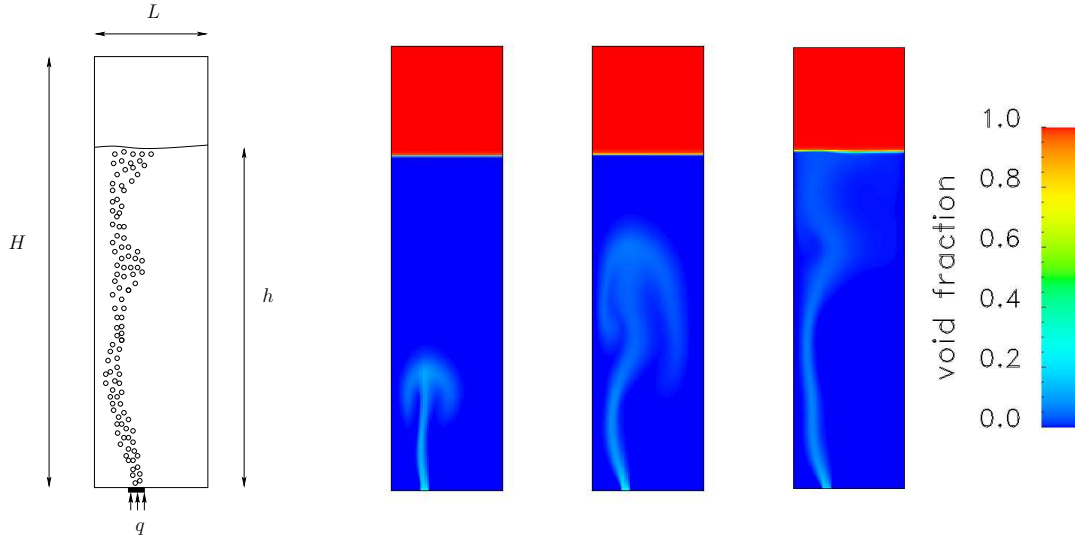


FIG. I.7 – Colonne à bulles de Becker : taux de vide à  $t = 2s$ ,  $t = 4s$  et  $t = 40s$ .

## I.7 Conclusion

Nous avons développé dans cette thèse des schémas de correction de pression pour les écoulements compressibles, ou, plus précisément, pour les équations de Navier-Stokes barotropes et le modèle de dérive.

Sur le plan de la discrétisation spatiale, ces schémas ont pour originalité de marier une discrétisation par éléments finis des équations de Navier-Stokes avec des techniques (voire des équations discrétisées par) volumes finis. Les espaces d'approximation utilisés sont des éléments finis de bas degré, non conformes (éléments finis de Crouzeix-Raviart ou Rannacher-Turek). Ils permettent une approximation naturelle des termes visqueux et sont intrinsèquement stables (*i.e.* ils vérifient la condition dite *inf-sup* ou condition de Babuska-Brezzi) dans la limite de l'incompressible. Ces algorithmes ont été développés dans le souci de satisfaire à trois propriétés de stabilité discrètes.

La première propriété est un analogue discret du théorème de l'énergie cinétique et donne la stabilité  $L^2$  de l'opérateur d'advection de vitesse (théorème I.4.1). Ce résultat a un caractère générique : il peut s'appliquer à tout autre opérateur d'advection discrétisé par volumes finis et peut être utilisé comme base pour démontrer des propriétés de stabilité pour d'autres schémas numériques (comme par exemple le schéma MAC). Dans notre cas, afin de pouvoir appliquer ce résultat, il est nécessaire de construire une approximation de type volumes fini de l'opérateur d'advection du bilan de quantité de mouvement basée sur un maillage dual (composé de mailles diamants).

Le deuxième résultat de type volumes finis permet de profiter de la stabilité induite par le travail des forces de pression (théorème I.4.2). Ce résultat peut s'appliquer ici car la pression est constante par maille. Grâce à cette propriété, les discrétisations par éléments finis du gradient de pression et du bilan de masse coïncident avec des approximations de type volumes finis. De plus, il est possible de choisir comme fonctions tests pour le bilan de masse des fonctions non-linéaires appliquées à l'inconnue. Par exemple, la fonction test  $[\rho f(\rho)]' = \log(\rho) + 1$  (dans le cas  $\varrho(p) = p$ ) appartient à l'espace d'approximation des pressions.

Le dernier résultat est un résultat de stabilité  $L^\infty$  pour une équation d'advection-diffusion non

linéaire (lemmes I.4.6 et I.4.7). Il est appliqué, dans notre cas, à l'équation de conservation de la phase gazeuse et permet de conserver l'inconnue dans ses bornes physiques, à savoir entre 0 et 1.

Si on considère une équation de conservation discrétisée par des volumes finis, du type :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K s_K - \rho_K^* s_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} s_\sigma + \dots [\text{éventuels termes de diffusion}] \dots = 0$$

le premier et le dernier de ces résultats ne sont valables que si est vérifiée la même condition de compatibilité entre les masses volumiques de début ( $\rho_K^*$ ) et fin de pas ( $\rho_K$ ) et les flux massiques  $F_{\sigma,K}$ , analogue au bilan de masse. La vérification de cette condition présente deux difficultés. La première tient au fait que dans un schéma à pas fractionnaires, le bilan de masse peut ne pas avoir été résolu au stade du schéma qui nous intéresse (*i.e.* la résolution de l'équation de transport de l'inconnue  $s$ ). Deux approches sont alors possibles : soit on ajoute une étape de prédiction de la masse volumique, soit on effectue un décalage en temps des masses volumiques, la condition de compatibilité étant, dans ce dernier cas, donnée par le bilan de masse au pas de temps précédent. Le premier choix permet un passage naturel à une discrétisation en temps d'ordre deux mais, contrairement au deuxième, ne garantit pas la conservativité du schéma. La deuxième difficulté est spécifique à l'étape de prédiction de vitesse. Elle est due au fait que l'opérateur d'advection pour la vitesse est discrétisé sur un maillage dual ; la condition de compatibilité devra ainsi être vérifiée sur ce même maillage alors que le bilan de masse au pas de temps précédent est résolu sur le maillage primal. Pour venir à bout de cette difficulté, nous proposons une construction de flux massiques qui garantit l'équivalence entre les bilans de masse sur les mailles primales et sur les mailles duales (lemme I.5.10).

Sur la base de ces trois résultats sur les volumes finis, nous construisons deux schémas de correction de pression, l'un pour la résolution des équations de Navier-Stokes compressibles, l'autre pour la résolution du modèle de dérive, tout deux stables indépendamment du nombre de Mach (ils dégénèrent en méthodes de projection incrémentales standard pour un écoulement incompressible). On démontre dans les deux cas la stabilité du schéma dans le sens de la décroissance de l'entropie discrète, concept à notre connaissance rarement associé aux méthodes de correction de pression. De plus, nous prouvons l'existence d'une solution à chacune des étapes. Pour des classes de problèmes génériques dans lesquelles entrent les étapes de de correction de pression de ces schémas, l'existence est prouvée par un argument de degré topologique. Ce résultat pourrait, moyennant une étude de convergence, être la base pour démontrer l'existence d'une solution au problème continu.

Afin de démontrer ces deux résultats dans le cas du modèle de dérive, une étape de correction de pression originale couplant les bilans de masse du mélange et de la phase gazeuse est construite. Ce couplage semble essentiel, au vu des expériences numériques, pour garantir la robustesse du schéma, en particulier dans le cas d'écoulements diphasiques avec de fortes différences de masse volumiques entre les deux phases.

Différentes idées, issues de ce travail, s'appliquent dans un cadre plus large que les écoulements diphasiques à phases dispersées, et de ce fait, sont d'ores et déjà utilisées quotidiennement dans le code de calcul ISIS dédié à la simulation des incendies et développé sur la base de la plateforme PELICANS au sein de l'IRSN : la discrétisation particulière de l'opérateur d'advection de vitesse, le décalage en temps des masse volumiques et l'approximation de type volumes finis des équations d'advection-diffusion non-linéaires qui garantit la conservation de l'inconnue dans ses bornes physiques.

Dans la suite, plusieurs développements sont envisageables afin d'améliorer les schéma numériques proposés dans cette thèse :

- le passage à une discrétisation spatiale du second ordre avec des techniques de type MUSCL ;

- l'utilisation de solveurs itératifs de type Krylov pour la résolution de l'étape de correction de pression : actuellement seuls les solveurs directs peuvent être utilisés, ce qui pose des problèmes de taille de mémoire ;
- le développement de techniques pour effectuer des calculs parallèles (un seul processeur est utilisé pour le moment).

Les algorithmes développés au cours de cette thèse sont spécifiques aux éléments finis de bas degré. L'extension à des éléments finis de plus haut degré semble pour le moment problématique, la difficulté étant d'utiliser, comme fonctions test pour le bilan de masse, la fonction obtenue en appliquant une fonction non-linéaire à l'inconnue.

De plus une étude de la convergence des deux schémas numériques proposés au cours de cette thèse est souhaitable. La convergence a d'ores et déjà été démontrée pour le problème de Stokes avec une loi d'état  $\varrho(p) = p$  dans [34]. L'étude du cas  $\varrho(p) = p^\gamma$ ,  $\gamma > 1$  est en cours. Cependant l'extension de ce résultat à Navier-Stokes et, plus ardu encore, au modèle de dérive semble difficile.

Enfin, plusieurs extensions du modèle physique sont envisagées : l'implémentation d'un modèle de turbulence, le passage au cas anisotherme et une meilleure simulation de l'interaction corium-béton en tenant compte des phénomènes d'ablation aux frontières du domaine. Ces extensions poseront des problèmes du point de vue numérique ; le dernier point, par exemple, suggère l'emploi d'une méthode "Arbitrary Lagrange Euler".



## Chapitre II

# An unconditionally stable pressure correction scheme for compressible barotropic Navier-Stokes equations

**Abstract.** We present in this paper a pressure correction scheme for barotropic compressible Navier-Stokes equations, which enjoys an unconditional stability property, in the sense that the energy and maximum-principle-based a priori estimates of the continuous problem also hold for the discrete solution. The stability proof is based on two independent results for general finite volume discretizations, both interesting for their own sake : the  $L^2$ -stability of the discrete advection operator provided it is consistent, in some sense, with the mass balance and the estimate of the pressure work by means of the time derivative of the elastic potential. The proposed scheme is built in order to match these theoretical results, and combines a fractional-step time discretization of pressure-correction type with a space discretization associating low order non-conforming mixed finite elements and finite volumes. Numerical tests with an exact smooth solution show the convergence of the scheme.

### II.1 Introduction

The problem addressed in this paper is the system of the so-called barotropic compressible Navier-Stokes equations, which reads :

$$\left\{ \begin{array}{l} \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \\ \frac{\partial}{\partial t}(\rho u) + \nabla \cdot (\rho u \otimes u) + \nabla p - \nabla \cdot \tau(u) = f_v \\ \rho = \varrho(p) \end{array} \right. \quad (\text{II.1})$$

where  $t$  stands for the time,  $\rho$ ,  $u$  and  $p$  are the density, velocity and pressure in the flow,  $f_v$  is a forcing term and  $\tau(u)$  stands for the shear stress tensor. The function  $\varrho(\cdot)$  is the equation of state used for the modelling of the particular flow at hand, which may be the actual equation of state of the fluid or may result from assumptions concerning the flow ; typically, laws as  $\varrho(p) = p^{1/\gamma}$ , where  $\gamma$  is a coefficient specific to the considered fluid, are obtained by making the assumption that the flow is isentropic. This system of equations is posed over  $\Omega \times (0, T)$ , where  $\Omega$  is a domain of  $\mathbb{R}^d$ ,  $d \leq 3$  supposed to be polygonal ( $d = 2$ ) or polyhedral ( $d = 3$ ), and the final time  $T$  is finite. It must be supplemented by boundary conditions and by an initial condition for  $\rho$  and  $u$ .

The development of pressure correction techniques for compressible Navier-Stokes equations may be traced back to the seminal work of Harlow and Amsden [41, 42] in the late sixties, who developed an iterative algorithm (the so-called ICE method) including an elliptic corrector step for the pressure. Later on, pressure correction equations appeared in numerical schemes proposed by several researchers, essentially in the finite-volume framework, using either a collocated [59, 20, 49, 60, 46, 54] or a staggered arrangement [11, 44, 45, 48, 7, 16, 69, 73, 74, 70, 72] of unknowns; in the first case, some corrective actions are to be foreseen to avoid the usual odd-even decoupling of the pressure in the low Mach number regime. Some of these algorithms are essentially implicit, since the final stage of a time step involves the unknown at the end-of-step time level; the end-of-step solution is then obtained by SIMPLE-like iterative processes [71, 48, 20, 49, 60, 46, 54]. The other schemes [44, 45, 59, 7, 16, 75, 69, 74, 70, 72] are predictor-corrector methods, where basically two steps are performed sequentially : first a semi-explicit decoupled prediction of the momentum or velocity (and possibly energy, for non-barotropic flows) and, second, a correction step where the end-of step pressure is evaluated and the momentum and velocity are corrected, as in projection methods for incompressible flows (see [14, 68] for the original papers, [53] for a comprehensive introduction and [38] for a review of most variants). The Characteristic-Based Split (CBS) scheme (see [56] for a recent review or [76] for the seminal paper), developed in the finite-element context, belongs to this latter class of methods.

Our aim in this paper is to propose and study a non-iterative pressure correction scheme for the solution of (II.1). In addition, this method is designed so as to be stable in the low Mach number limit, since our final goal is to apply it to simulate through a drift-flux approach a class of bubbly flows encountered in nuclear safety studies, where pure liquid (incompressible) and pure gaseous (compressible) zones may coexist. To this purpose, we use a low order mixed finite element approximation, which meets the two following requirements : to allow a natural discretization of the viscous terms and to provide a spatial discretization that is intrinsically stable (*i.e.* without the adjunction of stabilization terms to circumvent the so-called *inf-sup* or BB condition) in the incompressible limit.

In this work, a special attention is payed to stability issues. To be more specific, let us recall the *a priori* estimates associated to problem (II.1) with a zero forcing term, *i.e.* estimates which should be satisfied by any possible regular solution [52, 28, 57] :

$$\left| \begin{array}{ll}
 (i) & \rho(x, t) > 0, & \forall x \in \Omega, \forall t \in (0, T) \\
 (ii) & \int_{\Omega} \rho(x, t) \, dx = \int_{\Omega} \rho(x, 0) \, dx, & \forall t \in (0, T) \\
 (iii) & \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho(x, t) u(x, t)^2 \, dx + \frac{d}{dt} \int_{\Omega} \rho(x, t) P(\rho(x, t)) \, dx \\
 & \quad + \int_{\Omega} \tau(u(x, t)) : \nabla u(x, t) \, dx = 0, & \forall t \in (0, T)
 \end{array} \right. \tag{II.2}$$

In the latter relation,  $P(\cdot)$ , the "elastic potential", is a function derived from the equation of state, which satisfies :

$$P'(z) = \frac{\wp(z)}{z^2} \tag{II.3}$$

where  $\wp(\cdot)$  is the inverse function of  $\varrho(\cdot)$ , *i.e.* the function giving the pressure as a function of the density. The usual choice for  $P(\cdot)$  is, provided that this expression makes sense :

$$P(z) = \int_0^z \frac{\wp(s)}{s^2} \, ds \tag{II.4}$$



For these estimates to hold, the condition (II.2)-(i) must be satisfied by the initial condition ; note that a non-zero forcing term  $f_v$  in the momentum balance would make an additional term appear at the right hand side of relation (II.2)-(iii). This latter estimate is obtained from the second relation of (II.1) (*i.e.* the momentum balance) by taking the inner product by  $u$  and integrating over  $\Omega$ . This computation then involves two main arguments which read :

(i) Stability of the advection operator :

$$\int_{\Omega} \left[ \frac{\partial}{\partial t}(\rho u) + \nabla \cdot (\rho u \otimes u) \right] \cdot u \, dx = \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho |u|^2 \, dx \quad (\text{II.5})$$

(ii) Stability due to the pressure work :

$$\int_{\Omega} -p \nabla \cdot u \, dx = \frac{d}{dt} \int_{\Omega} \rho(x, t) P(\rho(x, t)) \, dx$$

Note that the derivation of both relations rely on the mass balance equation.

This paper is organized as follows.

We first derive a bound similar to (II.5)-(ii) for a given class of spatial discretizations ; the latter are introduced in section II.2.1 and the desired stability estimate (theorem II.2.1) is stated and proven in section II.2.2. We then show that this result allows to prove the existence of a solution for a fairly general class of discrete compressible flow problems. Section II.2 gathers this whole study, and constitutes the first part of this paper.

In a second part (section II.3), we turn to the derivation of a pressure correction scheme, the solution of which satisfies a discrete equivalent of the whole set of *a priori* estimates (II.2). To this purpose, besides theorem II.2.1, we need as a second key ingredient a discrete version of the bound (II.5)-(i) relative to the stability of the advection operator, which is stated and proven in section II.3.2 (theorem II.3.7). We then derive a fully discrete algorithm which is designed to meet the assumptions of these theoretical results, and establish its stability. Moreover, numerical experiments show that, for smooth solutions, this scheme converges as expected, namely with first order in time convergence for all the variables and first-to-second order in space in  $L^2$  and discrete  $L^2$  norm for the velocity and the pressure, respectively.

## II.2 Analysis of a class of discrete problems

The class of problems addressed in this section can be seen as the class of discrete systems obtained by space discretization by low-order non-conforming finite elements of continuous problems of the following form :

$$\left\{ \begin{array}{ll} A u + \nabla p = f_v & \text{in } \Omega \\ \frac{\varrho(p) - \rho^*}{\delta t} + \nabla \cdot (\varrho(p) u) = 0 & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{array} \right. \quad (\text{II.6})$$

where  $A$  stands for an abstract elliptic operator and the forcing term  $f_v$  and the density field  $\rho^*$  are known quantities. The unknowns of the problem are the velocity  $u$  and the pressure  $p$  ; the function  $\varrho(\cdot)$  stands for the equation of state. The domain  $\Omega$  is a polygonal ( $d = 2$ ) or polyhedral ( $d = 3$ ) open, bounded and connected subset of  $\mathbb{R}^d$ . Of course, at the continuous level, this statement of the problem should be completed by a precise definition of the functional spaces in which the velocity and the pressure are sought, together with regularity assumptions on the data. This is out of the scope here, since system (II.6) is only given to fix ideas ; indeed, the aim here is to prove some mathematical properties of the discrete problem, namely to establish some *a priori* estimates for

its solution and to prove that this nonlinear problem admits solutions for fairly general equations of state.

This section is organized as follows. We begin by describing the considered discretization and precisely stating the discrete problem at hand. Then we prove, for the chosen particular discretization, a fundamental result which is a discrete analogue of the elastic potential identity (II.5)-(ii). The next section is devoted to the proof of the existence of a solution, and we finally conclude by giving some practical examples of application of the abstract theory developed here.

### II.2.1 The discrete problem

Let  $\mathcal{M}$  be a decomposition of the domain  $\Omega$  either into convex quadrilaterals ( $d = 2$ ) or hexahedra ( $d = 3$ ) or in simplices. By  $\mathcal{E}$  and  $\mathcal{E}(K)$  we denote the set of all  $(d - 1)$ -edges  $\sigma$  of the mesh and of the element  $K \in \mathcal{M}$  respectively. The set of  $(d - 1)$ -edges included in the boundary of  $\Omega$  is denoted by  $\mathcal{E}_{\text{ext}}$  and the set of internal ones (*i.e.*  $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ ) is denoted by  $\mathcal{E}_{\text{int}}$ . The decomposition  $\mathcal{M}$  is supposed to be regular in the usual sense of the finite element literature (e.g. [15]), and, in particular,  $\mathcal{M}$  satisfies the following properties :  $\Omega = \bigcup_{K \in \mathcal{M}} \bar{K}$ ; if  $K, L \in \mathcal{M}$ , then  $\bar{K} \cap \bar{L}$  is reduced to the empty set, to a vertex or (if  $d = 3$ ) to a segment, or  $\bar{K} \cap \bar{L}$  is (the closure of) a common  $(d - 1)$ -edge of  $K$  and  $L$ , which is denoted by  $K|L$ . For each internal edge of the mesh  $\sigma = K|L$ ,  $n_{KL}$  stands for the normal vector of  $\sigma$ , oriented from  $K$  to  $L$ . By  $|K|$  and  $|\sigma|$  we denote the measure, respectively, of  $K$  and of the edge  $\sigma$ .

The space discretization relies either on the so-called "rotated bilinear element"/ $P_0$  introduced by Rannacher and Turek [61] for quadrilateral or hexahedric meshes, or on the Crouzeix-Raviart element (see [18] for the seminal paper and, for instance, [?, pp. 199–201] for a synthetic presentation) for simplicial meshes. The reference element  $\hat{K}$  for the rotated bilinear element is the unit  $d$ -cube (with edges parallel to the coordinate axes); the discrete functional space on  $\hat{K}$  is  $\tilde{Q}_1(\hat{K})^d$ , where  $\tilde{Q}_1(\hat{K})$  is defined as follows :

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

The reference element for the Crouzeix-Raviart is the unit  $d$ -simplex and the discrete functional space is the space  $P_1$  of affine polynomials. For both velocity elements used here, the degrees of freedom are determined by the following set of nodal functionals :

$$\{F_{\sigma,i}, \sigma \in \mathcal{E}(K), i = 1, \dots, d\}, \quad F_{\sigma,i}(v) = |\sigma|^{-1} \int_{\sigma} v_i d\gamma \quad (\text{II.7})$$

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

$$\begin{aligned} W_h = \{ & v_h \in L^2(\Omega) : v_h|_K \in W(K)^d, \forall K \in \mathcal{M}; \\ & F_{\sigma,i}(v_h) \text{ continuous across each edge } \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d; \\ & F_{\sigma,i}(v_h) = 0, \forall \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq i \leq d \} \end{aligned}$$

where  $W(K)$  is the space of functions on  $K$  generated by  $\tilde{Q}_1(\hat{K})$  through the  $Q_1$  mapping from  $\hat{K}$  to  $K$  for the Rannacher-Turek element and the space of affine functions on  $K$  for the Crouzeix-Raviart element. For both Rannacher-Turek and Crouzeix-Raviart discretizations, the pressure is approximated by the space  $L_h$  of piecewise constant functions :

$$L_h = \{q_h \in L^2(\Omega) : q_h|_K = \text{constant}, \forall K \in \mathcal{M}\}$$

Since only the continuity of the integral over each edge of the mesh is imposed, the velocities are 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  $c_{is} > 0$  possibly depending on the regularity of the shape of the cells but not on their size such that :

$$\forall p \in L_h, \quad \sup_{v \in W_h} \frac{\int_{\Omega, h} p \nabla \cdot v \, dx}{\|v\|_{1, b}} \geq c_{is} \|p - \bar{p}\|_{L^2(\Omega)}$$

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

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

From the definition (II.7), each velocity degree of freedom can be univoquely associated to an element edge. Therefore we shall use hereafter, somewhat improperly, the expression "velocity on the edge  $\sigma$ " to name the velocity vector defined by the degrees of freedom of the velocity components associated to  $\sigma$ . In addition, the velocity degrees of freedom are indexed by the number of the component and the associated edge, thus 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  $i^{\text{th}}$  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 considered finite elements, reads :

$$\varphi_\sigma^{(i)} = \varphi_\sigma e_i$$

where  $\varphi_\sigma$  is a scalar function. Similarly, each degree of freedom for the pressure is associated to a mesh  $K$ , and the set of pressure degrees of freedom is denoted by  $\{p_K, K \in \mathcal{M}\}$ .

For any  $K \in \mathcal{M}$ , let  $\rho_K^*$  be a quantity approximating a known density  $\rho^*$  on  $K$ . The family of real numbers  $(\rho_K^*)_{K \in \mathcal{M}}$  is supposed to be positive. The discrete problem considered in this section reads :

$$\left\{ \begin{array}{ll} a(u, \varphi_\sigma^{(i)}) - \int_{\Omega, h} p \nabla \cdot \varphi_\sigma^{(i)} = \int_\Omega f_v \cdot \varphi_\sigma^{(i)} \, dx, & \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d \\ \frac{|K|}{\delta t} (\varrho(p_K) - \rho_K^*) + \sum_{\sigma=K|L} v_{\sigma, K}^+ \varrho(p_K) - v_{\sigma, K}^- \varrho(p_L) = 0, & \forall K \in \mathcal{M} \end{array} \right. \quad (\text{II.8})$$

where  $v_{\sigma, K}^+$  and  $v_{\sigma, K}^-$  stand respectively for  $v_{\sigma, K}^+ = \max(v_{\sigma, K}, 0)$  and  $v_{\sigma, K}^- = -\min(v_{\sigma, K}, 0)$  with  $v_{\sigma, K} = |\sigma| u_\sigma \cdot n_{KL} = v_{\sigma, K}^+ - v_{\sigma, K}^-$ . The first equation is the standard finite element discretization of the first equation of (II.6), provided that the bilinear form  $a(\cdot, \cdot)$  is related to the operator  $A$  by a relation of the form :

$$a(v, w) = \int_\Omega Av \cdot w \, dx$$

where  $v$  and  $w$  are regular functions vanishing on the boundary (while this identity generally does not hold for functions of  $W_h$ ). Since the pressure is piecewise constant, the finite element discretization of the second relation of (II.6), *i.e.* the mass balance, is similar to a finite volume

formulation, in which we introduce the standard first-order upwinding. The bilinear form  $a(\cdot, \cdot)$  is supposed to be elliptic on  $W_h$ , *i.e.* to be such that the following property holds :

$$\exists c_a > 0 \text{ such that, } \forall v \in W_h, \quad a(v, v) \geq c_a \|v\|_*^2$$

where  $\|\cdot\|_*$  is a norm over  $W_h$ . We denote by  $\|\cdot\|^*$  its dual norm with respect to the  $L^2(\Omega)^d$  inner product, defined by :

$$\forall v \in W_h, \quad \|v\|^* = \sup_{w \in W_h} \frac{\int_{\Omega} v \cdot w \, dx}{\|w\|_*}$$

## II.2.2 On the pressure control induced by the pressure forces work

The aim of this subsection is to prove that the discretization at hand satisfies a stability bound which can be seen as the discrete analogue of equation (II.5)-(ii), which we recall here :

$$- \int_{\Omega} p \nabla \cdot u \, dx = \frac{d}{dt} \int_{\Omega} \rho P(\rho) \, dx, \quad \text{where } P(\cdot) \text{ is such that } P'(z) = \frac{\wp(z)}{z^2}$$

The formal computation which allows to derive this estimate in the continuous setting is the following. The starting point is the mass balance, which is multiplied by the derivative of  $z \mapsto z P(z)$  taken at  $\rho$ , denoted by  $[\rho P(\rho)]'$  :

$$[\rho P(\rho)]' \left[ \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) \right] = 0$$

This relation yields :

$$\frac{\partial [\rho P(\rho)]}{\partial t} + [\rho P(\rho)]' [u \cdot \nabla \rho + \rho \nabla \cdot u] = 0 \quad (\text{II.9})$$

And thus :

$$\frac{\partial [\rho P(\rho)]}{\partial t} + u \cdot \nabla [\rho P(\rho)] + [\rho P(\rho)]' \rho \nabla \cdot u = 0$$

Developping the derivative, we get :

$$\frac{\partial [\rho P(\rho)]}{\partial t} + \underbrace{u \cdot \nabla [\rho P(\rho)] + \rho P(\rho) \nabla \cdot u}_{\nabla \cdot (\rho P(\rho) u)} + \underbrace{\rho^2 [P(\rho)]' \nabla \cdot u}_{p \nabla \cdot u} = 0 \quad (\text{II.10})$$

and the result follows by integration in space, thanks to the fact that the velocity vanishes at the boundary. We are going here to reproduce this computation at the discrete level.

### Theorem II.2.1 (Stability due to the pressure work)

*Let us suppose that the equation of state  $\varrho(\cdot)$  is defined over  $[0, +\infty)$ . Let  $P(\cdot)$  be an elastic potential (*i.e.* a function satisfying (II.3)) such that the function  $f : (0, +\infty) \rightarrow \mathbb{R}$  defined by  $f(z) = z P(z)$  is once continuously differentiable and strictly convex. Let  $(p_K)_{K \in \mathcal{M}}$  satisfy the second relation of (II.8). For any  $K \in \mathcal{M}$ , we suppose that  $p_K > 0$  and we define  $\rho_K$  by  $\rho_K = \varrho(p_K)$ ; we also recall that, by assumption,  $\rho_K^* > 0$ . Then the following estimate holds :*

$$- \int_{\Omega, h} p \nabla \cdot u \, dx = \sum_{K \in \mathcal{M}} -p_K \sum_{\sigma=K|L} v_{\sigma, K} \geq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [\rho_K P(\rho_K) - \rho_K^* P(\rho_K^*)] \quad (\text{II.11})$$

**Proof.**

Let us write the divergence term in the discrete mass balance over  $K$  (*i.e.* the second relation of (II.8)) under the following form :

$$\sum_{\sigma=K|L} \rho_{\sigma} v_{\sigma,K}$$

where  $\rho_{\sigma}$  is either  $\rho_K$  if  $v_{\sigma,K} \geq 0$  or  $\rho_L$  if  $v_{\sigma,K} \leq 0$ . Multiplying this term by  $f'(\rho_K)$ , we obtain :

$$T_{\text{div},K} = f'(\rho_K) \sum_{\sigma=K|L} \rho_{\sigma} v_{\sigma,K} = f'(\rho_K) \left[ \sum_{\sigma=K|L} (\rho_{\sigma} - \rho_K) v_{\sigma,K} + \rho_K \sum_{\sigma=K|L} v_{\sigma,K} \right]$$

This latter form of  $T_{\text{div},K}$  may be compared to equation (II.9) : up to the multiplication by  $1/|K|$ , the first summation in the right hand side is the analogue of  $u \cdot \nabla \rho$  and the second one to  $\rho \nabla \cdot u$ . Developing the derivative of  $f(\cdot)$ , we then obtain a discrete analogue of the corresponding terms in relation (II.10) :

$$T_{\text{div},K} = f'(\rho_K) \sum_{\sigma=K|L} (\rho_{\sigma} - \rho_K) v_{\sigma,K} + \rho_K P(\rho_K) \sum_{\sigma=K|L} v_{\sigma,K} + \rho_K^2 P'(\rho_K) \sum_{\sigma=K|L} v_{\sigma,K} \quad (\text{II.12})$$

By definition (II.3) of  $P(\cdot)$ , the last term is equal to  $\rho_K \sum_{\sigma=K|L} v_{\sigma,K}$ . The process will be completed if we put the first two terms in divergence form. To this end, let us sum up the quantities  $T_{\text{div},K}$  over  $K \in \mathcal{M}$  and reorder the summation :

$$\sum_{K \in \mathcal{M}} T_{\text{div},K} = \sum_{K \in \mathcal{M}} \rho_K \sum_{\sigma=K|L} v_{\sigma,K} + \sum_{\sigma \in \mathcal{E}_{\text{int}}} T_{\text{div},\sigma} \quad (\text{II.13})$$

where, if  $\sigma = K|L$  :

$$T_{\text{div},\sigma} = v_{\sigma,K} \left[ f(\rho_K) + f'(\rho_K)(\rho_{\sigma} - \rho_K) - f(\rho_L) - f'(\rho_L)(\rho_{\sigma} - \rho_L) \right]$$

In this relation, there are two possible choices for the orientation of  $\sigma$ , *i.e.*  $K|L$  or  $L|K$  ; we suppose that the chosen orientation is such that  $v_{\sigma,K} \geq 0$ . Let  $\bar{\rho}_{\sigma}$  be defined by :

$$\left| \begin{array}{ll} \text{if } \rho_K \neq \rho_L : & f(\rho_K) + f'(\rho_K)(\bar{\rho}_{\sigma} - \rho_K) = f(\rho_L) + f'(\rho_L)(\bar{\rho}_{\sigma} - \rho_L) \\ \text{otherwise :} & \bar{\rho}_{\sigma} = \rho_K = \rho_L \end{array} \right. \quad (\text{II.14})$$

As the function  $f(\cdot)$  is supposed to be once continuously differentiable and strictly convex, the technical lemma II.2.3 proven hereafter applies and  $\bar{\rho}_{\sigma}$  is uniquely defined and satisfies  $\bar{\rho}_{\sigma} \in [\min(\rho_K, \rho_L), \max(\rho_K, \rho_L)]$ . By definition, the choice  $\rho_{\sigma} = \bar{\rho}_{\sigma}$  is such that the term  $T_{\text{div},\sigma}$  vanishes, which means that, in this case, we would indeed have obtained that the first two terms of equation (II.12) are a conservative approximation of the quantity  $\nabla \cdot \rho P(\rho) u$  appearing in equation (II.10), with the following expression for the flux :

$$F_{\sigma,K} = [f(\rho)]_{\sigma} v_{\sigma,K}, \quad \text{with :} \quad [f(\rho)]_{\sigma} = f(\rho_K) + f'(\rho_K)(\bar{\rho}_{\sigma} - \rho_K) = f(\rho_L) + f'(\rho_L)(\bar{\rho}_{\sigma} - \rho_L)$$

For any choice of  $\rho_{\sigma}$ , we have :

$$T_{\text{div},\sigma} = v_{\sigma,K} (\rho_{\sigma} - \bar{\rho}_{\sigma}) (f'(\rho_K) - f'(\rho_L))$$

With the orientation taken for  $\sigma$ , an upwind choice yields :

$$T_{\text{div},\sigma} = v_{\sigma,K} (\rho_K - \bar{\rho}_{\sigma}) (f'(\rho_K) - f'(\rho_L))$$

and, using the fact that  $f'(\cdot)$  is an increasing function since  $f(\cdot)$  is convex and that  $\min(\rho_K, \rho_L) \leq \bar{\rho}_\sigma \leq \max(\rho_K, \rho_L)$ , it is easily seen that  $T_{\text{div},\sigma}$  is non-negative.

Multiplying by  $f'(\rho_K)$  the mass balance over each cell  $K$  and summing for  $K \in \mathcal{M}$  thus yields, invoking equation (II.13) :

$$-\sum_{K \in \mathcal{M}} p_K \sum_{\sigma=K|L} v_{\sigma,K} = R + \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} f'(\rho_K) (\rho_K - \rho_K^*) \quad (\text{II.15})$$

where  $R$  is non-negative, and the result follows invoking once again the convexity of  $f(\cdot)$ . ■

**Remark II.2.2 (On a non-dissipative scheme)** *The preceding proof shows that, for a scheme to conserve the energy (i.e. to obtain a discrete equivalent of (II.2)-(iii)), besides other arguments, the choice of  $\bar{\rho}_\sigma$  given by (II.14) for the density at the edge of the control volume in the discretization of the flux in the mass balance seems to be mandatory; any other choice leads to an artificial (i.e. due to the numerical scheme) dissipation or production in the work of the pressure forces. Note however that, this discretization being essentially of centered type, the positivity of the density is not warranted in this case.*

In the course of the preceding proof, we used the following technical lemma.

**Lemma II.2.3**

Let  $g(\cdot)$  be a strictly convex and once continuously derivable real function over an open interval  $I \subset \mathbb{R}$ . Let  $\rho_1 \in I$  and  $\rho_2 \in I$  be two distinct real numbers. Then the following relation :

$$g(\rho_1) + g'(\rho_1)(\bar{\rho} - \rho_1) = g(\rho_2) + g'(\rho_2)(\bar{\rho} - \rho_2) \quad (\text{II.16})$$

uniquely defines the real number  $\bar{\rho}$ . In addition, we have  $\bar{\rho} \in [\min(\rho_1, \rho_2), \max(\rho_1, \rho_2)]$ .

**Proof.**

Without loss of generality, let us suppose that  $\rho_1 < \rho_2$ . Reordering equation (II.16), we get :

$$g(\rho_1) + g'(\rho_1)(\rho_2 - \rho_1) - g(\rho_2) = (\bar{\rho} - \rho_2) [g'(\rho_2) - g'(\rho_1)]$$

Because  $g(\cdot)$  is strictly convex,  $g'(\rho_2) - g'(\rho_1)$  does not vanish, and therefore the latter equation proves that  $\bar{\rho}$  is uniquely defined. In addition, for the same reason, the left hand side of this relation is negative and  $g'(\rho_2) - g'(\rho_1)$  is positive, thus we have  $\bar{\rho} < \rho_2$ . Reordering equation (II.16) yields :

$$g(\rho_2) + g'(\rho_2)(\rho_1 - \rho_2) - g(\rho_1) = (\bar{\rho} - \rho_1) [g'(\rho_1) - g'(\rho_2)]$$

which, considering the signs of the left hand side and of  $g'(\rho_1) - g'(\rho_2)$ , implies  $\bar{\rho} > \rho_1$ . ■

**II.2.3 Existence of a solution**

The aim of this section is to prove the existence of a solution to the discrete problem (II.8). It follows from a topological degree argument, linking by a homotopy the problem at hand to a linear system.

This section begins with a lemma which is used to obtain a positive lower bound for the pressure in the sequel.

**Lemma II.2.4**

Let us consider the following problem :

$$\forall K \in \mathcal{M}, \quad |K| \frac{\varphi_1(p_K) - \varphi_1(p_K^*)}{\delta t} + \sum_{\sigma=K|L} v_{\sigma,K}^+ \varphi_2(p_K) - v_{\sigma,K}^- \varphi_2(p_L) = 0 \quad (\text{II.17})$$

where  $\varphi_1(\cdot)$  is an increasing function and  $\varphi_2(\cdot)$  is a non-decreasing and non-negative function. Suppose that there exists  $\bar{p}$  such that :

$$\varphi_1(\bar{p}) + \delta t \varphi_2(\bar{p}) \max \left[ 0, \max_{K \in \mathcal{M}} \left( \frac{1}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right) \right] = \min_{K \in \mathcal{M}} [\varphi_1(p_K^*)] \quad (\text{II.18})$$

Then,  $\forall K \in \mathcal{M}$ ,  $p_K$  satisfies  $p_K \geq \bar{p}$ .

**Proof.**

Let us assume that there exists a cell  $\bar{K}$  such that  $p_{\bar{K}} = \min_{K \in \mathcal{M}} (p_K) < \bar{p}$ . Multiplying by  $\delta t/|\bar{K}|$  the relation (II.17) written for  $K = \bar{K}$ , we get :

$$\varphi_1(p_{\bar{K}}) + \frac{\delta t}{|\bar{K}|} \sum_{\sigma=\bar{K}|L} \left[ v_{\sigma,\bar{K}}^+ \varphi_2(p_{\bar{K}}) - v_{\sigma,\bar{K}}^- \varphi_2(p_L) \right] = \varphi_1(p_{\bar{K}}^*) \quad (\text{II.19})$$

Then, subtracting equation (II.18), we have :

$$\begin{aligned} & \varphi_1(p_{\bar{K}}) - \varphi_1(\bar{p}) + \frac{\delta t}{|\bar{K}|} \sum_{\sigma=\bar{K}|L} \left[ v_{\sigma,\bar{K}}^+ \varphi_2(p_{\bar{K}}) - v_{\sigma,\bar{K}}^- \varphi_2(p_L) \right] \\ & - \delta t \varphi_2(\bar{p}) \max \left[ 0, \max_{K \in \mathcal{M}} \left( \frac{1}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right) \right] = \varphi_1(p_{\bar{K}}^*) - \min_{K \in \mathcal{M}} [\varphi_1(p_K^*)] \geq 0 \end{aligned}$$

The previous relation can be written as  $T_1 + T_2 + T_3 \geq 0$  with :

$$\begin{aligned} T_1 &= \varphi_1(p_{\bar{K}}) - \varphi_1(\bar{p}) \\ T_2 &= \delta t \varphi_2(p_{\bar{K}}) \left[ \frac{1}{|\bar{K}|} \sum_{\sigma=\bar{K}|L} v_{\sigma,\bar{K}} \right] - \delta t \varphi_2(\bar{p}) \max \left[ 0, \max_{K \in \mathcal{M}} \left( \frac{1}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right) \right] \\ T_3 &= \frac{\delta t}{|\bar{K}|} \sum_{\sigma=\bar{K}|L} v_{\sigma,\bar{K}}^- (\varphi_2(p_{\bar{K}}) - \varphi_2(p_L)) \end{aligned}$$

Since  $\varphi_1(\cdot)$  is an increasing function and, by assumption,  $p_{\bar{K}} < \bar{p}$ , we have  $T_1 < 0$ . Similarly,  $0 \leq \varphi_2(p_{\bar{K}}) \leq \varphi_2(\bar{p})$  and the discrete divergence over  $\bar{K}$  (i.e.  $1/|\bar{K}| \sum_{\sigma=\bar{K}|L} v_{\sigma,\bar{K}}$ ) is necessarily smaller than the maximum of this quantity over the cells of the mesh, thus  $T_2 \leq 0$ . Finally, since, by assumption,  $p_{\bar{K}} \leq p_L$  for any neighbouring cell  $L$  of  $\bar{K}$ ,  $\varphi_2(\cdot)$  is a non-decreasing function and  $v_{\sigma,K}^- \geq 0$ ,  $T_3 \leq 0$ . We thus obtain a contradiction with the fact that  $T_1 + T_2 + T_3 \geq 0$ , which proves that  $p_K \geq \bar{p}$ ,  $\forall K \in \mathcal{M}$ .  $\blacksquare$

We now state the abstract theorem which will be used hereafter; this result follows from standard arguments of the topological degree theory (see [19] for an exposition of the theory and

[25] for another utilisation for the same objective as here, namely the proof of existence of a solution to a numerical scheme).

**Theorem II.2.5 (A result from the topological degree theory)**

Let  $N$  and  $M$  be two positive integers and  $V$  be defined as follows :

$$V = \{(x, y) \in \mathbb{R}^N \times \mathbb{R}^M \text{ such that } y > 0\}$$

where, for any real number  $c$ , the notation  $y > c$  means that each component of  $y$  is greater than  $c$ . Let  $b \in \mathbb{R}^N \times \mathbb{R}^M$  and  $f(\cdot)$  and  $F(\cdot, \cdot)$  be two continuous functions respectively from  $V$  and  $V \times [0, 1]$  to  $\mathbb{R}^N \times \mathbb{R}^M$  satisfying :

(i)  $F(\cdot, 1) = f(\cdot)$  ;

(ii)  $\forall \alpha \in [0, 1]$ , if  $v \in V$  is such that  $F(v, \alpha) = b$  then  $v \in W$ , where  $W$  is defined as follows :

$$W = \{(x, y) \in \mathbb{R}^N \times \mathbb{R}^M \text{ s.t. } \|x\| < C_1 \text{ and } \epsilon < y < C_2\}$$

with  $C_1, C_2$  and  $\epsilon$  three positive constants and  $\|\cdot\|$  a norm defined over  $\mathbb{R}^N$  ;

(iii) the topological degree of  $F(\cdot, 0)$  with respect to  $b$  and  $W$  is equal to  $d_0 \neq 0$ .

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

We are now in position to prove the existence of a solution to the discrete problem (II.8), for fairly general equations of states.

**Theorem II.2.6 (Existence of a solution)**

Let us suppose that the equation of state  $\varrho(\cdot)$  is such that :

1.  $\varrho(\cdot)$  is defined and increasing over  $[0, +\infty)$ ,  $\varrho(0) = 0$  and  $\lim_{z \rightarrow +\infty} \varrho(z) = +\infty$ ,
2. there exists an elastic potential  $P(\cdot)$  (i.e. a function satisfying (II.3)) such that the function  $f : (0, +\infty) \rightarrow \mathbb{R}$  defined by  $f(z) = zP(z)$  is once continuously differentiable, strictly convex and  $f(z) \geq -C_P$ ,  $\forall z \in (0, +\infty)$ , where  $C_P$  is a non-negative constant.

In addition, we recall that, in the discrete problem at hand (II.8),  $\rho_K^* > 0$ ,  $\forall K \in \mathcal{M}$ .

Then problem (II.8) admits at least one solution  $(u_\sigma, p_K)_{\sigma \in \mathcal{E}_{\text{int}}, K \in \mathcal{M}}$ , and any solution is such that,  $\forall K \in \mathcal{M}$ ,  $p_K$  is positive.

**Proof.**

Let  $N = d \text{ card}(\mathcal{E}_{\text{int}})$  and  $M = \text{card}(\mathcal{M})$ . We identify the space of discrete velocities  $W_h$  and pressures  $L_h$  to  $\mathbb{R}^N$  and  $\mathbb{R}^M$  respectively and, keeping the same notation for the discrete functions and the associated vectors of degree of freedom, we define  $V$  by :

$$V = \{(u, p) \in \mathbb{R}^N \times \mathbb{R}^M \text{ such that } p > 0\}$$

Let the mapping  $F : V \times [0, 1] \rightarrow \mathbb{R}^N \times \mathbb{R}^M$  be given by :

$$F(u, p, \alpha) = \begin{cases} v_{\sigma, i}, & \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d \\ q_K, & \forall K \in \mathcal{M} \end{cases}$$



with :

$$\begin{cases} v_{\sigma,i} = a(u, \varphi_{\sigma}^{(i)}) - \alpha \int_{\Omega,h} p \nabla \cdot \varphi_{\sigma}^{(i)} dx - \int_{\Omega} f_v \cdot \varphi_{\sigma}^{(i)} dx \\ q_K = \frac{|K|}{\delta t} [\varrho(p_K) - \varrho(p_K^*)] + \alpha \sum_{\sigma=K|L} v_{\sigma,K}^+ \varrho_{\alpha}(p_K) - v_{\sigma,K}^- \varrho(p_L) \end{cases} \quad (\text{II.20})$$

where,  $\forall K \in \mathcal{M}$ ,  $p_K^*$  is chosen such that  $\rho_K^* = \varrho(p_K^*)$ ; note that, by assumption,  $\varrho(\cdot)$  is one to one from  $(0, +\infty)$  to  $(0, +\infty)$ , so the previous definition makes sense. The problem  $F(u, p, 1) = 0$  is exactly system (II.8).

The present proof is obtained by applying theorem II.2.5 with  $b = 0$ ; we are thus going to show that any solution of  $F(u, p, \alpha) = 0$  satisfies suitable *a priori* estimates. To this purpose, we progress as follows. First, lemma II.2.4 shows that the pressure is positive, thus theorem II.2.1 applies, and we obtain a control on  $u$  in the discrete norm associated to  $a(\cdot, \cdot)$ , uniform with respect to  $\alpha$ . Since we work in a finite dimensional space, we then obtain a control on  $p$  in  $L^\infty$  by using the conservativity of the system of equations. For the same reason, the control on  $u$  yields a bound in  $L^\infty$  of the value of the discrete divergence, which is shown to allow, by lemma II.2.4, to bound  $p$  away from zero independently of  $\alpha$ . The proof finally ends by examining the properties of the system  $F(u, p, 0) = 0$ .

Step 1 :  $\alpha \in (0, 1]$ ,  $\|\cdot\|_*$  estimate for the velocity.

Applying lemma II.2.4 to the second equation  $F(u, p, \alpha) = 0$  (*i.e.* the relation obtained by setting  $q_K = 0$  in (II.20)) with  $\varphi_1(\cdot) = \varphi_2(\cdot) = \varrho(\cdot)$ , we get :

$$\forall K \in \mathcal{M}, \quad p_K \geq \bar{p}_\alpha \quad (\text{II.21})$$

where  $\bar{p}_\alpha$  is given by :

$$\varrho(\bar{p}_\alpha) = \frac{\min_{K \in \mathcal{M}} \varrho(p_K^*)}{1 + \delta t \max_{K \in \mathcal{M}} \left[ 0, \frac{\alpha}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right]}$$

Note that  $\bar{p}_\alpha$  is well defined since, by assumption,  $\varrho(\cdot)$  is one to one from  $(0, +\infty)$  to  $(0, +\infty)$ , and  $\bar{p}_\alpha > 0$ , for any discrete velocity field  $u$ . The pressure is thus positive. Setting now  $v_\sigma = 0$  in (II.20), multiplying the corresponding equation by  $u_{\sigma,i}$  and summing over  $\sigma \in \mathcal{E}_{\text{int}}$  and  $1 \leq i \leq d$  yields the following equation :

$$a(u, u) - \alpha \int_{\Omega,h} p \nabla \cdot u dx = \int_{\Omega} f_v \cdot u dx$$

Since the pressure is positive, by a computation very similar to the proof of theorem II.2.1, we see that, from the second relation of (II.20) with  $q_K = 0$  :

$$-\alpha \int_{\Omega,h} p \nabla \cdot u dx \geq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [\rho_K P(\rho_K) - \rho_K^* P(\rho_K^*)]$$

where  $\rho_K = \varrho(p_K)$ . By the stability of the bilinear form  $a(\cdot, \cdot)$  and Young's inequality, we thus get :

$$\underbrace{\frac{c_a}{2} \|u\|_*^2}_{T_1} + \underbrace{\frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| \rho_K P(\rho_K)}_{T_2} \leq \frac{1}{2c_a} \|f_v\|^{*2} + \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| \rho_K^* P(\rho_K^*) \quad (\text{II.22})$$

By assumption,  $T_2 \geq -C_p |\Omega|$  and we thus get the following estimate on the discrete norm of the velocity :

$$\|u\|_* \leq C_1 \tag{II.23}$$

where  $C_1$  only depends on the data of the problem, *i.e.* the bilinear form  $a(\cdot, \cdot)$ ,  $f_v$ ,  $\rho^*$ , the mesh and  $\delta t$  and not on  $\alpha$ .

Step 2 :  $\alpha \in (0, 1]$ ,  $L^\infty$  estimate for the pressure.

Let us now turn to the estimate of the pressure. By conservativity of the discrete mass balance, it is easily seen that :

$$\sum_{K \in \mathcal{M}} |K| \varrho(p_K) = \sum_{K \in \mathcal{M}} |K| \rho_K^*$$

Since each term in the sum on the left hand side is non-negative, we thus have :

$$\forall K \in \mathcal{M}, \quad \varrho(p_K) \leq \frac{1}{\min_{K \in \mathcal{M}} (|K|)} \sum_{K \in \mathcal{M}} |K| \rho_K^*$$

which, since by assumption  $\lim_{z \rightarrow +\infty} \varrho(z) = +\infty$ , yields :

$$\forall K \in \mathcal{M}, \quad p_K \leq C_2 \tag{II.24}$$

where  $C_2$  only depends on the data of the problem.

Step 3 :  $\alpha \in (0, 1]$ ,  $p$  bounded away from zero.

We now exploit the estimate (II.21). As  $\alpha \leq 1$ , we get :

$$\varrho(\bar{p}_\alpha) \geq \frac{\min_{K \in \mathcal{M}} \varrho(p_K^*)}{1 + \delta t \max_{K \in \mathcal{M}} \left[ 0, \frac{1}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right]}$$

and, by equivalence of norms in a finite dimensional space, the bound (II.23) also yields a bound in the  $L^\infty$  norm and, finally, an upper bound for the denominator of the fraction at the right hand side of this relation. We thus get, still since  $\varrho(\cdot)$  is increasing on  $(0, +\infty)$ , that,  $\forall \alpha \in (0, 1]$ ,  $\bar{p}_\alpha \geq \epsilon_1$ , and, finally :

$$\forall K \in \mathcal{M}, \quad p_K \geq \epsilon_1 \tag{II.25}$$

where  $\epsilon_1$  only depends on the data.

Step 4 : conclusion.

For  $\alpha = 0$ , the system  $F(u, p, 0) = 0$  reads :

$$\left| \begin{array}{ll} a(u, \varphi_\sigma^{(i)}) = \int_{\Omega} f_v \cdot \varphi_\sigma^{(i)} dx & \forall \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d \\ \varrho(p_K) = \varrho(p_K^*) & \forall K \in \mathcal{M} \end{array} \right.$$

Since  $\varrho(\cdot)$  is one to one from  $(0, +\infty)$  to  $(0, +\infty)$  and thanks to the stability of the bilinear form  $a(\cdot, \cdot)$ , this system has one and only one solution (which, for the pressure, reads of course  $p_K = p_K^*$ ,  $\forall K \in \mathcal{M}$ ), which satisfies :

$$\|u\|_* \leq C_3, \quad \epsilon_2 = \min_{K \in \mathcal{M}} p_K^* \leq p \leq \max_{K \in \mathcal{M}} p_K^* = C_4 \tag{II.26}$$

Let  $W$  be defined by :

$$W = \{(u, p) \in \mathbb{R}^N \times \mathbb{R}^M \text{ such that } \|u\|_* < 2 \max(C_1, C_3) \text{ and } \frac{1}{2} \min(\epsilon_1, \epsilon_2) < p < 2 \max(C_2, C_4)\}$$

We now need to prove that the topological degree  $d_0$  of  $F(\cdot, \cdot, 0)$  with respect to 0 and  $W$  is not zero. Let us first suppose that the function  $\varrho(\cdot)$  is continuously differentiable and that its derivative is positive over  $(0, +\infty)$ . The Jacobian matrix of the system  $F(u, p, 0) = 0$  is block diagonal : the first block, associated to the first relation, is constant (this part of the system is linear) and non-singular ; the second one, associated to the second relation, is diagonal, and each diagonal entry is equal to the derivative of  $\varrho(\cdot)$ , taken at the considered point. The determinant of this Jacobian matrix thus does not vanish for the solution of the system, and  $d_0 \neq 0$ . This proof can then be extended to a continuous increasing function  $\varrho(\cdot)$  by a regularization technique. Hence, finally, by inequalities (II.23), (II.24), (II.25) and (II.26), theorem II.2.5 applies, which concludes the proof.  $\blacksquare$

## II.2.4 Some cases of application

First of all, let us give some examples for the bilinear form  $a(\cdot, \cdot)$ , for which the theory developed in this work holds. The first of them is :

$$a(u, v) = \int_{\Omega} u \cdot v \, dx, \quad \|u\|_* = \|u\|_{L^2(\Omega)^d}, \quad \|f_v\|^* = \|f_v\|_{L^2(\Omega)^d}$$

This choice for  $a(\cdot, \cdot)$  yields a discrete Darcy-like problem which is, up to numerical integration technicalities, the projection step arising in the pressure correction scheme which is considered in the present paper (see section II.3). Note that, in this case, the boundary condition  $u \in H_0^1(\Omega)^d$  does not make sense at the continuous level ; in addition, the considered discretization is known to be not consistent enough to yield convergence (see remark II.3.12 hereafter) for the Darcy problem.

The bilinear form associated to the Stokes problem provides another example of application. It may read in this case :

$$a(u, v) = \int_{\Omega, h} \nabla u \cdot \nabla v \, dx, \quad \|u\|_* = |u|_{H^1(\Omega)^d}$$

or, without additional theoretical difficulties :

$$a(u, v) = \mu \int_{\Omega, h} \nabla u \cdot \nabla v \, dx + \frac{\mu}{3} \int_{\Omega, h} (\nabla \cdot u) (\nabla \cdot v) \, dx$$

this latter form, where the real number  $\mu > 0$  is the viscosity, corresponding to the physical shear stress tensor expression for a compressible flow of a constant viscosity Newtonian fluid.

In addition, consider a time step of a (semi-)implicit time discretization of the unsteady Navier-Stokes equations, in which case  $a(\cdot, \cdot)$  and  $f_v$  read :

$$a(u, v) = \frac{1}{\delta t} \int_{\Omega} \rho u \cdot v \, dx + \int_{\Omega, h} \nabla \cdot (\rho w \otimes u) \cdot v \, dx + \mu \int_{\Omega, h} \nabla u \cdot \nabla v \, dx + \frac{\mu}{3} \int_{\Omega, h} (\nabla \cdot u) (\nabla \cdot v) \, dx$$

$$f_v = \frac{1}{\delta t} \rho^* u^* + f_{v,0}$$

where  $f_{v,0}$  is the physical forcing term,  $\rho^*$  and  $u^*$  stand for known density and velocity fields and  $w$  is an advection field, which may be  $u$  itself or be derived from the velocity obtained at the previous time steps. Let us suppose that the following identity holds :

$$\frac{1}{\delta t} \int_{\Omega} (\rho u - \rho^* u^*) \cdot u \, dx + \int_{\Omega, h} \nabla \cdot (\rho w \otimes u) \cdot u \, dx \geq \frac{1}{2\delta t} \left[ \int_{\Omega} \rho |u|^2 - \int_{\Omega} \rho^* |u^*|^2 \right]$$

which is the discrete counterpart of equation (II.5)-(i). The algorithm considered in this paper provides an example where this condition is verified (see section II.3). Then the present theory applies with little modifications : in the proof of existence of theorem II.2.6, the right hand side of the preceding equation must be multiplied by the homotopy parameter  $\alpha$  (and thus this term vanishes at  $\alpha = 0$ , which yields the problem considered in step 4 above) ; the (uniform with respect to  $\alpha$ ) stability in step 1 stems from the diffusion term, and steps 2 and 3 remain unchanged.

Note finally that, in the steady state case, an additional constraint is needed for the problem to have a chance to be well posed, namely to impose the total mass  $M$  of fluid in the computational domain to a given value. This constraint can be simply enforced by solving an approximate mass balance which reads :

$$c(h) \left[ \rho - \frac{M}{|\Omega|} \right] + \nabla \cdot \rho u = 0$$

where  $|\Omega|$  stands for the measure of  $\Omega$ ,  $h$  is the spatial discretization step and  $c(h) > 0$  must tend to zero with  $h$ , fast enough to avoid any loss of consistency. With this form of the mass balance, the theory developed here directly applies to this case too, provided that the corresponding unsteady-like term is also introduced in the momentum balance equation.

Examining now the assumptions for the equation of state in theorem II.2.6, we see that our results hold with equations of state of the form :

$$\varrho(p) = p^{1/\gamma} \quad \text{or, equivalently} \quad \rho = p^\gamma, \quad \text{where } \gamma > 1$$

In this case, the elastic potential is given by equation (II.4), which yields :

$$P(\rho) = \frac{1}{\gamma-1} \rho^{\gamma-1}, \quad \rho P(\rho) = \frac{1}{\gamma-1} \rho^\gamma \quad (= \frac{1}{\gamma-1} p)$$

The same conclusion still holds with  $\gamma = 1$  (*i.e.*  $p = \rho$ ), with  $P(\rho) = \log(\rho)$  satisfying equation (II.3). The case  $\gamma > 1$  is for instance encountered for isentropic perfect gas flows, whereas  $\gamma = 1$  corresponds to the isothermal case. It is worth noting that this range of application is larger than what is known for the continuous case, for which the existence of a solution is known only in the case  $\gamma > d/2$  [52, 28, 57].

### II.3 A pressure correction scheme

In this section, we build a pressure correction numerical scheme for the solution of the compressible barotropic Navier-Stokes equations (II.1), based on the low order non-conforming finite element spaces used in the previous section, namely the Crouzeix-Raviart or Rannacher-Turek elements.

The presentation is organized as follows. First, we write the scheme in the time semi-discrete setting (section II.3.1). Then we prove a general stability estimate which applies to the discretization by a finite volume technique of the convection operator (section II.3.2). The proposed

scheme is built in such a way that the assumptions of this stability result hold (section II.3.3); this implies first a prediction of the density, as a non-standard first step of the algorithm and, second, a discretization of the convection terms in the momentum balance equation by a finite volume technique which is especially designed to this purpose. The discretization of the projection step (section II.3.4) also combines the finite element and finite volume methods, in such a way that the theory developed in section II.2 applies; in particular, the proposed discretization allows to take benefit of the pressure or density control induced by the pressure work, *i.e.* to apply theorem II.2.1. The remaining steps of the algorithm are described in section II.3.5 and an overview of the scheme is given in section II.3.6. The following section (section II.3.7) is devoted to the proof of the stability of the algorithm. Finally, we shortly address some implementation difficulties (section II.3.8), then we provide some numerical tests (section II.3.9) which are performed to assess the time and space convergence of the scheme.

### II.3.1 Time semi-discrete formulation

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 uniform. Let  $\delta t$  be the constant time step  $\delta t = t_{k+1} - t_k$  for  $k = 0, 1, \dots, n-1$ . In a time semi-discrete setting, the scheme considered in this paper reads :

$$1 - \text{Solve for } \tilde{\rho}^{n+1} : \quad \frac{\tilde{\rho}^{n+1} - \rho^n}{\delta t} + \nabla \cdot (\tilde{\rho}^{n+1} u^n) = 0 \quad (\text{II.27})$$

$$2 - \text{Solve for } \tilde{p}^{n+1} : \quad -\nabla \cdot \left( \frac{1}{\tilde{\rho}^{n+1}} \nabla \tilde{p}^{n+1} \right) = -\nabla \cdot \left( \frac{1}{\sqrt{\tilde{\rho}^{n+1} \tilde{\rho}^n}} \nabla p^n \right) \quad (\text{II.28})$$

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

$$\frac{\tilde{\rho}^{n+1} \tilde{u}^{n+1} - \rho^n u^n}{\delta t} + \nabla \cdot (\tilde{\rho}^{n+1} u^n \otimes \tilde{u}^{n+1}) + \nabla \tilde{p}^{n+1} - \nabla \cdot \tau(\tilde{u}^{n+1}) = f_v^{n+1} \quad (\text{II.29})$$

4 - Solve for  $\bar{u}^{n+1}, p^{n+1}, \rho^{n+1}$  :

$$\left\{ \begin{array}{l} \tilde{\rho}^{n+1} \frac{\bar{u}^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla (p^{n+1} - \tilde{p}^{n+1}) = 0 \\ \frac{\varrho(p^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho(p^{n+1}) \bar{u}^{n+1}) = 0 \\ \rho^{n+1} = \varrho(p^{n+1}) \end{array} \right. \quad (\text{II.30})$$

$$5 - \text{Compute } u^{n+1} \text{ given by : } \quad \sqrt{\rho^{n+1}} u^{n+1} = \sqrt{\tilde{\rho}^{n+1}} \bar{u}^{n+1} \quad (\text{II.31})$$

The first step is a prediction of the density, used for the discretization of the time derivative of the momentum. As remarked by Wesseling *et al* [7, 75], this step can be avoided when solving the Euler equations : in this case, the mass flowrate may be chosen as an unknown, using the explicit velocity as an advective field in the discretization of the convection term in the momentum balance ; the velocity is then updated by dividing by the density at the end of the time step. For viscous flows, if the discretization of the diffusion term is chosen to be implicit, both the mass flowrate and the velocity appear as unknowns in the momentum balance ; this seems to impede the use of this trick. Let us emphasize that the special way that step one is carried out (*i.e.* solving a discretization of the mass balance instead as, for instance, performing a Richardson's extrapolation) is crucial for the stability.

Likewise, the second step is a renormalization of the pressure the interest of which is clarified only by the stability analysis. A similar technique has already been introduced by Guermond and Quartapelle for variable density incompressible flows [37].

Step 3 consists in a classical semi-implicit solution of the momentum balance equation to obtain a predicted velocity.

Step 4 is a nonlinear pressure correction step, which degenerates in the usual projection step as used in incompressible flow solvers when the density is constant (e.g. [53]). Taking the divergence of the first relation of (II.30) and using the second one to eliminate the unknown velocity  $\bar{u}^{n+1}$  yields a non-linear elliptic problem for the pressure. This computation is formal in the semi-discrete formulation, but, of course, is necessarily made clear at the algebraic level, as described in section II.3.8. Once the pressure is computed, the first relation yields the updated velocity and the third one gives the end-of-step density.

Finally, step 5 is a renormalization of the velocity, once again useful for stability reasons.

### II.3.2 Stability of the advection operator : a finite-volume result

The aim of this section is to state and prove a discrete analogue to the stability identity (II.5)-(i), which may be written for any sufficiently regular functions  $\rho$ ,  $z$  and  $u$  as follows :

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

and which holds if the velocity  $u$  vanishes at the boundary of the computational domain  $\Omega$  and provided that the following balance is satisfied by  $\rho$  and  $u$  :

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0$$

As stated in introduction, applying this identity to each component of the velocity yields the central argument of the proof of the kinetic energy theorem.

The discrete analogue to this identity follows. This result is presented in a general algebraic setting, without making reference to any continuous partial differential equation (see remark II.3.8 hereafter for a clarification of this link) ; note however that, in the following relations, the sum of the fluxes is restricted to the internal edges of the mesh, which implicitly reflects the fact that the normal velocity is supposed to be zero at the boundary.

#### Theorem II.3.7 (Stability of the advection operator)

*Let  $(\rho_K^*)_{K \in \mathcal{M}}$  and  $(\rho_K)_{K \in \mathcal{M}}$  be two families of positive real numbers satisfying the following set of equations :*

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K - \rho_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} = 0 \quad (\text{II.32})$$

*where  $F_{\sigma,K}$  is a quantity associated to the edge  $\sigma$  and to the control volume  $K$  ; we suppose that, for any internal edge  $\sigma = K|L$ ,  $F_{\sigma,K} = -F_{\sigma,L}$ . Let  $(z_K^*)_{K \in \mathcal{M}}$  and  $(z_K)_{K \in \mathcal{M}}$  be two families of real numbers. For any internal edge  $\sigma = K|L$ , we define  $z_{\sigma}$  either by  $z_{\sigma} = \frac{1}{2}(z_K + z_L)$ , or by  $z_{\sigma} = z_K$  if  $F_{\sigma,K} \geq 0$  and  $z_{\sigma} = z_L$  otherwise. The first choice is usually referred to as the "centered choice", the second one as "the upwind choice". In both cases, the following stability property holds :*

$$\sum_{K \in \mathcal{M}} z_K \left[ \frac{|K|}{\delta t} (\rho_K z_K - \rho_K^* z_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} z_{\sigma} \right] \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \left[ \rho_K z_K^2 - \rho_K^* z_K^{*2} \right] \quad (\text{II.33})$$

**Proof.**

We write :

$$\sum_{K \in \mathcal{M}} z_K \left[ \frac{|K|}{\delta t} (\rho_K z_K - \rho_K^* z_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} z_\sigma \right] = T_1 + T_2$$

where  $T_1$  and  $T_2$  read :

$$T_1 = \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} z_K (\rho_K z_K - \rho_K^* z_K^*), \quad T_2 = \sum_{K \in \mathcal{M}} z_K \left[ \sum_{\sigma=K|L} F_{\sigma,K} z_\sigma \right]$$

The first term reads :

$$T_1 = \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [z_K^2 (\rho_K - \rho_K^*) + \rho_K^* z_K (z_K - z_K^*)]$$

Developping the last term by the identity  $a(a-b) = \frac{1}{2}(a^2 + (a-b)^2 - b^2)$ , we get :

$$T_1 = \underbrace{\sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} z_K^2 (\rho_K - \rho_K^*)}_{T_{1,1}} + \underbrace{\frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \rho_K^* (z_K^2 - z_K^{*2})}_{T_{1,2}} + \underbrace{\frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \rho_K^* (z_K - z_K^*)^2}_{T_{1,3}}$$

The last term, namely  $T_{1,3}$ , is always non-negative and can be seen as a dissipation associated to the backward time discretization of equation (II.33). We now turn to  $T_2$  :

$$T_2 = \underbrace{\sum_{K \in \mathcal{M}} z_K^2 \left[ \sum_{\sigma=K|L} F_{\sigma,K} \right]}_{T_{2,1}} + \underbrace{\sum_{K \in \mathcal{M}} z_K \left[ \sum_{\sigma=K|L} F_{\sigma,K} (z_\sigma - z_K) \right]}_{T_{2,2}}$$

The first term, namely  $T_{2,1}$ , will cancel with  $T_{1,1}$  by equation (II.32). The second term reads, developping as previously the quantity  $z_K (z_\sigma - z_K)$  :

$$T_{2,2} = -\frac{1}{2} \sum_{K \in \mathcal{M}} z_K^2 \left[ \sum_{\sigma=K|L} F_{\sigma,K} \right] - \underbrace{\frac{1}{2} \sum_{K \in \mathcal{M}} \left[ \sum_{\sigma=K|L} F_{\sigma,K} [(z_\sigma - z_K)^2 - z_\sigma^2] \right]}_{T_{2,3}}$$

Reordering the sum in the last term, we have, as  $F_{\sigma,K} = -F_{\sigma,L}$  :

$$T_{2,3} = \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} F_{\sigma,K} [(z_\sigma - z_K)^2 - (z_\sigma - z_L)^2]$$

This expression can easily be seen to vanish with the centered choice. With the upwind choice, supposing without loss of generality that we have chosen for the edge  $\sigma = K|L$  the orientation such that  $F_{\sigma,K} \geq 0$ , we get, as  $z_\sigma = z_K$  :

$$T_{2,3} = -\frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} F_{\sigma,K} (z_K - z_L)^2 \leq 0$$

We thus have, by equation (II.32) :

$$T_{2,2} \geq -\frac{1}{2} \sum_{K \in \mathcal{M}} z_K^2 \left[ \sum_{\sigma=K|L} F_{\sigma,K} \right] = \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} z_K^2 (\rho_K - \rho_K^*)$$

and thus :

$$T_1 + T_2 \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \left[ z_K^2 (\rho_K - \rho_K^*) + \rho_K^* (z_K^2 - z_K^{*2}) \right]$$

which concludes the proof. ■

**Remark II.3.8** Equation (II.32) can be seen as a discrete mass balance, with  $F_{\sigma,K}$  standing for the mass flux across the edge  $\sigma$ , and the right hand side of (II.33) may be derived by the multiplication by  $z_K$  and summation over the control volumes of the transport terms in a discrete balance equation for the quantity  $\rho z$ , reading :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K z_K - \rho_K^* z_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} z_\sigma + \dots [\text{possible diffusion terms}] \dots = 0$$

In this context, the relation (II.32) is known to be exactly the compatibility condition which ensures a discrete maximum principle for the solution  $z$  of this transport equation, provided that the upwind choice (or any monotone choice) is made for the expression of  $z_\sigma$  [51]. We proved here that the same compatibility condition ensures a  $L^2$  stability for  $\rho^{1/2} z$ .

### II.3.3 Space discretization of the density prediction and the momentum balance equation

The main difficulty in the discretization of the momentum balance equation is to build a discrete convection operator which enjoys the stability property (II.5)-(i). To this purpose, we derive for this term a finite volume discretization which satisfies the assumptions of theorem II.3.7.

The natural space discretization for the density is the same as for the pressure, *i.e.* piecewise constant functions over each element. This legitimates a standard mass lumping technique for the time derivative term, since no additional accuracy seems to have to be expected from a more complex numerical integration. Note that, for the Crouzeix-Raviart element in two dimensions, the mass matrix is genuinely diagonal. Let the quantity  $|D_\sigma|$  be defined as follows :

$$|D_\sigma| \stackrel{\text{def}}{=} \int_{\Omega} \varphi_\sigma \, dx > 0 \tag{II.34}$$

For any  $\sigma \in \mathcal{E}$  and any control volume  $K$  adjacent to  $\sigma$ , let  $D_{K,\sigma}$  be the cone of basis  $\sigma$  and having the mass center of  $K$  as opposite vertex. The volume  $D_{K,\sigma}$  is referred to as the half-diamond cell associated to  $\sigma$  and  $K$  (see figure II.1) and the measure of  $D_{K,\sigma}$  is denoted by  $|D_{K,\sigma}|$ . For  $\sigma \in \mathcal{E}$ , we define the diamond cell  $D_\sigma$  associated to  $\sigma$  by  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$  if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , and, if  $\sigma \in \mathcal{E}_{\text{ext}}$ , by  $D_\sigma = D_{K,\sigma}$  where  $K$  is the only control volume adjacent to  $\sigma$ . For the Crouzeix-Raviart element,  $|D_\sigma|$  can be identified to the measure of the diamond cell  $D_\sigma$  associated to  $\sigma$ . The same property holds for the Rannacher-Turek element in the case of rectangles ( $d = 2$ ) or cuboids ( $d = 3$ ), which are the only cases considered here, even though extensions to non-perpendicular grids are probably possible.

The discretization of the term  $\rho^n u^n$  thus leads, in the equations associated to the velocity on an internal edge  $\sigma$ , to an expression of the form  $\rho_\sigma^n u_\sigma^n$ , where  $\rho_\sigma^n$  results from an average of the



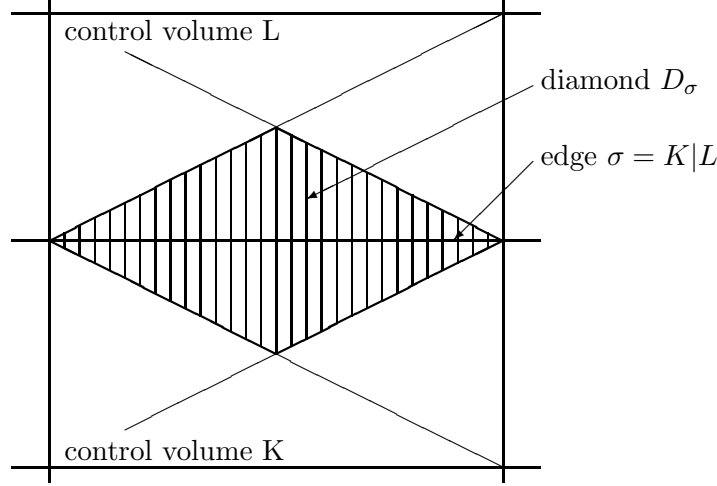


FIG. II.1 – Dual finite volume mesh : the so-called "diamond cells".

values taken by the density in the two elements adjacent to  $\sigma$ , weighted by the measure of the half-diamonds :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad |D_\sigma| \rho_\sigma^n = |D_{K,\sigma}| \rho_K^n + |D_{L,\sigma}| \rho_L^n \quad (\text{II.35})$$

This definition naturally extends to any external edge  $\sigma$  by  $\rho_\sigma^n = \rho_K^n$ , where  $K$  is the control volume which is adjacent to  $\sigma$ .

In order to satisfy the compatibility condition which was introduced in the previous section, a prediction of the density is first performed, by a finite volume discretization of the mass balance equation, taking the diamond cells as control volumes :

$$\frac{|D_\sigma|}{\delta t} (\tilde{\rho}_\sigma^{n+1} - \rho_\sigma^n) + \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma}^{n+1} = 0, \quad \forall \sigma \in \mathcal{E} \quad (\text{II.36})$$

where  $\mathcal{E}(D_\sigma)$  is the set of the edges of  $D_\sigma$  and  $F_{\varepsilon,\sigma}$  stands for the mass flux across  $\varepsilon$  outward  $D_\sigma$ . This latter quantity is expressed as follows :

$$F_\varepsilon^{n+1} = |\varepsilon| u_\varepsilon^n \cdot n_{\varepsilon,\sigma} \tilde{\rho}_\varepsilon^{n+1}$$

where  $|\varepsilon|$  is the measure of  $\varepsilon$ ,  $n_{\varepsilon,\sigma}$  is the normal to  $\varepsilon$  outward  $D_\sigma$ , the velocity  $u_\varepsilon^n$  is obtained by interpolation of  $u^n$  at the center of  $\varepsilon$  (using the standard finite element expansion) and  $\tilde{\rho}_\varepsilon^{n+1}$  is the density at the edge, calculated by the standard upwinding technique (*i.e.* either  $\tilde{\rho}_\sigma^{n+1}$  if  $u_\varepsilon^n \cdot n_{\varepsilon,\sigma} \geq 0$  or  $\tilde{\rho}_{\sigma'}^{n+1}$  otherwise, with  $\sigma'$  such that  $\varepsilon$  separates  $D_\sigma$  and  $D_{\sigma'}$ , which we denote by  $\varepsilon = D_\sigma|D_{\sigma'}$ ).

The discretization of the convection terms of the momentum balance equation is built from relation (II.36), according to the structure which is necessary to apply theorem II.3.7. This yields the following discrete momentum balance equation :

$$\begin{aligned} \frac{|D_\sigma|}{\delta t} (\tilde{\rho}_\sigma^{n+1} \tilde{u}_{\sigma,i}^{n+1} - \rho_\sigma^n u_{\sigma,i}^n) + \sum_{\substack{\varepsilon \in \mathcal{E}(D_\sigma), \\ \varepsilon = D_\sigma|D_{\sigma'}}} \frac{1}{2} F_{\varepsilon,\sigma}^{n+1} (\tilde{u}_{\sigma,i}^{n+1} + u_{\sigma',i}^{n+1}) + \int_{\Omega,h} \tau(\tilde{u}^{n+1}) : \nabla \varphi_\sigma^{(i)} dx \\ - \int_{\Omega,h} \tilde{p}^{n+1} \nabla \cdot \varphi_\sigma^{(i)} dx = \int_{\Omega} f_v^{n+1} \cdot \varphi_\sigma^{(i)}, \quad \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d \end{aligned} \quad (\text{II.37})$$

where  $\varphi_\sigma^{(i)}$  stands for the vector shape function associated to  $v_{\sigma,i}$ , which reads  $\varphi_\sigma^{(i)} = \varphi_\sigma e_i$  with  $e_i$  the  $i^{\text{th}}$  vector of the canonical basis of  $\mathbb{R}^d$  and  $\varphi_\sigma$  the scalar shape function, and where the notation  $\int_{\Omega,h}$  means  $\sum_{K \in \mathcal{M}} \int_K$ .

Note that, for Crouzeix-Raviart elements, a combined finite volume/finite element method, similar to the technique employed here for the discretization of the momentum balance, has already been analysed for a transient non-linear convection-diffusion equation by Feistauer and co-workers [1, 21, 29].

### II.3.4 Space discretization of the projection step

The fully discrete projection step of the proposed algorithm reads :

$$\left\{ \begin{array}{l} |D_\sigma| \frac{\tilde{\rho}_\sigma^{n+1}}{\delta t} (\bar{u}_{\sigma,i}^{n+1} - \tilde{u}_{\sigma,i}^{n+1}) - \int_{\Omega,h} (p^{n+1} - \tilde{p}^{n+1}) \nabla \cdot \varphi_\sigma^{(i)} dx = 0, \quad \forall \sigma \in \mathcal{E}_{\text{int}}, \text{ for } 1 \leq i \leq d \\ \frac{|K|}{\delta t} (\varrho(p_K^{n+1}) - \rho_K^n) + \sum_{\sigma=K|L} (v_{\sigma,K}^+)^{n+1} \varrho(p_K^{n+1}) - (v_{\sigma,K}^-)^{n+1} \varrho(p_L^{n+1}) = 0, \quad \forall K \in \mathcal{M} \end{array} \right. \quad (\text{II.38})$$

where  $(v_{\sigma,K}^+)^{n+1}$  and  $(v_{\sigma,K}^-)^{n+1}$  stand respectively for  $\max(v_{\sigma,K}^{n+1}, 0)$  and  $-\min(v_{\sigma,K}^{n+1}, 0)$  with  $v_{\sigma,K}^{n+1} = |\sigma| \bar{u}_\sigma^{n+1} \cdot n_{KL}$ . The first (vector) equation may be seen as the finite element discretization of the first relation of the projection step (II.30), with the same lumping of the mass matrix for the Rannacher-Turek element as in the prediction step. As the pressure is piecewise constant, the finite element discretization of the second relation of (II.30), *i.e.* the mass balance, is equivalent to a finite volume formulation, in which we introduce the standard first-order upwinding. Exploiting the expression of the velocity and pressure shape functions, the first set of relations of this system can be alternatively written as follows :

$$\begin{aligned} |D_\sigma| \frac{\tilde{\rho}_\sigma^{n+1}}{\delta t} (\bar{u}_\sigma^{n+1} - \tilde{u}_\sigma^{n+1}) \\ + |\sigma| [(p_L^{n+1} - \tilde{p}_L^{n+1}) - (p_K^{n+1} - \tilde{p}_K^{n+1})] n_{KL} = 0, \quad \forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L \end{aligned} \quad (\text{II.39})$$

or, in an algebraic setting :

$$\frac{1}{\delta t} M_{\tilde{\rho}^{n+1}} (\bar{u}^{n+1} - \tilde{u}^{n+1}) + B^t (p^{n+1} - \tilde{p}^{n+1}) = 0 \quad (\text{II.40})$$

In this relation,  $M_w$  stands for the diagonal mass matrix weighted by  $(w_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}$  (so, for  $1 \leq i \leq d$  and  $\sigma \in \mathcal{E}_{\text{int}}$ , the corresponding entry on the diagonal of  $M_{\tilde{\rho}^{n+1}}$  reads  $(M_{\tilde{\rho}^{n+1}})_{\sigma,i} = |D_\sigma| \tilde{\rho}_\sigma^{n+1}$ ),  $B^t$  is the matrix of  $\mathbb{R}^{(dN) \times M}$ , where  $N$  is the number of internal edges (*i.e.*  $N = \text{card}(\mathcal{E}_{\text{int}})$ ) and  $M$  is the number of control volumes in the mesh (*i.e.*  $M = \text{card}(\mathcal{M})$ ), associated to the gradient operator ; consequently, the matrix  $B$  is associated to the opposite of the discrete divergence operator. Throughout this section, we use the same notation for the discrete function (defined as usual in the finite element context by its expansion using the shape functions) and for the vector gathering the degrees of freedom ; so, in relation (II.40),  $\bar{u}$  (respectively  $\tilde{u}$ ) stands for the vector of  $\mathbb{R}^{dN}$  components  $\bar{u}_{\sigma,i}$  (respectively  $\tilde{u}_{\sigma,i}$ ),  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $1 \leq i \leq d$  and  $p$  (respectively  $\tilde{p}$ ) stands for the vector of  $\mathbb{R}^M$  of components  $p_K$  (respectively  $\tilde{p}_K$ ),  $K \in \mathcal{M}$ . Both forms (II.39) and (II.40) are used hereafter.

We have the following existence result.

**Proposition II.3.9**

Let the equation of state  $\rho(\cdot)$  be defined and increasing over  $[0, +\infty)$ , and be such that  $\rho(0) = 0$ ,  $\lim_{z \rightarrow +\infty} \rho(z) = +\infty$  and that there exists an elastic potential function  $P(\cdot)$  (i.e. a function satisfying (II.3)) such that the function  $z \mapsto zP(z)$  is bounded from below in  $(0, +\infty)$ , once continuously differentiable and strictly convex. Let us suppose that  $\tilde{\rho}_\sigma^{n+1} > 0$ ,  $\forall \sigma \in \mathcal{E}_{\text{int}}$ . Then the nonlinear system (II.38) admits at least one solution and any possible solution is such that  $p_K > 0$ ,  $\forall K \in \mathcal{M}$  (and thus  $\rho_K > 0$ ,  $\forall K \in \mathcal{M}$ ).

**Proof.**

The theory of section II.2 applies, with :

$$\|\bar{u}^{n+1}\|_*^2 = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma|}{\delta t} \tilde{\rho}_\sigma^{n+1} |\bar{u}_\sigma^{n+1}|^2$$

This yields both the existence of a solution and the positivity of the pressure. ■

In view of this result and of the form of the discrete density prediction (II.36), the property  $\tilde{\rho}^{n+1} > 0$  is satisfied by induction at any time step of the computation (provided, of course, that the initial density is positive everywhere).

We finish this section by some remarks concerning the projection step at hand.

**Lemma II.3.10**

The following identity holds for each discrete pressure  $q \in L_h$  :

$$\forall K \in \mathcal{M}, \quad (\text{B M}_{\tilde{\rho}^{n+1}}^{-1} \text{B}^t q)_K = \sum_{\sigma=K|L} \frac{1}{\tilde{\rho}_\sigma^{n+1}} \frac{|\sigma|^2}{|D_\sigma|} (q_K - q_L)$$

**Proof.**

Let  $q \in L_h$  be given. By relation (II.39), we have :

$$(\text{B}^t q)_{\sigma,i} = |\sigma| (q_L - q_K) n_{KL} \cdot e_i$$

Let  $1_K \in L_h$  be the characteristic function of  $K$ . Denoting by  $(\cdot, \cdot)$  the standard Euclidian inner product, by the previous relation and the definition of the lumped velocity mass matrix, we obtain :

$$\begin{aligned} (\text{B M}_{\tilde{\rho}^{n+1}}^{-1} \text{B}^t q, 1_K) &= (\text{M}_{\tilde{\rho}^{n+1}}^{-1} \text{B}^t q, \text{B}^t 1_K) = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \sum_{i=1}^d \frac{1}{\tilde{\rho}_\sigma^{n+1} |D_\sigma|} (\text{B}^t q)_{\sigma,i} (\text{B}^t 1_K)_{\sigma,i} \\ &= \sum_{\sigma=K|L} \sum_{i=1}^d \frac{1}{\tilde{\rho}_\sigma^{n+1} |D_\sigma|} [|\sigma| (q_L - q_K) n_{KL} \cdot e_i] [-|\sigma| n_{KL} \cdot e_i] \end{aligned}$$

which, remarking that  $\sum_{i=1}^d (n_{KL} \cdot e_i)^2 = 1$ , yields the result. ■

**Remark II.3.11 (On spurious pressure boundary conditions)** *In the context of projection methods for incompressible flow, it is known that spurious boundary conditions are to be imposed*

to the pressure in the projection step, in order to make the definition of this step of the algorithm complete. These boundary conditions are explicit when the process to derive the projection step is first to pose the elliptic problem for the pressure at the time semi-discrete level and then discretize it in space; for instance, with a constant density equal to one and prescribed velocity boundary conditions on  $\partial\Omega$ , the semi-discrete projection step would take the form :

$$\begin{cases} -\Delta (p^{n+1} - \tilde{p}^{n+1}) = -\frac{1}{\delta t} \nabla \cdot \tilde{u}^{n+1} & \text{in } \Omega \\ \nabla (p^{n+1} - \tilde{p}^{n+1}) \cdot n = 0 & \text{on } \partial\Omega \end{cases}$$

When the elliptic problem for the pressure is built at the algebraic level, the boundary conditions for the pressure are somehow hidden in the discrete operator  $\mathbf{B} \mathbf{M}^{-1} \mathbf{B}^t$ . Lemma II.3.10 shows that this matrix takes the form of a finite-volume Laplace discrete operator, with homogeneous Neumann boundary conditions, i.e. the same boundary conditions as in the time semi-discrete problem above stated.

**Remark II.3.12 (On the non-consistency of the discretization at hand for the Darcy problem)** Considering the semi-discrete problem (II.30), in the case of a constant density equal to one, one may expect to recover a consistent discretization of a Poisson problem with homogeneous Neumann boundary conditions, as stated above. The following example shows that this route is misleading. Let us take for the mesh a uniform square grid of step  $h$ . The coefficient  $|\sigma|^2/|D_\sigma|$  can be easily evaluated, and we obtain :

$$(\mathbf{B} \mathbf{M}^{-1} \mathbf{B}^t q)_K = d \sum_{\sigma=K|L} \frac{|\sigma|}{h} (q_K - q_L)$$

that is the usual finite volume Laplace operator, but multiplied by the space dimension  $d$ . This result is of course consistent with (and gives some insight in) the wellknown non-consistency of the Rannacher-Turek element for the Darcy problem; similar examples could also be given for simplicial grids, with the Crouzeix-Raviart element.

### II.3.5 Renormalization steps

The pressure renormalization (step 2 of the algorithm) reads, in an algebraic setting :

$$\mathbf{B} \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1} \mathbf{B}^t \tilde{p}^{n+1} = \mathbf{B} \mathbf{M}_{\sqrt{\tilde{\rho}^{n+1} \tilde{\rho}^n}}^{-1} \mathbf{B}^t p^n \quad (\text{II.41})$$

Note that, at the first time step, the quantity  $\tilde{\rho}^0$  must be defined; it can for instance be computed from the initial density (defined on the control volumes of  $\mathcal{M}$ ) by equation (II.35). In view of the expression of these operators provided by lemma II.3.10, this relation equivalently reads :

$$\sum_{\sigma=K|L} \frac{1}{\tilde{\rho}_\sigma^{n+1}} \frac{|\sigma|^2}{|D_\sigma|} (\tilde{p}_K^{n+1} - \tilde{p}_L^{n+1}) = \sum_{\sigma=K|L} \frac{1}{\sqrt{\tilde{\rho}_\sigma^{n+1} \tilde{\rho}_\sigma^n}} \frac{|\sigma|^2}{|D_\sigma|} (p_K^n - p_L^n), \quad \forall K \in \mathcal{M} \quad (\text{II.42})$$

As  $\mathbf{B}^t$  and  $\mathbf{B}$  stands for respectively the discrete gradient and (opposite of the) divergence operator, this system can be seen as a discretization of the semi-discrete expression of step 2; note however, as shown in remark II.3.12, that this discretization is non-consistent.

The velocity renormalization (step 5 of the algorithm) simply reads :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad \sqrt{\tilde{\rho}_\sigma^{n+1}} u_\sigma^{n+1} = \sqrt{\tilde{\rho}_\sigma^{n+1}} \tilde{u}_\sigma^{n+1} \quad \text{or} \quad \mathbf{M}_{\sqrt{\tilde{\rho}^{n+1}}} u^{n+1} = \mathbf{M}_{\sqrt{\tilde{\rho}^{n+1}}} \tilde{u}^{n+1} \quad (\text{II.43})$$

### II.3.6 An overview of the algorithm

To sum up, the algorithm considered in this section is the following one :

1. Prediction of the density – The density on the edges at  $t^n$ ,  $(\rho_\sigma^n)_{\sigma \in \mathcal{E}}$ , being given by (II.35), compute  $(\tilde{\rho}^{n+1})_{\sigma \in \mathcal{E}}$  by the upwind finite volume discretization of the mass balance over the diamond cells (II.36).
2. Renormalization of the pressure – Compute a renormalized pressure  $(\tilde{p}_K^{n+1})_{K \in \mathcal{M}}$  by equation (II.42).
3. Prediction of the velocity – Compute  $(\tilde{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  by equation (II.37), obtained by a finite volume discretization of the transport terms over the diamond cells and a finite element discretization of the other terms.
4. Projection step – Compute  $(\bar{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  and  $(p_K^{n+1})_{K \in \mathcal{M}}$  from equation (II.38), obtained by a finite element discretization of the velocity correction equation and an upwind finite volume discretization of the mass balance (over the elements  $K \in \mathcal{M}$ ).
5. Renormalization of the velocity – Compute  $(u_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  from equation (II.43).

The existence of a solution to step 4 is proven above ; the other problems are linear, and their well-posedness follows from standard coercivity arguments, using the fact that the discrete densities (*i.e.*  $\rho^n$  and  $\tilde{\rho}^{n+1}$ ) are positive, provided that this property is satisfied by the initial condition.

### II.3.7 Stability analysis

In this section, we use the following discrete norm and semi-norm :

$$\begin{aligned}
 \forall v \in W_h, \quad \|v\|_{h, \tilde{\rho}}^2 &= \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \tilde{\rho}_\sigma |v_\sigma|^2 \\
 \forall q \in L_h, \quad |q|_{h, \tilde{\rho}}^2 &= \sum_{\sigma \in \mathcal{E}_{\text{int}}, \sigma=K|L} \frac{1}{\tilde{\rho}_\sigma} \frac{|\sigma|^2}{|D_\sigma|} (q_K - q_L)^2
 \end{aligned} \tag{II.44}$$

where  $\tilde{\rho} = (\tilde{\rho}_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}$  is a family of positive real numbers. The function  $\|\cdot\|_{h, \tilde{\rho}}$  defines a norm over  $W_h$ , and  $|\cdot|_{h, \tilde{\rho}}$  can be seen as a weighted version (with mesh dependent weights) of the  $H^1$  semi-norm classical in the finite volume context [26]. The links between this latter semi-norm and the problem at hand are clarified in the following lemma, which is a straightforward consequence of lemma II.3.10.

#### **Lemma II.3.13**

*The following identity holds for each discrete pressure  $q \in L_h$  :*

$$(\text{B M}_{\tilde{\rho}}^{-1} \text{B}^t q, q) = |q|_{h, \tilde{\rho}}^2$$

We are now in position to state the stability of the scheme under consideration.

**Theorem II.3.14 (Stability of the scheme)**

Let the equation of state  $\varrho(\cdot)$  be defined and increasing over  $[0, +\infty)$ , and be such that  $\varrho(0) = 0$ ,  $\lim_{z \rightarrow +\infty} \varrho(z) = +\infty$  and that there exists an elastic potential function  $P(\cdot)$  (i.e. a function satisfying (II.3)) such that the function  $z \mapsto z P(z)$  is bounded from below in  $(0, +\infty)$ , once continuously differentiable and strictly convex. Let  $(\tilde{u}^n)_{0 \leq n \leq N}$ ,  $(u^n)_{0 \leq n \leq N}$ ,  $(p^n)_{0 \leq n \leq N}$  and  $(\rho^n)_{0 \leq n \leq N}$  be the solution to the considered scheme, with a zero forcing term. Then the following bound holds for  $0 \leq n < N$  :

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^{n+1}}^2 + \int_{\Omega} \rho^{n+1} P(\rho^{n+1}) \, dx + \delta t \sum_{k=1}^{n+1} \int_{\Omega, h} \nabla \tilde{u}^k : \tau(\tilde{u}^k) \, dx + \frac{\delta t^2}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 \\ \leq \frac{1}{2} \|u^0\|_{h, \rho^0}^2 + \int_{\Omega} \rho^0 P(\rho^0) \, dx + \frac{\delta t^2}{2} |p^0|_{h, \tilde{\rho}^0}^2 \end{aligned} \quad (\text{II.45})$$

**Proof.**

Multiplying each equation of the step 3 of the scheme (II.37) by the corresponding unknown (i.e. the corresponding component of the velocity  $\tilde{u}^{n+1}$  on the corresponding edge  $\sigma$ ) and summing over the edges and the components yields, by virtue of the stability of the discrete advection operator (theorem II.3.7) :

$$\frac{1}{2\delta t} \|\tilde{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2 - \frac{1}{2\delta t} \|u^n\|_{h, \rho^n}^2 + \int_{\Omega, h} \tau(\tilde{u}^{n+1}) : \nabla \tilde{u}^{n+1} \, dx - \int_{\Omega, h} \tilde{p}^{n+1} \nabla \cdot \tilde{u}^{n+1} \, dx \leq 0 \quad (\text{II.46})$$

On the other hand, reordering equation (II.40) and multiplying by  $M_{\tilde{\rho}^{n+1}}^{-1/2}$  (recall that  $M_{\tilde{\rho}^{n+1}}$  is diagonal), we obtain :

$$\frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \bar{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t p^{n+1} = \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \tilde{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t \tilde{p}^{n+1}$$

Squaring this relation gives :

$$\begin{aligned} \left( \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \bar{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t p^{n+1}, \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \bar{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t p^{n+1} \right) = \\ \left( \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \tilde{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t \tilde{p}^{n+1}, \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}}^{1/2} \tilde{u}^{n+1} + M_{\tilde{\rho}^{n+1}}^{-1/2} B^t \tilde{p}^{n+1} \right) \end{aligned}$$

which reads :

$$\begin{aligned} \frac{1}{\delta t^2} (M_{\tilde{\rho}^{n+1}} \bar{u}^{n+1}, \bar{u}^{n+1}) + (M_{\tilde{\rho}^{n+1}}^{-1} B^t p^{n+1}, B^t p^{n+1}) + \frac{2}{\delta t} (\bar{u}^{n+1}, B^t p^{n+1}) = \\ \frac{1}{\delta t^2} (M_{\tilde{\rho}^{n+1}} \tilde{u}^{n+1}, \tilde{u}^{n+1}) + (M_{\tilde{\rho}^{n+1}}^{-1} B^t \tilde{p}^{n+1}, B^t \tilde{p}^{n+1}) + \frac{2}{\delta t} (\tilde{u}^{n+1}, B^t \tilde{p}^{n+1}) \end{aligned}$$

Multiplying by  $\delta t/2$ , remarking that,  $\forall v \in W_h$ ,  $(M_{\tilde{\rho}^{n+1}} v, v) = \|v\|_{h, \tilde{\rho}^{n+1}}^2$  and that, thanks to lemma II.3.13,  $\forall q \in L_h$ ,  $(M_{\tilde{\rho}^{n+1}}^{-1} B^t q, B^t q) = (B M_{\tilde{\rho}^{n+1}}^{-1} B^t q, q) = |q|_{h, \tilde{\rho}^{n+1}}^2$ , we get :

$$\begin{aligned} \frac{1}{2\delta t} \|\bar{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2 + \frac{\delta t}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 + (\bar{u}^{n+1}, B^t p^{n+1}) \\ - \frac{1}{2\delta t} \|\tilde{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2 - \frac{\delta t}{2} |\tilde{p}^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 - (\tilde{u}^{n+1}, B^t \tilde{p}^{n+1}) = 0 \end{aligned} \quad (\text{II.47})$$

The quantity  $-(\tilde{u}^{n+1}, \mathbf{B}^t \tilde{p}^{n+1})$  is nothing more than the opposite of the term  $\int_{\Omega, h} \tilde{p}^{n+1} \nabla \cdot \tilde{u}^{n+1} \, dx$  appearing in (II.46), so, when summing (II.46) and (II.47), these terms cancel, leading to :

$$\begin{aligned} \frac{1}{2\delta t} \|\bar{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2 - \frac{1}{2\delta t} \|u^n\|_{h, \rho^n}^2 + \int_{\Omega, h} \tau(\tilde{u}^{n+1}) : \nabla \tilde{u}^{n+1} \, dx \\ + \frac{\delta t}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 - \frac{\delta t}{2} |\tilde{p}^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 + (\bar{u}^{n+1}, \mathbf{B}^t p^{n+1}) \leq 0 \end{aligned}$$

Finally,  $(\bar{u}^{n+1}, \mathbf{B}^t p^{n+1})$  is precisely the pressure work which can be bounded by the time derivative of the elastic potential, as stated in theorem II.2.1 :

$$\begin{aligned} \frac{1}{2\delta t} \|\bar{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2 + \int_{\Omega, h} \tau(\tilde{u}^{n+1}) : \nabla \tilde{u}^{n+1} \, dx + \frac{\delta t}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 - \frac{\delta t}{2} |\tilde{p}^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 \\ + \frac{1}{\delta t} \int_{\Omega} \rho^{n+1} P(\rho^{n+1}) \, dx \leq \frac{1}{2\delta t} \|u^n\|_{h, \rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} \rho^n P(\rho^n) \, dx \end{aligned} \quad (\text{II.48})$$

The proof then ends by using the renormalization steps (step 2 and 5 of the algorithm). Step 2 reads in an algebraic setting :

$$\mathbf{B} \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1} \mathbf{B}^t \tilde{p}^{n+1} = \mathbf{B} \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{M}_{\tilde{\rho}^n}^{-1/2} \mathbf{B}^t p^n$$

Multiplying by  $\tilde{p}^{n+1}$ , we obtain :

$$\left( \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1}, \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1} \right) = \left( \mathbf{M}_{\tilde{\rho}^n}^{-1/2} \mathbf{B}^t p^n, \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1} \right)$$

and thus, by Cauchy-Schwartz inequality :

$$\left( \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1}, \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1} \right) \leq \left( \mathbf{M}_{\tilde{\rho}^n}^{-1/2} \mathbf{B}^t p^n, \mathbf{M}_{\tilde{\rho}^n}^{-1/2} \mathbf{B}^t p^n \right)^{1/2} \left( \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1}, \mathbf{M}_{\tilde{\rho}^{n+1}}^{-1/2} \mathbf{B}^t \tilde{p}^{n+1} \right)^{1/2}$$

This relation yields  $|\tilde{p}^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 \leq |p^n|_{h, \tilde{\rho}^n}^2$ . In addition, step 5 of the algorithm gives  $\|u^{n+1}\|_{h, \rho^{n+1}}^2 = \|\bar{u}^{n+1}\|_{h, \tilde{\rho}^{n+1}}^2$ . Using these two relations in (II.48), we get :

$$\begin{aligned} \frac{1}{2\delta t} \|u^{n+1}\|_{h, \rho^{n+1}}^2 + \int_{\Omega, h} \tau(\tilde{u}^{n+1}) : \nabla \tilde{u}^{n+1} \, dx + \frac{\delta t}{2} |p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2 + \frac{1}{\delta t} \int_{\Omega} \rho^{n+1} P(\rho^{n+1}) \, dx \\ \leq \frac{1}{2\delta t} \|u^n\|_{h, \rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} \rho^n P(\rho^n) \, dx + \frac{\delta t}{2} |p^n|_{h, \tilde{\rho}^n}^2 \end{aligned}$$

and the estimate of theorem II.3.14 follows by summing over the time steps. ■

**Remark II.3.15 (Entropy conservation)** *In this paper, we employ the terminology of the mathematical analysis of the compressible Navier-Stokes equations, as can be found in [52, 28, 57]. In this context, the function  $P(\cdot)$  is referred to as the elastic potential, and the total energy of the system is the sum of the kinetic energy and of the integral of  $\rho P(\rho)$  over the fluid domain. In the literature concerned with hyperbolic problems, the same quantity is referred to as the entropy of the system, and theorem II.3.14 states that the considered pressure correction scheme preserves the entropy.*

**Remark II.3.16 (On the upwinding of the mass balance discretization, the *inf-sup* stability of the discretization and the appearance of spurious pressure wiggles.)** *In the scheme considered in this section, the upwinding in the discretization of mass balance controls the onset of density oscillations. As long as the pressure and the density are linked by an increasing function, that is as long as the flow remains compressible with a reasonable equation of state, it is probably sufficient to prevent pressure oscillations. Besides, the fourth term of the left hand side of (II.45), i.e. the term involving  $|p^{n+1}|_{h, \tilde{\rho}^{n+1}}^2$ , provides a control on the discrete  $H^1$ -like semi-norm of the pressure, at least for large time steps, and therefore also produces an additional pressure smearing. However, it comes up in the analysis as the composition of the discrete divergence with the discrete gradient; consequently, one will obtain such a smoothing effect only for *inf-sup* stable discretizations. Note also that, even for steady state problems, some authors recommend the use of stable approximation space pairs to avoid pressure wiggles [10, 30].*

**Remark II.3.17 (On a different projection step)** *Some authors propose a different projection step [7, 75], which reads in the time semi-discrete setting :*

$$\left\{ \begin{array}{l} \frac{\varrho(p^{n+1}) u^{n+1} - \tilde{\rho}^{n+1} \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - \tilde{p}^{n+1}) = 0 \\ \frac{\varrho(p^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho(p^{n+1}) u^{n+1}) = 0 \end{array} \right.$$

*Considering this system, one may be tempted by the following line of thought : choosing  $q^{n+1} = \varrho(p^{n+1}) u^{n+1}$  as variable, taking the discrete divergence of the first equation and using the second one will cause the convection term of the mass balance to disappear from the discrete elliptic problem for the pressure, whatever the discretization of this term (and, in particular, the choice of the density at the edges) may be. Consequently, the equation for the pressure will be free of the nonlinearities induced by the upwinding and the dependency of the convected density on the pressure, while one still may hope to obtain a positive upwind (with respect to the density) scheme. In fact, this last point is incorrect. To be valid, it would necessitate that, from any solution  $(q^{n+1}, p^{n+1})$ , one be able to compute a velocity field  $u^{n+1}$  by dividing  $q^{n+1}$  by the density of the control volume located upstream with respect to  $u^{n+1}$ . Unfortunately, it is not always possible to obtain this upstream value ; for instance, if for two neighbouring control volumes  $K$  and  $L$ ,  $\rho_K < 0$ ,  $\rho_L > 0$  and  $q^{n+1} \cdot n_{K|L} > 0$ , neither the choice of  $K$  nor  $L$  for the upstream control volume is valid. Consequently, with this discretization, we are no longer able to guarantee the positivity of  $\rho$  or the absence of oscillations. However, as explained above, if the density remains positive, we will have a smearing of pressure or density wiggles due to the fact that the discretization is *inf-sup* stable.*

### II.3.8 Implementation

The implementation of the first three steps (II.27)-(II.29) of the numerical scheme is standard, and we therefore only describe here in details the fourth step, that is the projection step. The precise algebraic formulation of the system (II.30) reads :

$$\left\{ \begin{array}{l} \frac{1}{\delta t} M_{\tilde{\rho}^{n+1}} (\bar{u}^{n+1} - \tilde{u}^{n+1}) + B^t (p^{n+1} - \tilde{p}^{n+1}) = 0 \\ \frac{1}{\delta t} R (\varrho(p^{n+1}) - \rho^n) - B Q_{\rho^{n+1}}^{\text{up}} \bar{u}^{n+1} = 0 \end{array} \right. \quad (\text{II.49})$$

where  $M_{\tilde{\rho}^{n+1}}$  and  $Q_{\rho^{n+1}}^{\text{up}}$  are two diagonal matrices; for the first one, we recall that the entry corresponding to an edge  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$  is computed by multiplying the measure of the



diamond associated to  $\sigma$  by the predicted density (at the edge center)  $\tilde{\rho}_\sigma^{n+1}$ ; in the second one, the same entry is obtained by just taking the density at  $t^{n+1}$  in the element located upstream of  $\sigma$  with respect to  $\bar{u}^{n+1}$ , *i.e.* either  $\varrho(p_K^{n+1})$  or  $\varrho(p_L^{n+1})$ . Note that these definitions can be extended in a straightforward way for the boundary edges, if the velocity is not prescribed to zero on the boundary of the computational domain. The matrix  $\mathbf{R}$  is diagonal and, for any  $K \in \mathcal{M}$ , its entry  $\mathbf{R}_K$  is the measure of the element  $K$ . For the sake of simplicity, we suppose for the moment that the equation of state is linear :

$$\varrho(p^{n+1}) = \frac{\partial \varrho}{\partial p} p^{n+1}$$

The elliptic problem for the pressure is obtained by multiplying the first relation of (II.49) by  $\mathbf{B} \mathbf{Q}_{\rho^{n+1}}^{\text{up}} (\mathbf{M}_{\bar{\rho}^{n+1}})^{-1}$  and using the second one. This equation reads :

$$\left[ \mathbf{L} + \frac{\partial \varrho}{\partial p} \frac{1}{\delta t^2} \mathbf{R} \right] p^{n+1} = \mathbf{L} \tilde{p}^{n+1} + \frac{1}{\delta t^2} \mathbf{R} \rho^n + \frac{1}{\delta t} \mathbf{B} \mathbf{Q}_{\rho^{n+1}}^{\text{up}} \tilde{u}^{n+1} \quad (\text{II.50})$$

where  $\mathbf{L} = \mathbf{B} \mathbf{Q}_{\rho^{n+1}}^{\text{up}} (\mathbf{M}_{\bar{\rho}^{n+1}})^{-1} \mathbf{B}^t$  can be viewed, for the discretization at hand, as a finite volume discrete approximation of the Laplace operator with Neumann boundary conditions (when the velocity is prescribed at the boundary), weighted by a mesh-dependent coefficient and the densities ratio (see remarks II.3.11 and II.3.12). We recall that, by a calculation similar to the proof of lemma II.3.10, this matrix can be evaluated directly in the "finite volume way", by the following relation, valid for each element  $K$  :

$$(\mathbf{L} p^{n+1})_K = \sum_{\sigma=K|L} \frac{\rho_{\text{up},\sigma}}{\tilde{\rho}_\sigma^{n+1}} \frac{|\sigma|^2}{|D_\sigma|} (p_K^{n+1} - p_L^{n+1})$$

where  $\rho_{\text{up},\sigma}$  stands for the upwind density associated to the edge  $\sigma$ . Provided that  $p^{n+1}$  is known, the first relation of (II.49) gives us the updated value of the velocity :

$$\bar{u}^{n+1} = \tilde{u}^{n+1} - \delta t (\mathbf{M}_{\bar{\rho}^{n+1}})^{-1} \mathbf{B}^t (p^{n+1} - \tilde{p}^{n+1}) \quad (\text{II.51})$$

In order to preserve the positivity of the density, it is necessary to use the value of the density upwinded with respect to  $\bar{u}^{n+1}$  in the mass balance; therefore, equations (II.50) and (II.51) are not decoupled, in contrast with what happens in usual projection methods. We thus implement the following iterative algorithm :

Initialization :  $p_0^{n+1} = \tilde{p}^{n+1}$  and  $\bar{u}_0^{n+1} = \tilde{u}^{n+1}$

Step 4.1 – Solve for  $p_{k+1/2}^{n+1}$  :

$$\left[ \mathbf{L} + \frac{\partial \varrho}{\partial p} \frac{1}{\delta t^2} \mathbf{R} \right] p_{k+1/2}^{n+1} = \mathbf{L} \tilde{p}^{n+1} + \frac{1}{\delta t^2} \mathbf{R} \rho^n + \frac{1}{\delta t} \mathbf{B} \mathbf{Q}_{\rho^{n+1}}^{\text{up}} \tilde{u}^{n+1}$$

where the density in  $\mathbf{L}$  and  $\mathbf{Q}_{\rho^{n+1}}^{\text{up}}$  is evaluated at  $p_k^{n+1}$  and the upwinding in  $\mathbf{Q}_{\rho^{n+1}}^{\text{up}}$  is performed with respect to  $\bar{u}_k^{n+1}$

Step 4.2 – Compute  $p_{k+1}^{n+1}$  as  $p_{k+1}^{n+1} = \alpha p_{k+1/2}^{n+1} + (1 - \alpha) p_k^{n+1}$

Step 4.3 – Compute  $\bar{u}_{k+1}^{n+1}$  as :

$$\bar{u}_{k+1}^{n+1} = \tilde{u}^{n+1} - \delta t (\mathbf{M}_{\bar{\rho}^{n+1}})^{-1} \mathbf{B}^t (p_{k+1}^{n+1} - \tilde{p}^{n+1})$$

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

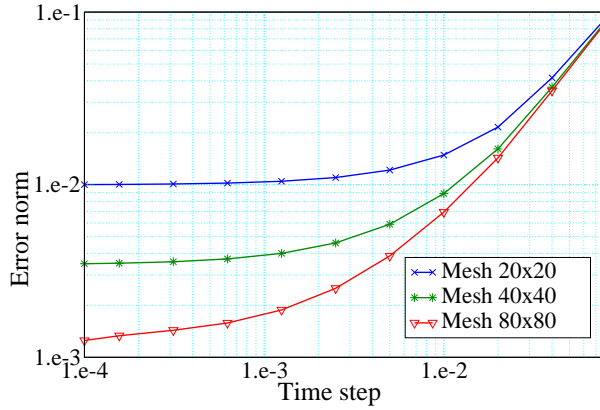


FIG. II.2 – Velocity error as a function of the time step.

The second step of the previous algorithm is a relaxation step which can be performed to ensure convergence; however, in the tests presented hereafter, we use  $\alpha = 1$  and obtain convergence in a few iterations (typically less than 5). When the equation of state is nonlinear, step 4.1 is replaced by one iteration of Newton's algorithm.

### II.3.9 Numerical experiments

In this section, we describe numerical experiments which are performed to assess the behaviour of the pressure correction scheme presented in this paper, in particular the convergence rate with respect to the space and time discretizations.

With  $\Omega = (0, 1) \times (-\frac{1}{2}, \frac{1}{2})$ , we choose for the momentum and density the following expressions :

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

These functions satisfy the mass balance equation; for the momentum balance, we add the corresponding right hand side. In this latter equation, the divergence of the stress tensor is given by :

$$\nabla \cdot \tau(u) = \mu \Delta u + \frac{\mu}{3} \nabla \nabla \cdot u, \quad \mu = 10^{-2}$$

and, in the discrete momentum balance equation (II.37), we use instead of  $\int_{\Omega, h} \tau(\tilde{u}^{n+1}) : \nabla \varphi_\sigma^{(i)} dx$  the expression :

$$\mu \int_{\Omega, h} \left[ \nabla \tilde{u}^{n+1} : \nabla \varphi_\sigma^{(i)} + \frac{1}{3} (\nabla \cdot \tilde{u}^{n+1}) (\nabla \cdot \varphi_\sigma^{(i)}) \right] dx$$

which ensures that this term is dissipative. The pressure is linked to the density by the following equation of state :

$$p = \wp(\rho) = \frac{\rho - 1}{\gamma \text{Ma}^2}, \quad \gamma = 1.4, \text{Ma} = 0.5$$

where the parameter Ma corresponds to the characteristic Mach number.

We use in these tests a special numerical integration of the forcing term of the momentum balance, which is designed to ensure that the discretization of a gradient is indeed a discrete

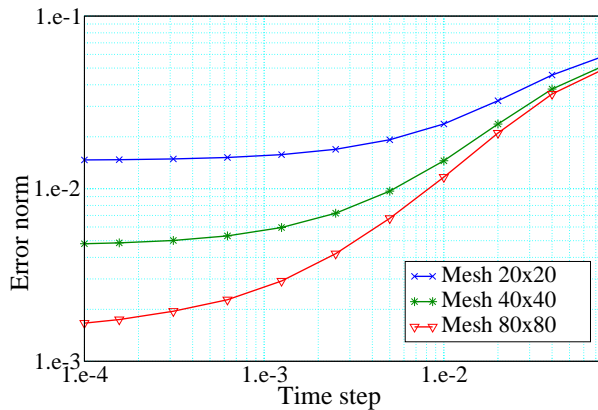


FIG. II.3 – Pressure error as a function of the time step.

gradient (*i.e.* if the forcing term  $f_v$  can be recast under the form  $f_v = \nabla g_v$ , the discrete right hand side of the momentum balance belongs to the range of  $B^t$ ).

Velocity and pressure errors obtained at  $t = 0.5$ , in respectively  $L^2$  and discrete  $L^2$  norms and as a function of the time step, are drawn on respectively figure II.2 and figure II.3, for  $20 \times 20$ ,  $40 \times 40$  and  $80 \times 80$  uniform meshes (so using the Rannacher-Turek element). For large time steps, these curves show a decrease corresponding to approximately a first order convergence in time for the velocity and the pressure, until a plateau is reached, due to the fact that errors are bounded from below by the residual spatial discretization error. For both velocity and pressure, the value of the errors on this plateau show a space convergence order (in  $L^2$  norm) between 1 and 2.

## II.4 Conclusion

We presented in this paper a numerical scheme for the barotropic Navier-Stokes compressible equations, based on a pressure-correction time stepping algorithm. For the space discretization, it combines low-order non-conforming mixed finite elements with finite volumes ; in the incompressible limit, one recovers a classical projection scheme based on an *inf-sup* stable pair of approximation spaces for the velocity and the pressure. This scheme is proven to enjoy an unconditional stability property : irrespectively of the time step, the discrete solution obeys the *a priori* estimates associated to the continuous problem, *i.e.* strict positivity of the density, bounds in  $L^\infty$ -in-time norm of the quantity  $\int_\Omega \rho u^2 dx$  and  $\int_\Omega \rho P(\rho) dx$  and in  $L^2$ -in-time norm of the viscous dissipation  $\int_\Omega \tau(u) : \nabla u dx$ . To our knowledge, this result is the first one of this type for barotropic compressible flows.

However, the scheme presented here is by no means "the ultimate scheme" for the solution to the compressible Navier-Stokes equations. It should rather be seen as an example of application (and probably one of the less sophisticated ones) of the mathematical arguments developed to obtain stability, namely theorems II.2.1 (discrete elastic potential identity) and II.3.7 (stability of the advection operator), and our hope is that these two ingredients could be used as such or adapted in the future to study other algorithms. For instance, a computation close to the proof of theorem II.3.14 (and even simpler) would yield the stability of the fully implicit scheme ; adding to this latter algorithm a prediction step for the density (as performed here) would also allow to linearize (once again as performed here) the convection operator without loss of stability. A stable

pressure-correction scheme avoiding this prediction step can also be obtained, and is currently under tests at IRSN for the computation of compressible bubbly flows. Besides these variants, less diffusive schemes should certainly be sought. Finally, the proposed scheme is currently the object of more in-depth numerical studies including, in particular, problems admitting less smooth solutions than the test presented here.

## Chapitre III

# On a discretization of phases mass balance in segregated algorithms for the drift-flux model

**Abstract.** We address in this paper a parabolic equation used to model the phases 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 non-linear term of the form  $\nabla \cdot \rho \varphi(y) u_r$ , where  $y$  is the unknown mass fraction,  $\rho$  stands for the density,  $\varphi(\cdot)$  is a regular function such that  $\varphi(0) = \varphi(1) = 0$  and  $u_r$  is a (non-necessarily divergence free) velocity field. We propose a finite-volume scheme for the numerical approximation of this equation, with a discretization of the non-linear term based on monotone flux functions [26]. Under the classical assumption [51] 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.

### III.1 Introduction

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

$$\left\{ \begin{array}{l} \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \\ \frac{\partial \rho u}{\partial t} + \nabla \cdot (\rho u \otimes u) + \nabla p - \nabla \cdot \tau = f \\ \frac{\partial \rho y}{\partial t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \varphi(y) u_r) = \nabla \cdot (D \nabla y) \\ \rho = \rho(p, y) \end{array} \right. \quad (\text{III.1})$$

where  $t$  stands for the time,  $u$  for the fluid velocity,  $p$  for the pressure,  $\rho$  for the fluid density. The tensor  $\tau$  is the viscous part of the stress tensor, given by the following expression :

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

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

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

and, in this case, this term is dissipative (*i.e.* for any regular velocity field  $u$  vanishing on the boundary, the integral of  $\nabla \cdot \tau(u) \cdot u$  over the computational domain is non-negative). The function  $\varrho(\cdot, \cdot)$ , which gives the density as an explicit function of  $y$  and the pressure, is the equation of state of the considered fluid. The nonlinear function  $\varphi(\cdot)$  is such that  $\varphi \in C^1([0, 1], \mathbb{R})$  and  $\varphi(0) = \varphi(1) = 0$ , and, since  $y$  is supposed for physical reasons to satisfy  $0 \leq y \leq 1$ ,  $\varphi(\cdot)$  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 supposed to be posed over  $\Omega$ , an open bounded connected subset 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  $\rho$ ,  $u$  and  $y$ .

Several physical problems enter this framework. For instance, taking for  $y$  the gas mass fraction, and for  $\rho$  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  $\varphi(\cdot)$  is given by  $\varphi(y) = y(1 - y)$ . 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 (III.1), one faces at least two difficulties. First, the unknown  $y$  can be expected, from both physical and mathematical reasons, to remain in the interval  $[0, 1]$ . This suggests to build a numerical scheme, here a finite volume scheme, which reproduces these features at the discrete level. This is performed by combining a monotone flux approach [26, section 21] for the term  $\nabla \cdot (\rho \varphi(y) u_r)$  with the argument introduced by Larrouturou [51] : 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 technique involving an auxiliary linear dual problem.

Second, even if the model at hand represents a compressible flow, in fact the liquid is almost incompressible, so that zones may appear in the flow where the velocity of acoustic waves is very large, and the Mach number 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. To this purpose, we use a fractional step algorithm issued from the incompressible flow numerics, namely from the class of finite element projection methods. For a description of this kind of schemes for incompressible flow, see *e.g.* [38, 53] and references herein; an extension to barotropic Navier-Stokes equations close to the scheme developed here can be found in [31], 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 Babuska-Brezzi condition (*e.g.* [9]). 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 a 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 built as follows. The finite volume scheme for the advection-diffusion of  $y$  is first described and analysed in section III.2, then the fractional step algorithm for the solution of the whole problem is presented in section III.3. Numerical tests are reported in section III.4; first a problem exhibiting an analytical solution allows to assess convergence properties of the discretization, then two physical situations are addressed, first a phase separation problem under gravity, then the flow of a sedimenting dilute suspension in a channel flow with an obstacle.

## III.2 Discretization for the nonlinear advection-diffusion equation

In this section, we consider the balance equation of the independent unknown field  $y$ , which reads in the semi-discrete form, with a first order backward Euler time discretization :

$$\frac{\rho y - \rho^* y^*}{\delta t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \varphi(y) u_r^*) = \nabla \cdot (D \nabla y) \quad \text{in } \Omega \times (0, T) \quad (\text{III.4})$$

where the density fields  $\rho$  and  $\rho^*$ , the mass fraction  $y^*$  and the field  $u_r^*$  are here supposed to be known quantities. In addition, for the sake of simplicity, we assume in this section that both the velocity field and  $u_r^*$  vanish on the boundary  $\partial\Omega$  of the computational domain and that the fluid density satisfies the mass balance equation, *i.e.* that the density and velocity fields satisfy, in the semi-discrete time setting :

$$\frac{\rho - \rho^*}{\delta t} + \nabla \cdot (\rho u) = 0 \quad (\text{III.5})$$

Again to simplify the presentation, we assume that  $y$  obeys a homogeneous Neumann condition on the whole boundary; however, it is clear from the subsequent developments that similar results would hold if  $y$  was prescribed on the boundary, provided that the prescribed value lies in the interval  $[0, 1]$ .

We begin this section by describing the considered space discretization and precisely setting the discrete problem at hand. An admissible finite volume mesh of  $\Omega$  involves :

- (i) a family  $\mathcal{M}$  of control volumes, which are convex disjoint polygons ( $d = 2$ ) or polyhedrons ( $d = 3$ ) included in  $\Omega$  such that  $\bar{\Omega} = \bigcup_{K \in \mathcal{M}} \bar{K}$ .
- (ii) a family  $\mathcal{E}$  of bounded subsets of hyperplanes of  $\mathbb{R}^d$  included in  $\bar{\Omega}$ , which are the edges ( $d = 2$ ) or faces ( $d = 3$ ) of the control volumes. The set of edges or faces included in the boundary of  $\Omega$  is denoted by  $\mathcal{E}_{\text{ext}}$  and the set of internal ones (*i.e.*  $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ ) is denoted by  $\mathcal{E}_{\text{int}}$ . If  $K, L \in \mathcal{M}$ , we suppose that either  $\bar{K} \cap \bar{L} = \emptyset$  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  $K|L$ .
- (iii) a family  $\mathcal{P} = (x_K)_{K \in \mathcal{M}}$  of points of  $\Omega$  such that  $x_K \in \bar{K}$  for all  $K \in \mathcal{M}$  and, if  $\sigma = K|L$ ,  $x_K \neq x_L$  and the straight line going through  $x_K$  and  $x_L$  is orthogonal to  $\sigma$ .

The following notations are used hereafter. The set of edges ( $d = 2$ ) or faces ( $d = 3$ ) of an element  $K$  of  $\mathcal{M}$  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 of  $\sigma$ , 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  $\sigma$ , respectively. For any  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}(K)$ , we denote by  $d_{K,\sigma}$  the Euclidean distance between  $x_K$  and  $\sigma$ . 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  $d_\sigma$  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  $X_{\mathcal{M}}$  the space of functions which are piecewise constant on each control volume  $K \in \mathcal{M}$ .

Throughout this paper, for any real number  $a$ , we define  $a^+ = \max(a, 0)$  and  $a^- = -\min(a, 0)$ , so  $a = a^+ - a^-$  and both  $a^+$  and  $a^-$  are non-negative.

With the previous notations, the discrete problem considered in this section reads :

$$\left| \begin{array}{l} \forall K \in \mathcal{M}, \\ \frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} [F_{\sigma,K}^+ y_K - F_{\sigma,K}^- y_L] \\ + \sum_{\sigma=K|L} [G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K)] + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (y_K - y_L) = 0 \end{array} \right. \quad (\text{III.6})$$

where the following notations are used.

The quantity  $F_{\sigma,K}$ , associated to the edge  $\sigma$  and to the control volume  $K$ , is such that, for any internal edge  $\sigma = K|L$ ,  $F_{\sigma,K} = -F_{\sigma,L}$ . In addition, we suppose that  $(\rho_K)_{K \in \mathcal{M}}$ ,  $(\rho_K^*)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  satisfy the following condition :

$$\forall K \in \mathcal{M}, \quad \rho_K > 0, \quad \rho_K^* > 0 \quad \text{and} \quad \frac{|K|}{\delta t} (\rho_K - \rho_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} = 0 \quad (\text{III.7})$$

This abstract setting is sufficient for the developments of this section ; however, equation (III.7) may be interpreted as a finite volume discretization of the mass balance (III.5), where  $\rho_K^*$  and  $\rho_K$  stands for the approximation of the density in  $K$  at two consecutive time steps,  $F_{\sigma,K}$  for the mass flux from  $K$  through the edge  $\sigma$ .

The quantity  $G_{\sigma,K}$  is given by :

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

where  $\rho_\sigma$  stands either for  $\rho_\sigma = \frac{1}{2}(\rho_K + \rho_L)$  (centered choice), or for  $\rho_\sigma = \rho_K$  if  $F_{\sigma,K} \geq 0$  and  $\rho_\sigma = \rho_L$  otherwise (upwind choice).

Finally, the function  $g(\cdot, \cdot)$  is called a numerical monotone flux function and is defined in the following way (see [26] for the theory and some examples).

### Definition III.2.1

*[Numerical monotone flux function] Let the function  $g(\cdot, \cdot) \in C(\mathbb{R}^2, \mathbb{R})$  satisfy the following assumptions :*

1.  $g(a, b)$  is non-decreasing with respect to  $a$  and non-increasing with respect to  $b$ , for any real numbers  $a$  and  $b$ ,
2.  $g(\cdot, \cdot)$  is Lipschitz continuous with respect to both variables over  $\mathbb{R}$ ,
3.  $g(a, a) = \varphi(a)$ , for any  $a \in \mathbb{R}$ .

*Then  $g(\cdot, \cdot)$  is said to be a numerical monotone flux function.*

The result proven in this section is the following.



**Theorem III.2.2 (Existence, uniqueness and  $L^\infty$  bounds for a discrete solution)**

Let us suppose that  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  are such that the condition (III.7) is satisfied and that  $g(\cdot, \cdot)$  is a numerical monotone flux function such that  $\varphi(x) = g(x, x)$  vanishes for  $x \leq 0$  and  $x \geq 1$ . Then, if  $y_K^* \in [0, 1]$ ,  $\forall K \in \mathcal{M}$ , there exists a unique solution to the considered discrete problem (III.6), which verifies  $y_K \in [0, 1]$ ,  $\forall K \in \mathcal{M}$ .

This theorem summarizes a series of lemmata, which are exposed in the following subsections : first (section III.2.1), we prove an *a priori*  $L^\infty$  estimate for  $y$ , precisely speaking 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 III.2.2); finally, this latter is shown to be unique (section III.2.3).

**III.2.1 An  $L^\infty$  stability property**

From a physical point of view, for instance thinking of the field  $y$  as a mass fraction, it seems natural for  $y$  to satisfy an “ $L^\infty$  stability property”, more precisely speaking to remain at any time in the  $[0, 1]$  interval. The aim of this section is to prove that this property holds for the solution of the scheme (III.6), provided that (III.7) holds at each time step and that the initial condition for  $y$  takes its values in  $[0, 1]$ .

Let us first review the proof for the continuous problem :

$$\frac{\partial \rho y}{\partial t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \varphi(y) u_r) = \nabla \cdot (D \nabla y)$$

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

$$-\int_{\Omega} \frac{\partial \rho y}{\partial t} y^- - \int_{\Omega} \nabla \cdot (\rho y u) y^- - \int_{\Omega} \nabla \cdot (\rho \varphi(y) u_r) y^- - \int_{\Omega} D \nabla y \cdot \nabla y^- = 0 \quad (\text{III.8})$$

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

$$T_{\text{adv}} = -\int_{\Omega} \frac{\partial \rho y}{\partial t} y^- - \int_{\Omega} \nabla \cdot (\rho y u) y^-$$

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

$$T_{\text{adv}} = \int_{\Omega} \left[ \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) \right] [y^-]^2 - \int_{\Omega} \left[ \rho y^- \frac{\partial y}{\partial t} + (\rho u y^-) \cdot \nabla y \right]$$

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

$$-\int_{\Omega} \left[ \rho y^- \frac{\partial y}{\partial t} + (\rho u y^-) \cdot \nabla y \right] = \frac{1}{2} \int_{\Omega} \left[ \rho \frac{\partial (y^-)^2}{\partial t} + (\rho u) \cdot \nabla (y^-)^2 \right]$$

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

$$\begin{aligned} T_{\text{adv}} &= \frac{1}{2} \int_{\Omega} \left[ \rho \frac{\partial (y^-)^2}{\partial t} - (y^-)^2 \nabla \cdot (\rho u) \right] \\ &= \frac{1}{2} \int_{\Omega} \left[ \rho \frac{\partial (y^-)^2}{\partial t} + (y^-)^2 \frac{\partial \rho}{\partial t} \right] = \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho (y^-)^2 \end{aligned} \quad (\text{III.9})$$

Substituting the term  $T_{\text{adv}}$  in the relation (III.8) yields :

$$\frac{1}{2} \frac{\partial}{\partial t} \int_{\Omega} \rho (y^-)^2 - \int_{\Omega} \nabla \cdot (\rho \varphi(y) u_r) y^- + \int_{\Omega} D |\nabla y^-|^2 = 0$$

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

$$\frac{1}{2} \frac{\partial}{\partial t} \int_{\Omega} \rho (y^-)^2 = -D \int_{\Omega} (\nabla y^-)^2 \leq 0$$

Thus  $y$  is non-negative, provided that the initial condition for  $y$  is non-negative. Considering the equation satisfied by  $y' = 1 - y$ , 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 to relation (III.9) ; this is achieved by the following lemma. As the mass balance plays a central role in the continuous setting, it is natural that the condition (III.7) also does in the discrete one, which indeed is the case.

**Lemma III.2.3**

Let  $(\rho_K)_{K \in \mathcal{M}}$ ,  $(\rho_K^*)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  be such that the condition (III.7) is satisfied. Then the following stability property holds :

$$- \sum_{K \in \mathcal{M}} y_K^- \left[ \frac{|K|}{\delta t} (\rho_K y_K - \rho_K^* y_K^*) + \sum_{\sigma=K|L} [F_{\sigma,K}^+ y_K - F_{\sigma,K}^- y_L] \right] \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* ((y_K^*)^-)^2] \quad (\text{III.10})$$

**Proof.**

To make the following computations easier, we introduce the following notation :

$$F_{\sigma,K}^+ y_K - F_{\sigma,K}^- y_L = F_{\sigma,K} y_{\sigma} \quad \text{with } y_{\sigma} = y_K \text{ if } F_{\sigma,K} \geq 0 \text{ and } y_{\sigma} = y_L \text{ otherwise}$$

The left-hand side of (III.10) may be written as :

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

where  $T_1$  and  $T_2$  read :

$$T_1 = - \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} y_K^- (\rho_K y_K - \rho_K^* y_K^*), \quad T_2 = - \sum_{K \in \mathcal{M}} y_K^- \left[ \sum_{\sigma=K|L} F_{\sigma,K} y_{\sigma} \right]$$

To express the first term, we first use the fact that, when  $y_K^-$  is non-zero,  $y_K = -y_K^-$ , then we split it as :

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

We thus have :

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

As the last term is always non-negative, we get :

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

Then, expanding the last term by the identity  $2a(a-b) = a^2 + (a-b)^2 - b^2$ , we obtain :

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

The last sum is always non-negative, and we finally get for  $T_1$  :

$$T_1 \geq \underbrace{\sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*)}_{T_{1,1}} + \frac{1}{2} \underbrace{\sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \rho_K^* [(y_K^-)^2 - ((y_K^*)^-)^2]}_{T_{1,2}} \quad (\text{III.11})$$

We now turn to  $T_2$  :

$$T_2 = \underbrace{\sum_{K \in \mathcal{M}} (y_K^-)^2 \left[ \sum_{\sigma=K|L} F_{\sigma,K} \right]}_{T_{2,1}} - \underbrace{\sum_{K \in \mathcal{M}} y_K^- \left[ \sum_{\sigma=K|L} F_{\sigma,K} (y_\sigma + y_K^-) \right]}_{T_{2,2}} \quad (\text{III.12})$$

Expanding the quantity  $y_K^- (y_\sigma + y_K^-)$  by the identity  $2a(a+b) = a^2 - b^2 + (a+b)^2$ , the second term reads :

$$-T_{2,2} = -\frac{1}{2} \sum_{K \in \mathcal{M}} (y_K^-)^2 \sum_{\sigma=K|L} F_{\sigma,K} - \frac{1}{2} \underbrace{\sum_{K \in \mathcal{M}} \left[ \sum_{\sigma=K|L} F_{\sigma,K} [(y_K^- + y_\sigma)^2 - y_\sigma^2] \right]}_{T_{2,3}}$$

Reordering the sum in the last term and using the property  $F_{\sigma,K} = -F_{\sigma,L}$ , we have :

$$T_{2,3} = -\frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} F_{\sigma,K} [(y_K^- + y_\sigma)^2 - (y_L^- + y_\sigma)^2]$$

Supposing without loss of generality that we have chosen for the edge  $\sigma = K|L$  the orientation such that  $F_{\sigma,K} \geq 0$ , we get, as  $y_\sigma = y_K$  :

$$T_{2,3} = - \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} F_{\sigma,K} \begin{cases} \left[ -\frac{1}{2}(y_L^- + y_K)^2 \right] & \text{if } y_K \leq 0 \\ \left[ -(y_L^-)^2 - 2y_K y_L^- \right] & \text{otherwise} \end{cases}$$

thus  $T_{2,3}$  is non-negative. Hence, by equation (III.7), we have :

$$-T_{2,2} \geq -\frac{1}{2} \sum_{K \in \mathcal{M}} (y_K^-)^2 \sum_{\sigma=K|L} F_{\sigma,K} = \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) \quad (\text{III.13})$$

Summing (III.11) and (III.12), remarking that  $T_{2,1}$  cancels with  $T_{1,1}$  by equation (III.7) and using (III.13), we finally obtain :

$$\begin{aligned} T &\geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \rho_K^* [(y_K^-)^2 - ((y_K^*)^-)^2] + \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} (y_K^-)^2 (\rho_K - \rho_K^*) \\ &= \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K (y_K^-)^2 - \rho_K^* ((y_K^*)^-)^2] \end{aligned}$$

which concludes the proof.  $\blacksquare$

We are now in position to prove that  $y$  remains in the interval  $[0, 1]$ ; this result is given in the following two lemmas.

**Lemma III.2.4 (Non-negativity of  $y$ )**

Let  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  be such that the condition (III.7) is satisfied and  $g(\cdot, \cdot)$  be a numerical monotone flux function such that  $\varphi(x) = g(x, x)$  vanishes for  $x \leq 0$ . Then, if  $y_K^* \geq 0$ ,  $\forall K \in \mathcal{M}$ , the discrete solution of (III.6) also verifies  $y_K \geq 0$ ,  $\forall K \in \mathcal{M}$ .

**Proof.**

As in the continuous case, the starting point is to multiply the equation by  $-y^-$ , which, in the discrete case, consists in multiplying relation (III.6) by  $-y_K^-$  and summing over the control volumes. We get :

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

The previous relation can be written as  $T_{\text{adv}} + T_{\text{nl}} + T_{\text{dif}} = 0$  with :

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

By virtue of lemma III.2.3, the first term can be estimated as follows :

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

Reordering the sum in the term  $T_{\text{nl}}$ , we have :

$$T_{\text{nl}} = \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} -y_K^- [G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K)] - y_L^- [G_{\sigma,L}^+ g(y_L, y_K) - G_{\sigma,L}^- g(y_K, y_L)]$$

Supposing, without loss of generality, that we have choosen for the edge  $\sigma = K|L$  the orientation such that  $G_{\sigma,K} \geq 0$ , so  $G_{\sigma,K}^+ = G_{\sigma,L}^- = G_{\sigma,K}$  and  $G_{\sigma,K}^- = G_{\sigma,L}^+ = 0$ , we get :

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

If both  $y_K$  and  $y_L$  are non-negative, this term vanishes. If  $y_L \leq 0$ , as  $g(y_L, y_L) = \varphi(y_L) = 0$ , we have :

$$T_{\text{nl}} = \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} -G_{\sigma,K} (y_K^- - y_L^-) [g(y_K, y_L) - g(y_L, y_L)]$$

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

$$T_{\text{nl}} = \sum_{\sigma \in \mathcal{E}_{\text{int}} (\sigma=K|L)} -G_{\sigma,K} (y_K^- - y_L^-) [g(y_K, y_L) - g(y_K, y_K)]$$

which is also non-negative, since  $g(\cdot, \cdot)$  is non-increasing with respect to the second argument. Let us now turn to the third term. Reordering the sum, we have :

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

which is also non-negative. Finally, we have :

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

and thus, if  $(y_K^*)_K \geq 0$ ,  $\forall K \in \mathcal{M}$  :

$$\frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} \rho_K (y_K^-)^2 \leq 0$$

which concludes the proof. ■

### Lemma III.2.5 ( $y$ bounded by 1)

Let  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  be such that the condition (III.7) is satisfied and  $g(\cdot, \cdot)$  be a numerical monotone flux function such that  $\varphi(x) = g(x, x)$  vanishes for  $x \geq 1$ . Then, if  $y_K^* \leq 1$ ,  $\forall K \in \mathcal{M}$ , the discrete solution of (III.6) also verifies  $y_K \leq 1$ ,  $\forall K \in \mathcal{M}$ .

#### Proof.

Let us consider the equation verified by  $1 - y$ . Relation (III.6) equivalently reads :

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

By the condition (III.7), the right-hand side of this equation vanishes. In addition, let  $g'(\cdot, \cdot)$  be the function defined by  $g'(1-a, 1-b) = -g(a, b)$ . This function is non-decreasing with respect to the first variable and non-increasing with respect to the second one. Moreover,  $g'(x, x) = \varphi'(x) = \varphi(1-x)$  vanishes for  $x \leq 0$ , as  $\varphi(x)$  vanishes for  $x \geq 1$ . Thus the assumptions of the preceding proposition hold and,  $\forall K \in \mathcal{M}$ ,  $(1-y)_K$  is non-negative, which concludes the proof. ■

**Remark III.2.6** ( $y > 0$  or  $y < 1$ ) *These results can be slightly sharpened in the following way. Let us suppose that,  $\forall K \in \mathcal{M}$ ,  $y_K^* > 0$ . Then  $y$  satisfies the strict inequality  $\forall K \in \mathcal{M}$ ,  $y_K > 0$ . To prove this result, let us suppose that there exists  $K \in \mathcal{M}$  such that  $y_K = 0$ , and show that this assumption leads to a contradiction. Replacing  $y_K$  by zero in the equation (III.6) of the scheme, we get :*

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

The first term is by assumption negative, while, since  $y \geq 0$ , the second and last ones are non-positive. Since the function  $s \mapsto g(s, 0)$  is increasing and  $g(0, 0) = 0$ ,  $g(y_L, 0) \geq 0$ ; from a similar argument,  $g(0, y_L) \leq 0$ , and the third term also is non-positive, which contradicts the fact that the whole sum vanishes. By applying this result to  $1-y$ , we similarly prove that, if  $y^* < 1$ ,  $y < 1$ .

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.

### III.2.2 Existence for the approximate solution

The existence of a solution to the scheme (III.6) is obtained through a so-called “topological degree” argument. For the sake of completeness, we recall this result in the following theorem (see [22, chapter 5] for an exposition of the theory and [25, 31] for other uses for the same objective as here, namely the proof of existence of a solution to a numerical scheme).

#### **Theorem III.2.7 (Application of the topological degree, finite dimensional case)**

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

- (i)  $F(\cdot, 1) = f(\cdot)$ ;
- (ii)  $\forall \alpha \in [0, 1]$ , if  $v$  is such that  $F(v, \alpha) = b$  then  $v \in W$ , where  $W$  is defined as follows :

$$W = \{v \in V \text{ such that } \|v\| < R\}$$

where  $R$  is a positive real number independent of  $\alpha$  and  $\|\cdot\|$  is a norm defined over  $V$  ;  
 (iii) the topological degree of  $F(\cdot, 0)$  with respect to  $b$  and  $W$  is equal to  $d_0 \neq 0$ .  
 Then the topological degree of  $F(\cdot, 1)$  with respect to  $b$  and to  $W$  is also equal to  $d_0 \neq 0$  ;  
 consequently, there exists at least a solution  $v \in W$  such that  $f(v) = 0$ .

**Lemma III.2.8 (Existence of a discrete solution)**

Let us suppose that  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  are such that the condition (III.7) is satisfied and that  $g(\cdot, \cdot)$  is a numerical monotone flux function such that  $\varphi(x) = g(x, x)$  vanishes for  $x \leq 0$  and  $x \geq 1$ . Then, if  $y_K^* \in [0, 1]$ ,  $\forall K \in \mathcal{M}$ , there exists a solution to the considered discrete problem (III.6).

**Proof.**

Here  $f(v) = 0$  represents the nonlinear system (III.6) and we are going to build a function  $F(\cdot, \cdot)$  suitable to apply theorem III.2.7. Such an application  $F : X_{\mathcal{M}} \times [0, 1] \rightarrow V$  is given by  $F(y, \alpha) = q_K$ ,  $\forall K \in \mathcal{M}$ , with :

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

The function  $F(\cdot, \cdot)$  is continuous from  $X_{\mathcal{M}} \times [0, 1]$  to  $X_{\mathcal{M}}$ . In addition, it is easy to see that propositions III.2.4 and III.2.5 apply to the solution to the equation  $F(y, \alpha) = 0$ , for  $0 \leq \alpha \leq 1$ . Hence, any solution to this equation belongs to :

$$W = \{y \in X_{\mathcal{M}}, \text{ such that } \max_{K \in \mathcal{M}} y_K < 2\}$$

Moreover, because of the estimate  $y_K \leq 1$ ,  $\forall K \in \mathcal{M}$ , the linear system  $F(y, 0) = 0$  has a unique solution, which belongs to  $W$ ; the topological degree of  $F(\cdot, 0)$  with respect to  $W$  and 0 is thus either 1 or  $-1$  (precisely speaking, the sign of the determinant of the matrix associated to the linear system  $F(y, 0) = 0$ ). Applying the theorem III.2.7, the proof is complete.  $\blacksquare$

**III.2.3 Uniqueness of the approximate solution**

The uniqueness of the solution to the scheme (III.6) is obtained through a dual problem method. First, we introduce this technique in the semi-discrete time setting. Let  $y$  and  $y'$  be two solutions of the problem (III.4) satisfying :

$$\frac{\rho y - \rho^* y^*}{\delta t} + \nabla \cdot (\rho y u) + \nabla \cdot (\rho \varphi(y) u_r) = \nabla \cdot (D \nabla y) \quad (\text{III.14})$$

and

$$\frac{\rho y' - \rho^* y^*}{\delta t} + \nabla \cdot (\rho y' u) + \nabla \cdot (\rho \varphi(y') u_r) = \nabla \cdot (D \nabla y') \quad (\text{III.15})$$

Then, subtracting (III.15) from (III.14), we have :

$$\frac{\rho \delta y}{\delta t} + \nabla \cdot (\rho \delta y u) + \nabla \cdot \left( \rho \frac{\varphi(y) - \varphi(y')}{\delta y} \delta y u_r \right) = \nabla \cdot (D \nabla \delta y)$$

where  $\delta y$  is given by  $\delta y = y - y'$ . Multiplying by a test function  $\psi$  and integrating on  $\Omega$  yields, for all  $\psi$  :

$$\frac{1}{\delta t} \int_{\Omega} \rho \delta y \psi + \int_{\Omega} \nabla \cdot (\rho \delta y u) \psi + \int_{\Omega} \nabla \cdot \left( \rho \frac{\varphi(y) - \varphi(y')}{\delta y} \delta y u_r \right) \psi + D \int_{\Omega} \nabla \delta y \cdot \nabla \psi = 0 \quad (\text{III.16})$$

Then, we define the following dual problem :

$$\forall \psi, \quad U'(\bar{y}, \psi) = \int_{\Omega} \delta y \psi \quad (\text{III.17})$$

where  $U'(\bar{y}, \psi)$  is given by :

$$U'(\bar{y}, \psi) = \frac{1}{\delta t} \int_{\Omega} \rho \bar{y} \psi + \int_{\Omega} \nabla \cdot (\rho \psi u) \bar{y} + \int_{\Omega} \nabla \cdot \left( \rho \frac{\varphi(y) - \varphi(y')}{\delta y} \psi u_r \right) \bar{y} + D \int_{\Omega} \nabla \bar{y} \cdot \nabla \psi$$

Under some regularity assumptions, the dual problem (III.17) is known to satisfy the maximum principle (e.g. [36, 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 :

$$\frac{1}{\delta t} \int_{\Omega} \rho \bar{y} \delta y + \int_{\Omega} \nabla \cdot (\rho \delta y u) \bar{y} + \int_{\Omega} \nabla \cdot \left( \rho \frac{\varphi(y) - \varphi(y')}{\delta y} \delta y u_r \right) \bar{y} + D \int_{\Omega} \nabla \bar{y} \cdot \nabla \delta y = \int_{\Omega} (\delta y)^2$$

But, by equation (III.16), the left-hand side of the previous relation is equal to zero. Thus, we have :

$$\int_{\Omega} (\delta y)^2 = 0$$

Thereby  $\delta y = 0$ , and we conclude to the uniqueness of the solution. The proof that we give for the discrete problem is adapted from this technique.

**Lemma III.2.9 (Uniqueness of the discrete solution)**

Let us suppose that  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  are such that the condition (III.7) is satisfied and that  $g(\cdot, \cdot)$  is a numerical monotone flux function. Then there exists a unique solution  $y \in X_{\mathcal{M}}$  to the discrete equation (III.6).

**Proof.**

Let  $y$  and  $y'$  be two solutions of (III.6). Then, substracting the equation verified by  $y$  and  $y'$ , we get for all  $K \in \mathcal{M}$  :

$$\begin{aligned} \frac{|K|}{\delta t} \rho_K (y_K - y'_K) + \sum_{\sigma=K|L} [F_{\sigma,K}^+ (y_K - y'_K) - F_{\sigma,K}^- (y_L - y'_L)] \\ + \sum_{\sigma=K|L} G_{\sigma,K}^+ [g(y_K, y_L) - g(y'_K, y'_L)] - G_{\sigma,K}^- [g(y_L, y_K) - g(y'_L, y'_K)] \\ + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_{\sigma}} [(y_K - y'_K) - (y_L - y'_L)] = 0 \end{aligned}$$

Extending the definition of the quotient  $[g(a, \cdot) - g(b, \cdot)] / [a - b]$  and  $[g(\cdot, a) - g(\cdot, b)] / [a - b]$  to the case  $a = b$ , which is possible because  $g(\cdot, \cdot)$  is supposed to be Lipschitz continuous, the nonlinear term can be recast to obtain :

$$\begin{aligned} \frac{|K|}{\delta t} \rho_K (y_K - y'_K) + \sum_{\sigma=K|L} [F_{\sigma,K}^+ (y_K - y'_K) - F_{\sigma,K}^- (y_L - y'_L)] \\ + \sum_{\sigma=K|L} G_{\sigma,K}^+ \left[ \frac{g(y_K, y_L) - g(y'_K, y_L)}{y_K - y'_K} (y_K - y'_K) + \frac{g(y'_K, y_L) - g(y'_K, y'_L)}{y_L - y'_L} (y_L - y'_L) \right] \\ + \sum_{\sigma=K|L} -G_{\sigma,K}^- \left[ \frac{g(y_L, y_K) - g(y'_L, y_K)}{y_L - y'_L} (y_L - y'_L) + \frac{g(y'_L, y_K) - g(y'_L, y'_K)}{y_K - y'_K} (y_K - y'_K) \right] \\ + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_{\sigma}} [(y_K - y'_K) - (y_L - y'_L)] = 0 \end{aligned} \quad (\text{III.18})$$



We now introduce the following discrete variational problem :

$$\begin{aligned}
 & \text{Find } \bar{y} \in X_{\mathcal{M}} \text{ such that, } \forall \psi \in X_{\mathcal{M}}, \quad \sum_{K \in \mathcal{M}} U'_K(\bar{y}, \psi) = \sum_{K \in \mathcal{M}} H_K(\psi) \quad \text{where :} \\
 & U'_K(\bar{y}, \psi) = \frac{|K|}{\delta t} \rho_K \bar{y}_K \psi_K + \bar{y}_K \sum_{\sigma=K|L} [F_{\sigma,K}^+ \psi_K - F_{\sigma,K}^- \psi_L] \\
 & \quad + \bar{y}_K \sum_{\sigma=K|L} G_{\sigma,K}^+ \left[ \frac{g(y_K, y_L) - g(y'_K, y_L)}{y_K - y'_K} \psi_K + \frac{g(y'_K, y_L) - g(y'_K, y'_L)}{y_L - y'_L} \psi_L \right] \\
 & \quad + \bar{y}_K \sum_{\sigma=K|L} -G_{\sigma,K}^- \left[ \frac{g(y_L, y_K) - g(y'_L, y_K)}{y_L - y'_L} \psi_L + \frac{g(y'_L, y_K) - g(y'_L, y'_K)}{y_K - y'_K} \psi_K \right] \\
 & \quad + D \bar{y}_K \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} [\psi_K - \psi_L] \\
 & H_K(\psi) = |K| (y_K - y'_K) \psi_K
 \end{aligned} \tag{III.19}$$

Let us suppose for the moment that such an  $\bar{y}$  exists. Then taking as test function  $\psi = y - y'$  in this variational problem, we get by equation (III.18) :

$$\sum_{K \in \mathcal{M}} |K| (y_K - y'_K)^2 = \sum_{K \in \mathcal{M}} U'_K(\bar{y}, y - y') = 0$$

In order to conclude, we thus have to prove that there exists a unique solution to the dual problem (III.19). To this purpose, we prove that the discrete operator  $A$  associated to this problem is an M-matrix, *i.e.* that the dual problem satisfies a discrete maximum principle. 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 extra-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 :

$$A_{KK} = \frac{|K|}{\delta t} \rho_K + \sum_{\sigma=K|L} F_{\sigma,K}^+ + G_{\sigma,K}^+ \frac{g(y_K, y_L) - g(y'_K, y_L)}{y_K - y'_K} - G_{\sigma,K}^- \frac{g(y'_L, y_K) - g(y'_L, y'_K)}{y_K - y'_K} + D \frac{|\sigma|}{d_\sigma}$$

As the density is positive and  $g(\cdot, \cdot)$  is non-decreasing with respect to its first argument and non-increasing with respect to its second argument,  $A_{KK}$  is positive. Non-diagonal entries on the same line are given by :

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

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

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

which is positive, since  $\rho_K > 0$ . Thus  $A$  is an M-matrix and problem (III.19) admits a unique solution, which completes the proof.  $\blacksquare$

**Remark III.2.10** *Let us introduce the (non-linear) function  $\mathcal{F}(\cdot)$  from  $X_{\mathcal{M}}$  to  $X_{\mathcal{M}}$  corresponding to the equations of the scheme, i.e. such that the (non-linear) system (III.6) is equivalent to  $\mathcal{F}(y) = 0$ . By rather sophisticated arguments, the solution of the dual problem  $\bar{y}$  introduced in the above proof can be proven to be bounded in the discrete  $H^1$  norm usual in the finite volume setting. We thus obtain a stability relation for the operator defined by  $\delta\mathcal{F}(y, y') = \mathcal{F}(y) - \mathcal{F}(y')$ , which could be used for the error analysis of finite volume schemes applied to elliptic equations involving non-linear terms of the same form as here; this work is in progress.*

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

In this section, we address the solution of the full system (III.1). We begin by describing the proposed scheme in the semi-discrete time setting; this allows us to describe how the density is discretized in time, which is one of the difficulties of this development. Then the spatial discretization for each step of the algorithm are given. We conclude with an overview of the scheme and of its properties.

From this point, we will suppose that the function  $\varphi(\cdot)$  reads  $\varphi(x) = x(1 - x)$ , which is the form taken by drift terms in balance equations governing dispersed two-phase flows; however, the following developments apply to any function  $\varphi(\cdot)$  satisfying the assumptions of theorem III.2.2.

#### III.3.1 The semi-discretized in time 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 uniform. Let  $\delta t$  be the constant time step  $\delta t = t_{n+1} - t_n$  for  $n = 0, 1, \dots, N - 1$ . A numerical scheme for the solution of the full system (III.1) is obtained by complementing the scheme presented in the preceding section by an incremental projection method. Writing this algorithm in a semi-discrete time setting, this yields the following three steps scheme :

1 - solve for  $y^{n+1}$  :

$$\begin{aligned} \frac{\rho^n y^{n+1} - \rho^{n-1} y^n}{\delta t} + \nabla \cdot (\rho^n y^{n+1} u^n) \\ + \nabla \cdot (\rho^n y^{n+1} (1 - y^{n+1}) u_r^n) - \nabla \cdot (D \nabla y^{n+1}) = 0 \end{aligned} \quad (\text{III.20})$$

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

$$\frac{\rho^n \tilde{u}^{n+1} - \rho^{n-1} u^n}{\delta t} + \nabla \cdot (\rho^n u^n \otimes \tilde{u}^{n+1}) + \nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) = f^{n+1} \quad (\text{III.21})$$

3 - solve for  $p^{n+1}$ ,  $u^{n+1}$  and  $\rho^{n+1}$  :

$$\left\{ \begin{aligned} \rho^n \frac{u^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - p^n) &= 0 \\ \frac{\varrho(p^{n+1}, y^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho(p^{n+1}, y^{n+1}) u^{n+1}) &= 0 \\ \rho^{n+1} &= \varrho(p^{n+1}, y^{n+1}) \end{aligned} \right. \quad (\text{III.22})$$

After a computation of the unknown  $y$  (step 1), step 2 consists in a semi-implicit solution of the momentum equation to obtain a predicted velocity. Step 3 is a nonlinear pressure correction step, which degenerates in the usual projection step used in incompressible flow solvers when the density is constant (e.g. [53]). Taking the divergence of the first relation of (III.22) and using the second

one to eliminate the unknown velocity  $u^{n+1}$  yields a non-linear elliptic problem for the pressure. This computation is formal in the semi-discrete formulation, but, of course, is necessarily made clear at the algebraic level, as described in section III.3.3. Once the pressure is computed, the first relation yields the updated velocity and the third one gives the end-of-step density.

The main difficulty to design a fractional step algorithm for the discretization of the problem at hand lies in the approximation of the density. Indeed, we have to meet two requirements : first, to satisfy the compatibility condition (condition (III.7)) when computing  $y$ , second to ensure the conservativity of the scheme. The first point has been shown in the preceding section to be necessary to a reliable computation of the unknown  $y$ , and, from 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 from our experience, using a non-conservative scheme for the approximation of  $y$  leads to large errors. This is specially important when the equation of state is strongly non-linear, as for flows involving phases of very different densities. In this latter case, reasonable results could not be obtained with the non-conservative alternative described in remark III.3.11 below, even by drastically limitating the time-step.

To meet both requirements, we use here a time-shift of the density : in the advection terms of both the computation of  $y$  and the prediction of the velocity, 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 will be assessed in numerical experiments.

Finally, by a computation similar to the proof of lemma III.2.3, the condition (III.7) also guarantees an  $L^2$ -stability property for the advection operator ; the chosen time-discretization for the density thus also yields the  $L^2$ -stability of the advection of the velocity, *i.e.* the discrete counterpart of the following relation, which 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 :

$$\int_{\Omega} \left[ \frac{\partial \rho u}{\partial t} + \nabla \cdot (\rho u) \right] \cdot u = \frac{d}{dt} \int_{\Omega} \frac{1}{2} \rho |u|^2 \quad (\text{III.23})$$

This result, in the discrete case, is proven in [31] and recalled here in theorem III.3.12. Although we do not know whether this stability property is really useful in the applications addressed here, it has been found to be essential for convection dominant flows [3]. It is also one of the ingredients used in [31] to derive a pressure correction scheme for compressible barotropic flows which conserves the entropy of the system.

**Remark III.3.11 (On another time-discretization of the density)** *To satisfy condition (III.7), another way to proceed has already been proposed [2, 31]. 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 step reads :*

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

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  $y$  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 [2]. 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  $\rho^{-1}$  and  $u^{-1}$  are suitable approximations for the initial density and the velocity, respectively.

### III.3.2 Spatial discretization of the momentum balance equation

The spatial discretization of the momentum and mass balance equations relies on the so-called "rotated bilinear element" introduced by Rannacher and Turek for quadrilateral or hexahedric meshes [61], or on the Crouzeix-Raviart element for simplicial meshes [18]. For both elements, the pressure is piecewise constant over each cell, and the approximation space is the same for each component of the velocity. For the Rannacher-Turek element, in two, respectively three dimensions, it is spanned on the reference element by the set of functions  $\{1, x_1, x_2, x_1^2 - x_2^2\}$ , respectively  $\{1, x_1, x_2, x_3, x_1^2 - x_2^2, x_1^2 - x_3^2\}$ ; the mapping to a generic element is performed by the standard Q1 mapping. For the Crouzeix-Raviart element, the approximation space is the space of piecewise affine functions. Only the continuity of the integral over each edge (in two dimensions) or face (in three dimensions) of the mesh is imposed, so the velocities are discontinuous through each edge or face; the discretization is thus non-conforming in  $H^1(\Omega)$ . The linear forms defining the velocity degrees of freedom for each element are the integral over each edge or face. Each degree of freedom can then be univoquely associated to an edge or face, excluding the external ones if we suppose that the velocity obeys a homogeneous Dirichlet boundary condition, which is done here to simplify the presentation. We take benefit of this relationship to index the degrees of freedom for the velocity, the set of which for any discrete velocity field  $v$  thus reads  $\{v_\sigma, \sigma \in \mathcal{E}_{\text{int}}\}$ , each  $v_\sigma$  being a vector of  $\mathbb{R}^d$ , of components  $v_{\sigma,i}$ ,  $i = 1, d$ . The degrees of freedom for the pressure are denoted by  $\{p_K, K \in \mathcal{M}\}$ .

These pairs of approximation spaces for the velocity and the pressure are *inf-sup* stable. They do not satisfy the discrete Korn lemma, but this can be cured by adding a stabilization term built from the jumps of the velocity across each edge or face [8]; however, this problem is known to essentially affect the Crouzeix-Raviart element. It is not addressed in the numerical experiments presented here : indeed, each time the divergence of the stress tensor is not set under the coercive form (*i.e.*  $\nabla \cdot \tau(u) = \mu \Delta u + \mu/3 \nabla \nabla \cdot u$ ), the chosen element is the Rannacher-Turek one. If not specified, the original form of the stress tensor (III.2) is used. No stabilization strategy has been implemented. Finally let us remark that, even at the continuous level, in the compressible case (*i.e.* for non-divergence-free velocity fields) and when the viscosity is not constant, Korn's lemma may not hold and, worse, the viscous term may be not dissipative.

The main difficulty in the use of these elements is to obtain, in the discrete case, the  $L^2$  stability induced by the advection terms, that is, in other words, to build a discretization satisfying a discrete analogue of relation (III.23). To this purpose, we follow an idea developped in [3]; we review here its main arguments for the sake of completeness. The first step is to derive a finite-volume-like discretization of the convection operator, in order to make use of the following result [31].

**Theorem III.3.12 (Stability of finite-volume advection operators)**

Let  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(F_{\sigma,K})_{K \in \mathcal{M}, \sigma=K|L}$  be three families of real numbers such that the condition (III.7) is satisfied and,  $\forall \sigma = K|L \in \mathcal{E}_{\text{int}}, F_{\sigma,K} = -F_{\sigma,L}$ . Let  $(z_K^*)_{K \in \mathcal{M}}$  and  $(z_K)_{K \in \mathcal{M}}$  be two families of real numbers. For any internal edge  $\sigma = K|L$ , we define  $z_\sigma$  either by  $z_\sigma = \frac{1}{2}(z_K + z_L)$ , or by  $z_\sigma = z_K$  if  $F_{\sigma,K} \geq 0$  and  $z_\sigma = z_L$  otherwise. The first choice is referred to as the "centered choice", the second one as "the upwind choice". In both cases, the following stability property holds :

$$\sum_{K \in \mathcal{M}} z_K \left[ \frac{|K|}{\delta t} (\rho_K z_K - \rho_K^* z_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} z_\sigma \right] \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K z_K^2 - \rho_K^* z_K^{*2}]$$

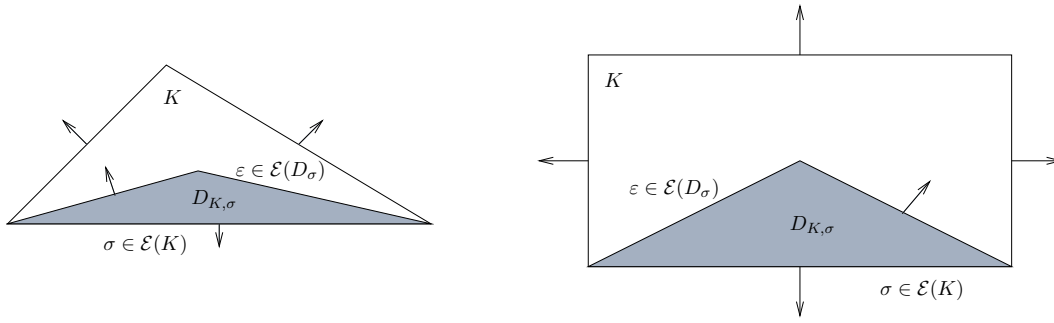


FIG. III.1 – Diamond-cells for the Crouzeix-Raviart and Rannacher-Turek element.

For each internal edge  $\sigma = K|L$ , let  $D_{K,\sigma}$  be the conic volume having  $\sigma$  for basis and the mass center of the mesh  $K$  as additional vertex. The volume  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$  is referred to as the "diamond cell" associated to  $\sigma$  and  $D_{K,\sigma}$  is the half-diamond cell associated to  $\sigma$  and  $K$  (see figure III.1). As the pressure, the density is approximated by piecewise constant functions over each element. In addition, for the Crouzeix-Raviart element and, for the Rannacher-Turek element, when the mesh is a rectangle (in two dimensions) or a cuboid (in three dimensions), the integral of the shape function associated to the edge  $\sigma$  over the element  $K$  is the measure of the half-diamond cell  $D_{K,\sigma}$ . Thus, the application of the mass lumping to the terms  $\rho^n u^{n+1}$  in the equations corresponding to the velocity on the edge  $\sigma$ , leads to an expression of the form  $\rho_\sigma^n u_\sigma^{n+1}$ , where  $\rho_\sigma^n$  results from an average of the values taken by the density in the two elements adjacent to  $\sigma$ , weighted by the measure of the half-diamonds :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad |D_\sigma| \rho_\sigma^n = |D_{K,\sigma}| \rho_K^n + |D_{L,\sigma}| \rho_L^n \quad (\text{III.24})$$

where  $|D_\sigma|$  is the measure of the diamond cell  $D_\sigma$ ,  $|D_{K,\sigma}|$  and  $|D_{L,\sigma}|$  are the measures of the half-diamond cells associated respectively to  $\sigma$  and  $K$  and to  $\sigma$  and  $L$ . This suggests the following discretization for the advection term :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad \text{advection term} \sim \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma}^n u_\varepsilon^{n+1} \quad (\text{III.25})$$

where  $\mathcal{E}(D_\sigma)$  is the set of the edges of  $D_\sigma$ ,  $u_\varepsilon^{n+1}$  is a centered approximation of  $u^{n+1}$  on  $\varepsilon$  and  $F_{\varepsilon,\sigma}^n$  is expressed as follows :

$$F_{\varepsilon,\sigma}^n = |\varepsilon| (\widetilde{\rho u})|_\varepsilon \cdot n_{\varepsilon,\sigma}$$

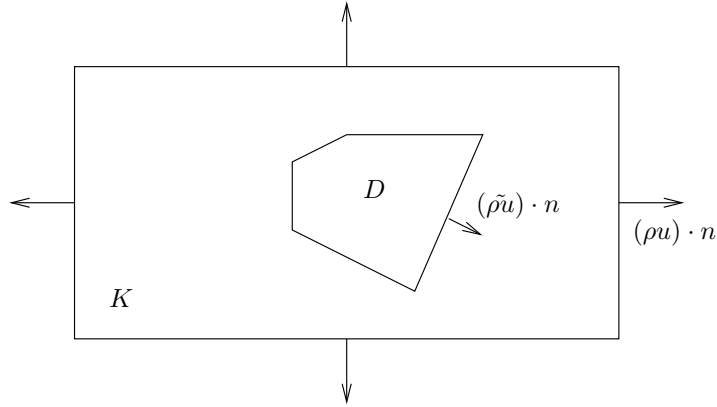


FIG. III.2 – sub-volume of  $K$ .

where  $|\varepsilon|$  is the measure of  $\varepsilon$  and  $n_{\varepsilon,\sigma}$  is the normal to  $\varepsilon$  outward  $D_\sigma$ .

The second step is now to build an approximation of the massic velocity  $(\widetilde{\rho u})|_\varepsilon$  which satisfies the compatibility condition of theorem III.3.12 (in fact, the discrete mass balance over the diamond cells). To this goal, we use the following result (see [3] for its (elementary) proof).

**Lemma III.3.13 (Mass-balance in a sub-volume of a mesh)**

Let  $K \in \mathcal{M}$  and assume that there exists two real numbers  $\rho^*$  and  $\rho$ , constant over  $K$  and a family of fluxes  $[|\sigma| (\rho u)_\sigma \cdot n_\sigma]_{\sigma \in \mathcal{E}(K)}$  such that :

$$\frac{|K|}{\delta t} (\rho - \rho^*) + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma \cdot n_\sigma = 0 \quad (\text{III.26})$$

Let  $D$  be a subset of  $K$ , of boundary  $\partial D$  (see figure III.2). Let  $w$  be a massic velocity field, supposed to be regular over  $K$ , such that the quantity  $\nabla \cdot w$  is constant over  $K$  and such that :

$$\int_K \nabla \cdot w = \int_{\partial K} w \cdot n_{\partial K} = \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma \cdot n_\sigma$$

where  $\partial K$  and  $n_{\partial K}$  stand for the boundary of  $K$  and the normal vector to  $\partial K$  outward to  $K$ , respectively. 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  $\partial D$  outward to  $D$ .

Since the pressure is constant over each mesh, the mass balance at the previous time step will take the following finite volume form :

$$\frac{|K|}{\delta t} (\rho^n - \rho^{n-1}) + \sum_{\sigma \in \mathcal{E}(K)} |\sigma| (\rho u)_\sigma^n \cdot n_\sigma = 0$$

where  $|\sigma| (\rho u)_\sigma^n \cdot n_\sigma$  is the mass flux through  $\sigma$ , which will be precisely defined in the next section. Using the previous lemma, to obtain the massic velocities, we must build on each cell  $K$  a field  $\widetilde{\rho u}$

with constant divergence and such that :

$$\forall \sigma \in \mathcal{E}(K), \quad \int_{\sigma} \widetilde{\rho u} \cdot n_{\sigma} = |\sigma| (\rho u)_{\sigma}^n \cdot n_{\sigma} \quad (\text{III.27})$$

The massic velocity through the edge or face  $\varepsilon$  of the diamond cell,  $(\widetilde{\rho u})|_{\varepsilon}$ , will then be computed by integrating  $\widetilde{\rho u}$  over  $\varepsilon$ , which will yield a discrete mass balance over both half-diamond cells ; it is then easy to check that summing these relations will in turn give the desired compatibility condition. Such a field  $\widetilde{\rho u}$  is derived for the Crouzeix-Raviart element by direct (*i.e.* using the standard expansion of the Crouzeix-Raviart elements) interpolation of the quantities  $((\rho u)_{\sigma}^n)_{\sigma \in \mathcal{E}(K)}$  :

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

where  $\varphi_{\sigma}$  is the Crouzeix-Raviart basis function associated to the velocity node of  $\sigma$ . For the Rannacher-Turek element, when the mesh is a rectangle or a cuboid, it is obtained by the following interpolation formula :

$$\widetilde{\rho u}(x) = \sum_{\sigma \in \mathcal{E}(K)} \alpha_{\sigma}(x \cdot n_{\sigma}) [(\rho u)_{\sigma} \cdot n_{\sigma}] n_{\sigma}$$

where the  $\alpha_{\sigma}(\cdot)$  are affine interpolation functions which are determined in such a way that the relations (III.27) hold. Extension for more general grids is underway.

Finally, standard finite elements techniques are used to discretize the terms  $\nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1})$ , to obtain the following discrete momentum balance equation :

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

$$\begin{aligned} \frac{|D_{\sigma}|}{\delta t} (\tilde{\rho}_{\sigma}^n \tilde{u}_{\sigma,i}^{n+1} - \rho_{\sigma}^{n-1} u_{\sigma,i}^n) + \sum_{\substack{\varepsilon \in \mathcal{E}(D_{\sigma}), \\ \varepsilon = D_{\sigma} | D_{\sigma'}}} \frac{1}{2} F_{\varepsilon, \sigma}^n (\tilde{u}_{\sigma,i}^{n+1} + \tilde{u}_{\sigma',i}^{n+1}) \\ + a_d(\tilde{u}^{n+1}, \varphi_{\sigma}^{(i)}) - \int_{\Omega, h} p^n \nabla \cdot \varphi_{\sigma}^{(i)} = \int_{\Omega} f^{n+1} \cdot \varphi_{\sigma}^{(i)} \end{aligned} \quad (\text{III.28})$$

where  $\varphi_{\sigma}^{(i)}$  is the vector shape function associated to the velocity degree of freedom related to  $\sigma$ , which reads  $\varphi_{\sigma} e^{(i)}$  with  $e^{(i)}$  the  $i^{\text{th}}$  vector of the canonical basis of  $\mathbb{R}^d$  and  $\varphi_{\sigma}$  the scalar shape function. The bilinear form  $a_d(\cdot, \cdot)$  represent the viscous term and, for all discrete velocity fields  $v$  and  $w$ ,  $a_d(v, w)$  is defined as follows :

$$a_d(v, w) = \begin{cases} \mu \int_{\Omega, h} \left[ \nabla v : \nabla w + \frac{1}{3} \nabla \cdot v \nabla \cdot w \right] & \text{if (III.3) holds (case of constant viscosity),} \\ \int_{\Omega, h} \tau(v) : \nabla w & \text{with } \tau \text{ given by (III.2) otherwise.} \end{cases}$$

### III.3.3 Spatial discretization of the projection step

The discretization of the first equation of the pressure correction step is consistent with the momentum balance one ; a mass lumping is performed for the unsteady term and a standard finite element formulation is used for the gradient of the pressure increment :

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

As the pressure is piecewise constant, the discrete gradient operator takes the form of the transposed of the finite volume standard discretization of the divergence (based on the finite element mesh, and not on the diamond cells) and can be rewritten as follows :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \quad \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (u_\sigma^{n+1} - \tilde{u}_\sigma^{n+1}) + |\sigma| [(p_L^{n+1} - p_L^n) - (p_K^{n+1} - p_K^n)] n_{KL} = 0 \quad (\text{III.29})$$

For the same reason, the approximation of the time derivative of the density in the mass balance will also look like a finite volume term. This point suggests a finite volume discretization of this latter equation, which reads :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} [\varrho(p_K^{n+1}, y_K^{n+1}) - \rho_K^n] + \sum_{\sigma=K|L} F_{\sigma,K}^{n+1} = 0 \quad (\text{III.30})$$

To ensure the positivity of the density, we use an upwinding technique for the convection term, then the mass flux across  $\sigma$ ,  $F_{\sigma,K}^{n+1}$ , is expressed as follows :

$$F_{\sigma,K}^{n+1} = |\sigma| (\rho^{n+1} u^{n+1})|_\sigma \cdot n_\sigma = (v_{\sigma,K}^{n+1})^+ \varrho(p_K^{n+1}, y_K^{n+1}) - (v_{\sigma,K}^{n+1})^- \varrho(p_L^{n+1}, y_L^{n+1})$$

where  $(v_{\sigma,K}^{n+1})^+$  and  $(v_{\sigma,K}^{n+1})^-$  stands respectively for  $\max(v_{\sigma,K}^{n+1}, 0)$  and  $-\min(v_{\sigma,K}^{n+1}, 0)$  with  $v_{\sigma,K}^{n+1} = |\sigma| u_\sigma^{n+1} \cdot n_{KL}$ .

Under some restrictive assumptions for the equation of state, we can prove that this projection step admits one solution.

**Lemma III.3.14**

*Let us suppose that the equation of state satisfies the following assumptions :*

- (1)  $\forall y \in [0, 1]$ , the function  $p \mapsto \varrho(p, y)$  is defined and increasing on  $[0, +\infty)$  and is such that  $\varrho(0, y) = 0$  and  $\lim_{p \rightarrow +\infty} \varrho(p, y) = +\infty$ .
- (2) The function  $p \mapsto \varrho(p, y)$  tends to infinity uniformly with respect to  $y$ , i.e. :

$$\exists M_0 > 0 \text{ such that } \forall y \in [0, 1], \forall M > M_0, \varrho(p, y) \leq M \Rightarrow p \leq c(M)$$

*where the positive real number  $c(M)$  is independent of  $y$ .*

*Then the non-linear algebraic system (III.29)-(III.30) admits at least a solution.*

**Proof.**

The proof of this lemma consists in an application of the Brouwer fixed point theorem (e.g. [22, chapter 5]), the idea of which can be found in [27]. Let  $H(\cdot)$  be the function, from  $V_{\mathcal{M}} \times X_{\mathcal{M}}$  to itself, where  $V_{\mathcal{M}}$  stands for the approximation space for the velocity, defined by  $(u, p) \mapsto (v, q)$  where  $v$  and  $q$  are obtained as follows :

- $q$  is solution of the mass balance (III.30) with fluxes evaluated with the velocity field  $u$ . The fact that this solution exists relies on the following argument : if we choose in (III.30) the density  $\rho$  as unknown, this equation becomes linear and admits a unique solution, by the fact that the density is everywhere positive because of the upwinding; then, by assumption (1) of the lemma,  $\forall K \in \mathcal{M}$ , knowing  $\rho_K$  and  $y_K$ , we compute  $q_K$  by the equation of state.
- $v$  is then solution of (III.29) taking  $q$  for the pressure.



Any fixed point of the function  $H(\cdot)$  is a solution to the system (III.29)-(III.30). By conservativity of the finite volume equation (III.30), we easily see that the integral of the density  $\varrho(p_K^{n+1}, y_K^{n+1})$  is bounded (in fact the same as the integral of  $\rho^n$ ), which, as the density is positive, yields an  $L^\infty$  estimate for the density, thus, by assumption (2), also for the pressure (say  $\max_{K \in \mathcal{M}} q_K \leq c_p$ ), and, finally, an  $L^2$  estimate for the pressure. Exploiting the variational form of (III.29), we then obtain that  $\|v\|_{\rho^n}$  is bounded, say  $\|v\|_{\rho^n} \leq c_u$ , where :

$$\|v\|_{\rho^n} = \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma^n |v|^2$$

By application of the Brouwer theorem, the application  $H(\cdot)$  thus admits a fixed point in the convex  $\mathcal{C}$  defined by :

$$\mathcal{C} = \{(v, q) \in X_{\mathcal{M}} \times V_{\mathcal{M}} \text{ such that } \|v\|_{\rho^n} \leq c_u \text{ and } \max_{K \in \mathcal{M}} q_K \leq c_p\}$$

■

**Remark III.3.15** *For the numerical tests presented in section III.4, we use the classical equation of state :*

$$\rho = \frac{1}{y/\rho_g + (1-y)/\rho_\ell}$$

where the gas density  $\rho_g$  depends on the pressure but the liquid density  $\rho_\ell$  is supposed to be constant. This equation of state does not satisfy the assumption of the preceding lemma ; indeed, we even have  $\rho \leq \rho_\ell/(1-y)$ . Since the total mass in the system (i.e. the integral of the density) is constant by conservativity, it is easy to imagine situations where system (III.29)-(III.30) has no solution : for instance, we can suppose that the total mass is equal to  $2|\Omega|\rho_\ell$  and that the gass mass fraction  $y$  is everywhere lower than  $1/2$ . This problem can easily be cured by assuming a slight compressibility of the liquid. However, in our tests, this was never necessary.

Let us now combine the two algebraic relations (III.29) and (III.30) to build a discrete elliptic problem for the pressure. To this purpose, let us introduce the algebraic formulation of this system :

$$\begin{cases} \frac{1}{\delta t} M_{\rho^n} (u^{n+1} - \tilde{u}^{n+1}) + B^t (p^{n+1} - p^n) = 0 \\ \frac{1}{\delta t} R(\varrho(p^{n+1}, y^{n+1}) - \rho^n) - B Q_{\rho_{\text{up}}^{n+1}} u^{n+1} = 0 \end{cases} \quad (\text{III.31})$$

In the first relation,  $M_{\rho^n}$  stands for the diagonal mass matrix weighted by the density at  $t^n$  (at edges or faces center)  $\rho_\sigma^n$  (i.e. the  $d$  diagonal entries of  $M_{\rho^n}$  associated to the edge  $\sigma$  are given by  $|D_\sigma| \rho_\sigma^n$ ) ;  $B^t$  is the matrix of  $\mathbb{R}^{(dN) \times M}$ , where  $N$  is the number of internal edges (i.e.  $N = \text{card}(\mathcal{E}_{\text{int}})$ ) and  $M$  is the number of control volumes in the mesh (i.e.  $M = \text{card}(\mathcal{M})$ ), associated to the gradient operator ; consequently, the matrix  $B$  is associated to the opposite of the divergence operator. In the second relation,  $Q_{\rho_{\text{up}}^{n+1}}$  is a diagonal matrix, the entry of which corresponding to an edge  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , is obtained by simply taking the density at  $t^{n+1}$  in the element located upstream of  $\sigma$  with respect to  $u^{n+1}$ , i.e. either  $\varrho(p_K^{n+1}, y_K^{n+1})$  or  $\varrho(p_L^{n+1}, y_L^{n+1})$ . The matrix  $R$  is diagonal and, for any  $K \in \mathcal{M}$ , its entry  $R_{K,K}$  is the measure of the element  $K$ . The elliptic problem for the pressure is obtained by multiplying the first relation of (III.31) by  $B Q_{\rho_{\text{up}}^{n+1}} (M_{\tilde{\rho}^{n+1}})^{-1}$  and using the second one. This equation reads :

$$L p^{n+1} + \frac{1}{\delta t^2} R \varrho(p^{n+1}, y^{n+1}) = L p^n + \frac{1}{\delta t^2} R \rho^n + \frac{1}{\delta t} B Q_{\rho_{\text{up}}^{n+1}} \tilde{u}^{n+1} \quad (\text{III.32})$$

where  $L = B Q_{\rho^{n+1}}^{\text{up}} (M_{\rho^n})^{-1} B^t$  can be evaluated, as seen in [31], in the "finite volume way" as an approximation of the Laplace operator by the following relation, valid for each element  $K$  :

$$(L p^{n+1})_K = \sum_{\sigma=K|L} \frac{\rho_{\text{up},\sigma}}{\rho_\sigma^n} \frac{|\sigma|^2}{|D_\sigma|} (p_K - p_L)$$

where  $\rho_{\text{up},\sigma}$  stands for the upwind density associated to the edge  $\sigma$ . Note that, even if  $\forall \sigma \in \mathcal{E}_{\text{int}}$ ,  $\rho_{\text{up},\sigma} = \rho_\sigma^n$ , the operator  $L$  differs from the usual finite volume Laplace operator (in the case of cuboids, by a factor  $d$  [31]), which is linked to the fact that neither the Rannacher-Turek element nor the Crouzeix-Raviart one provides a consistent approximation of the Darcy problem. Provided that  $p^{n+1}$  is known, the first relation of (III.31) gives us the updated value of the velocity :

$$u^{n+1} = \tilde{u}^{n+1} - \delta t (M_{\rho^n})^{-1} B^t (p^{n+1} - p^n) \quad (\text{III.33})$$

As, to preserve the positivity of the density, we want to use in the mass balance the value of the density upwinded with respect to  $u^{n+1}$ , equations (III.32) and (III.33) are not decoupled, by contrast with what happens in usual projection methods. We thus implement the following iterative algorithm :

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

Step 4.1 – Solve for  $p_{k+1}^{n+1}$  :

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

where the density in  $L$  and  $Q_{\rho^{n+1}}^{\text{up}}$  is evaluated at  $p_k^{n+1}$  and  $y^{n+1}$  and the upwinding in  $Q_{\rho^{n+1}}^{\text{up}}$  is performed with respect to  $u_k^{n+1}$

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

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

### III.3.4 Spatial discretization for the nonlinear advection diffusion equation

The discretization of this step is performed by the scheme detailed in section III.2, which we recall here :

$\forall K \in \mathcal{M}$ ,

$$\begin{aligned} |K| \frac{\rho_K^n y_K^{n+1} - \rho_K^{n-1} y_K^n}{\delta t} + \sum_{\sigma=K|L} [(F_{\sigma,K}^n)^+ y_K^{n+1} - (F_{\sigma,K}^n)^- y_L^{n+1}] \\ + \sum_{\sigma=K|L} [(G_{\sigma,K}^n)^+ g(y_K^{n+1}, y_L^{n+1}) - (G_{\sigma,K}^n)^- g(y_L^{n+1}, y_K^{n+1})] \\ + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_\sigma} (y_K^{n+1} - y_L^{n+1}) = 0 \end{aligned} \quad (\text{III.34})$$

where  $F_{\sigma,K}^n$  is the mass flux used for the mass balance equation at the previous time step (*i.e.* at time  $t^n$ ) and the density used to evaluate the quantity  $G_{\sigma,K}^n$  is the same density as for  $F_{\sigma,K}^n$  (so results from an upwind choice with respect to the velocity at time  $t^n$ ). The function  $\varphi(\cdot)$  is supposed to be given here by  $\varphi(x) = x(1-x)$ , and we choose as numerical monotone flux function  $g(a,b) = a - b^2$ .

### III.3.5 An overview of the algorithm

To sum up, the algorithm considered in this section is the following one :

1. Computation of  $y$  – Compute  $(y_K^{n+1})_{K \in \mathcal{M}}$  from equation (III.34), obtained by an upwind finite volume discretization. The discretization of transport terms are built according to the structure which is necessary to apply theorem III.2.2. This yields the existence and the uniqueness of the discrete solution together with the fact that it lies in the interval  $[0, 1]$ .
2. Prediction of the velocity – Compute  $(\tilde{u}_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  by equation (III.28), obtained by a finite volume discretization of the transport terms over the diamond cells and a finite element discretization of the other terms, the density on the edges at  $t^n$ ,  $(\rho_\sigma^n)_{\sigma \in \mathcal{E}_{\text{int}}}$ , being given by (III.24). The mass fluxes at the edges or faces of the diamond cells are built in such a way that the mass balance over the diamond cells derives from the mass balance over the primal cells. Then we have the structure which is necessary to apply theorem III.3.12. This yields an  $L^2$  stability property for the advection operator of the momentum balance.
3. Projection step – Compute  $(u_\sigma^{n+1})_{\sigma \in \mathcal{E}_{\text{int}}}$  and  $(p_K^{n+1})_{K \in \mathcal{M}}$  from equations (III.29) and (III.30), obtained by a finite element discretization of the velocity correction equation and an upwind finite volume discretization of the mass balance (over the elements  $K \in \mathcal{M}$ ). This upwind discretization ensures the positivity of the density.

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

#### **Theorem III.3.16 (Properties of the numerical scheme)**

Let  $(u^n)_{1 \leq n \leq N}$ ,  $(p^n)_{1 \leq n \leq N}$ ,  $(\rho^n)_{1 \leq n \leq N}$  and  $(y^n)_{1 \leq n \leq N}$  be the solution to the considered scheme, with a positive initial condition for  $\rho$  and an initial condition for  $y$  taking its values in  $[0, 1]$ . Then, the scheme enjoys the following properties, for all  $n \leq N$  :

(i) Let us assume that the assumptions for the equation of state of lemma III.3.14 hold and that the viscous term is dissipative (*i.e.* for any discrete velocity field  $v$ ,  $a_d(v, v) \geq 0$ ); then each step of the scheme has a solution.

(ii) positivity of the density :

$$\rho_K^n \geq 0, \quad \forall K \in \mathcal{M}$$

(iii)  $L^\infty$  stability property :

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

(iv) conservativity property :

$$\sum_{K \in \mathcal{M}} |K| \rho_K^{n-1} y_K^n = cste, \quad \sum_{K \in \mathcal{M}} |K| \rho^n = cste$$

#### **Proof.**

The existence of a solution to the balance equation for  $y$  is stated in theorem III.2.2. For the velocity prediction step, under the assumption that the viscous term is dissipative, it is a consequence of

the stability of the advection operator. The existence of a solution to the projection step is given by lemma III.3.14.

The second property follows from the upwind discretization of the mass balance in the projection step (III.22), the third property from theorem III.2.2, and the last one from the conservativity of the discretization. ■

### III.4 Numerical results

In this section, we present three numerical tests performed to assess the behaviour of the above described fractional step scheme. In all these tests, we compute the flow of an isothermal two-phase mixture of non-miscible liquid and gas. In this case, the model (III.1) is the so-called drift-flux model, that is a mixture model that takes into account the relative velocity  $u_r$  between the liquid and the gas phase (the so-called drift velocity), for which a phenomenologic relation must be supplied. Then,  $\rho$  stands for the mixture density and takes the general form :

$$\rho = (1 - \alpha_g)\rho_\ell + \alpha_g \varrho_g(p) \quad (\text{III.35})$$

where  $\alpha_g$  stands for the void fraction and  $\varrho_g(p)$  yields the gas density as a function of the pressure ; in the perfect gas approximation and for a constant temperature,  $\varrho_g(p)$  is simply proportional to the pressure :

$$\varrho_g(p) = \frac{p}{RT} \quad (\text{III.36})$$

where  $R$  is the gas constant and  $T$  is the absolute temperature. We assume further that the liquid phase can be considered as incompressible, then the liquid density  $\rho_\ell$  is constant. Introducing the mass gas fraction  $y$  in (III.35) by using the relation  $\alpha_g \varrho_g = \rho y$  leads to the following equation of state :

$$\varrho(p, y) = \frac{\rho_g \rho_\ell}{\rho_\ell y + (1 - y) \rho_g} \quad (\text{III.37})$$

As in the preceding section, the nonlinear function  $\varphi(\cdot)$  is given by  $\varphi(y) = y(1 - y)$  and the flux function by  $g(a, b) = a - b^2$ .

#### III.4.1 Assessing the convergence against an analytic solution

We first assess the convergence rate of the proposed scheme with respect to space and time discretizations, by dealing with a case where an analytic solution can be exhibited.

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

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

The pressure and the gas mass fraction are linked to the density by the equation of state (III.37), where the liquid density  $\rho_\ell$  is set at  $\rho_\ell = 5$  and the product  $RT$  in the equation of state of the gas (III.36) 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 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 ; for the gas mass fraction and momentum balance, we add the corresponding right-hand side. In this latter equation, we suppose that the divergence of the stress tensor is given by :

$$\nabla \cdot \tau(u) = \mu \Delta u + \frac{\mu}{3} \nabla \nabla \cdot u, \quad \mu = 10^{-2}$$

and we use the corresponding expression for the bilinear form  $a_d(\cdot, \cdot)$ .

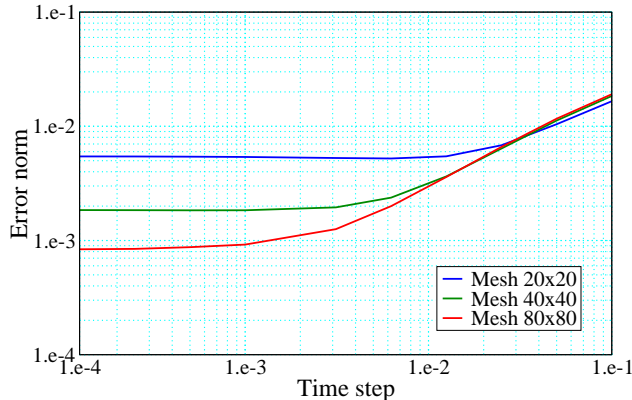


FIG. III.3 – Rannacher-Turek element - Velocity error as a function of the time step.

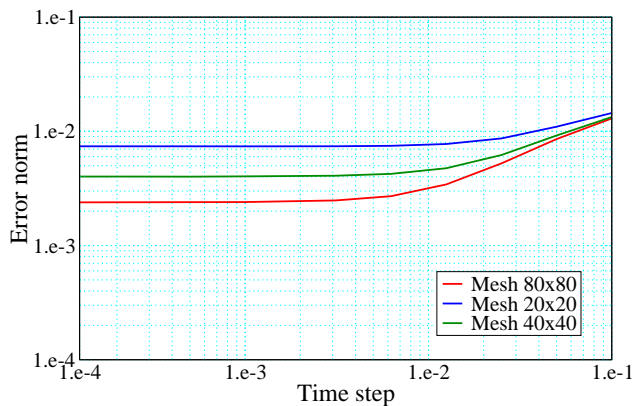


FIG. III.4 – Crouzeix-Raviart element – Velocity error as a function of the time step.

Velocity, pressure and gas mass fraction errors obtained at  $t = 0.5$  as a function of the time step are drawn, for the Rannacher-Turek element, on figure III.3, figure III.5 and III.7, respectively and, for the Crouzeix-Raviart element, on figure III.4, figure III.6 and III.8, 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 the Rannacher-Turek element, computations are made with  $20 \times 20$ ,  $40 \times 40$  and  $80 \times 80$  uniform meshes. For the Crouzeix-Raviart one, 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  $80^\circ$ , according to the pattern given in [6, figure 5 – bbbb]. The first splitting of the domain yields  $20 \times 20$ ,  $40 \times 40$  and  $80 \times 80$  uniform grids. 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

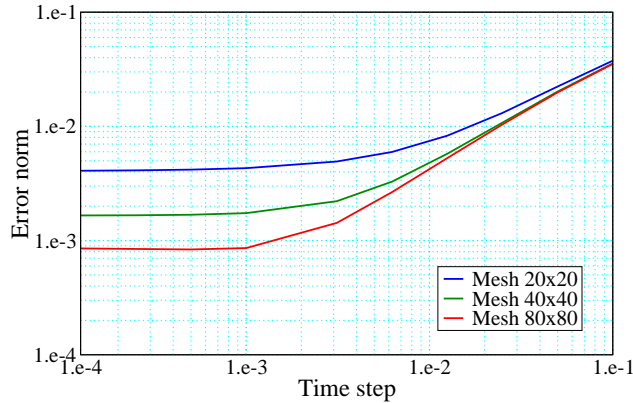


FIG. III.5 – Rannacher-Turek element – Pressure error as a function of the time step.

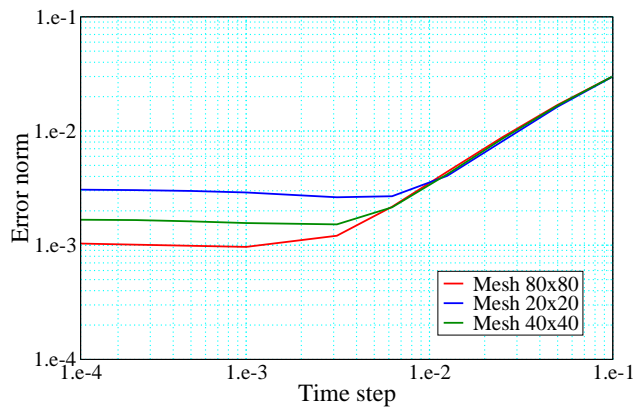


FIG. III.6 – Crouzeix-Raviart element – Pressure error as a function of the time step.

spatial discretization error. The value of the errors on this plateau then show 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, results seem to be significantly more accurate with the Rannacher-Turek element.

### III.4.2 A phase separation problem

We now present numerical results obtained for a phase separation problem, with data inspired from a classical benchmark test for the simulation of two-phase flows [17, 33, 66] with two-fields models (*i.e.* models considering separate balance equations for each phase). The considered physical domain is a vertical tube of length  $L = 7.5 m$ , 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 both 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 non-physical (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,

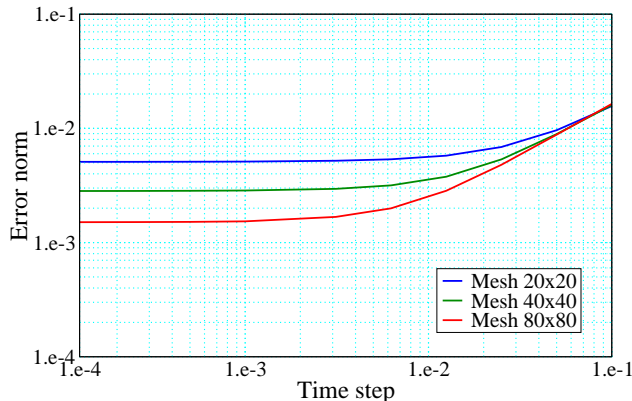


FIG. III.7 – Rannacher-Turek element – Gas mass fraction error as a function of the time step.

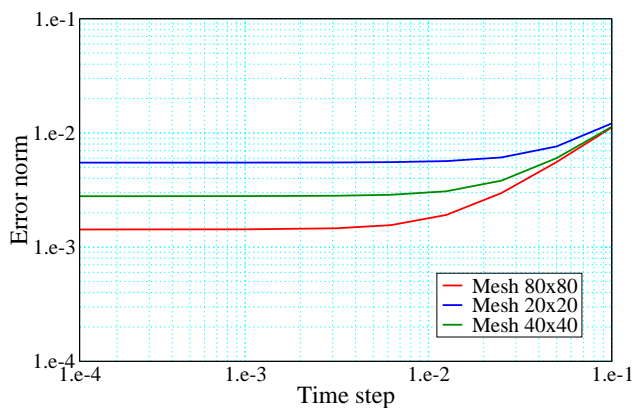


FIG. III.8 – Crouzeix-Raviart element – Gas mass fraction error as a function of the time step.

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  $\mu$  are set to zero. At the top and the bottom boundaries, both the velocity and the relative velocity are prescribed to zero.

For this test case, we use a regular meshing 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 one. We perform a special numerical integration of the forcing term of the momentum balance (here the gravity), which is designed to ensure that the discretization of a gradient is indeed a discrete gradient (*i.e.* if there exists a function  $\psi$  such that the forcing term  $f$  can be recast under the form  $f = \nabla\psi$ , the discrete right-hand side of the momentum balance belongs to the range of the discrete gradient  $\mathbf{B}^\dagger$ ). Note that, with the finite elements used here, thanks to the special discrete convection operator proposed in this paper and with a zero viscosity, the degrees of freedom for the velocity which are tangent to the edges are fully decoupled from the remainder of the computation; in particular, this applies here for the vertical velocities on the vertical boundaries, which are just set to zero, as the horizontal ones.

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

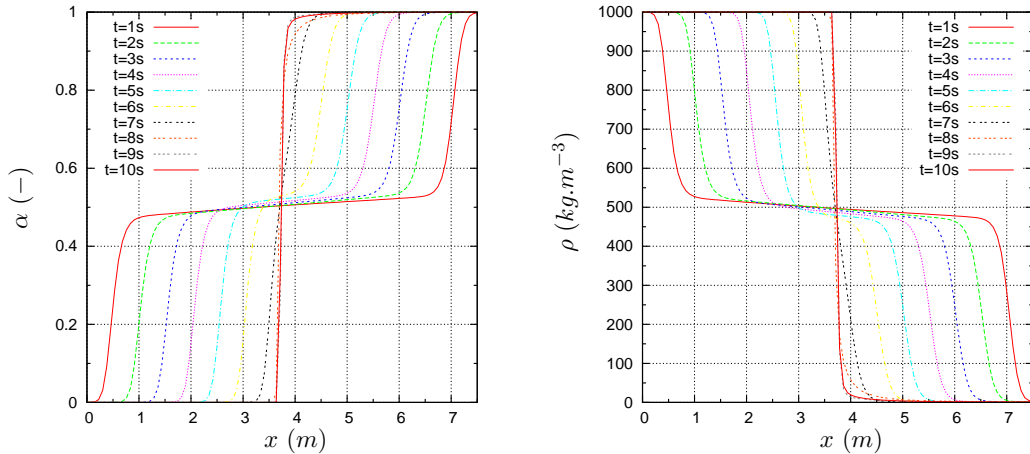


FIG. III.9 – void fraction and density profiles for the phase separation problem.

$\delta t = 10^{-2} s$  are reported on figure III.9. However, probably because of the very steep 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 of the obtained profiles are observed when decreasing the time step down to  $\delta t = 5.10^{-4} s$ .

### III.4.3 Flow of a sedimenting dilute suspension over a rectangular bump

In this section, we study the sedimentation of a dilute suspension in a two-dimensional channel with a rectangular bump (see figure III.10). The channel is  $D = 2 m$  high, the step height is  $H = 1 m$ , the length of the channel is  $W_I + W_S + W_O = 25 m$ , and the bump starts at  $W_I = 5 m$  and finishes at  $W_I + W_S = 7 m$ ; the remaining length after the bump is thus  $W_O = 18 m$ . Inlet conditions are  $u_\infty = 10 m.s^{-1}$  and  $y_\infty = 0.9$ .

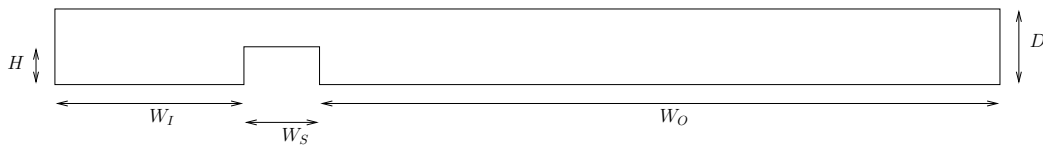


FIG. III.10 – geometry of the computational domain for the channel problem

The velocity is prescribed to zero at both the upper and the lower wall, and the pressure at the outflow section is fixed to the ambient pressure  $p = 10^5 Pa$ . The initial conditions for the velocity and the gas mass fraction are set to the inlet conditions, while the initial pressure is the usual ambient pressure ( $p = 10^5 Pa$ ) and the initial densities are  $\rho_\ell = 1000 kg.m^{-3}$  and  $\rho_g = p/RT = 1.2 kg.m^{-3}$ , for water and air respectively. The viscosity  $\mu$  is supposed to obey Batchelor and Green's law [4] :

$$\mu = \mu_0 \left[ 1 + \frac{5}{2}(1 - \alpha_g) + 7.6(1 - \alpha_g)^2 \right]$$

The coefficient  $\mu_0$  is chosen such that the Reynolds number  $Re$  is equal to 75, where  $Re$  is defined



as follows :

$$Re = \frac{H u_\infty \rho_\infty}{\mu_0}$$

The equation of state for the mixture is the same than in the previous tests and the diffusion coefficient  $D$  is zero.

A characteristic Mach number for this flow is defined by :

$$M_\infty = u_\infty \frac{\alpha_{g,\infty}}{\sqrt{y_\infty} \sqrt{RT}}$$

where, from the above expression for  $\rho_g$ , the gas constant  $R$  and the absolute temperature  $T$  are such that  $RT$  is approximately equal to  $9.10^4 \text{ J.mol}^{-1}$  and the inlet void fraction  $\alpha_{g,\infty}$  is determined from the inlet gas mass fraction and the phase densities by the relation  $\alpha_{g,\infty} = y_\infty \rho_\infty / \rho_g$ . The value of the Mach number is thus equal to  $M_\infty = 0.03$ , which corresponds to a very weakly compressible flow.

The computational domain is meshed with about 110 000 rectangular cells. In the horizontal direction, the space step is smaller near the bump and equal to  $\delta x_1 = 0.01 \text{ m}$ , and increases when moving away from the bump, up to  $\delta x_1 = 0.03 \text{ m}$  at the inlet and outlet sections. In the vertical direction, the mesh is uniform and  $\delta x_2 = 1/80 \text{ m}$ .

Figure III.11 displays the obtained isolines of the streamfunction, *i.e.*, as in divergence free flows, the function  $\psi$  obeying the partial differential equation  $-\Delta\psi = \partial u_1 / \partial x_2 - \partial u_2 / \partial x_1$ . Note that, as the flow is slightly compressible, the isolines of  $\psi$  only approximately match the trajectories of the particles in the flow.

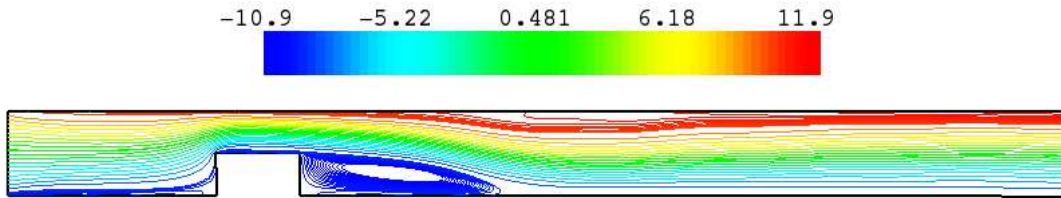


FIG. III.11 – Isolines of the streamfunction.

The gas mass fraction is represented on figure III.12. Due to the vertical fall of the liquid, almost only gas is left close to the top walls. On the opposite, an accumulation of liquid is observed in the stagnation zone of the flow in the wake of the bump and near the bottom wall. This liquid boundary layer is thinner in the vicinity of the reattachment point. The gas mass fraction in the flow remains in the interval  $[0.127, 0.9999]$ .

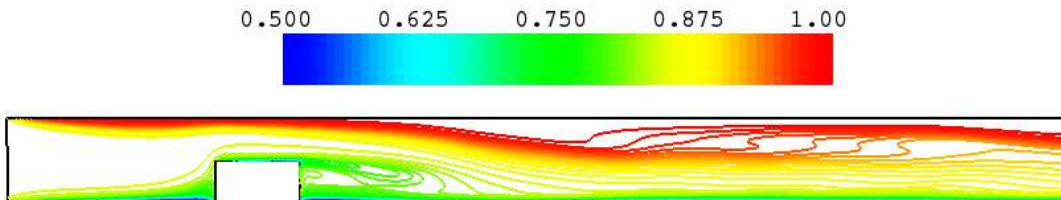


FIG. III.12 – Gas mass fraction in the flow.

### III.5 Conclusion

In this paper, we first address a parabolic equation used to model the phases mass balance in two-phase flows. This equation differs from the mass balance for chemical species in compressible multi-component flows studied in [51] by the addition of a non-linear term of the form  $\nabla \cdot \rho \varphi(y) u_r$ , where  $y$  is the unknown,  $\varphi(\cdot)$  is a regular function such that  $\varphi(0) = \varphi(1) = 0$  and  $u_r$  is a general (in particular non-necessarily divergence free) velocity field. We propose a finite-volume scheme for its numerical approximation, and we prove the existence and uniqueness of the solution, together with the fact that it remains within physical bounds, *i.e.* within the interval  $[0, 1]$ . As in [51], the necessary condition for this  $L^\infty$  stability is that the discretization of the convection operator must be such that it vanishes for constant  $y$ , which can be seen as a particular discrete mass balance equation. The second ingredient of this scheme is a discretization of the non-linear term based on the notion of monotone flux functions [26]. This work extends the theory developed in [51] in two directions : it copes with a new non-linear term, and introduces different techniques well-suited for non-linear problem, first to prove  $L^\infty$  *a priori* estimates for the solution (recasting the equation under variational form and choosing the possible negative part of the solution as test function), to prove its existence (by a topological degree argument) and its uniqueness (introducing an auxiliary "linear dual problem").

Second, we propose a discretization by a fractional step method for the set of equations composed of Navier-Stokes equations (*i.e.* mass and momentum balance) and the phases mass balance. This formulation decouples the solution of this latter (as previously mentioned, performed by a finite volume method), and the solution of the Navier-Stokes equations, performed by a pressure correction method based on low degree non-conforming finite elements. This algorithm meets two essential requirements : its is conservative, and the above-mentioned stability condition for the 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 [2] developed at IRSN and devoted to the modelling of reacting flows ; in this context, it demonstrates very satisfactory stability properties, even for cases where very steep variations of the density appear, and for which numerical difficulties are often reported in the literature. Finally, let us also mention that the proposed numerical scheme degenerates to classical projection methods 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 fractional step approach (all the equations are decoupled) and, consequently, the less time-consuming one, and this choice should probably 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 of which seems to need a drastic reduction of the time step. 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 phases mass balance but coupling this latter equation to the projection step, is now under development and shows promising results.

## Chapitre IV

# An entropy preserving finite-element/finite-volume pressure correction scheme for the drift-flux model

**Abstract.** We present in this paper a pressure correction scheme for the drift-flux model combining finite element and finite volume discretizations, which is shown to enjoy essential stability features of the continuous problem : the scheme is conservative, the unknowns are kept within their physical bounds and, in the homogeneous case (*i.e.* when the drift velocity vanishes), the discrete entropy of the system decreases ; in addition, when using for the drift velocity a closure law which takes the form of a Darcy-like relation, the drift term becomes dissipative. Finally, the present algorithm preserves a constant pressure and a constant velocity through moving interfaces between phases. To ensure the stability as well as to obtain this latter property, a key ingredient is to couple the mass balance and the transport equation for the dispersed phase in an original pressure correction step. The existence of a solution to each step of the algorithm is proven ; in particular, the existence of a solution to the pressure correction step is derived as a consequence of a more general existence result for discrete problems associated to the drift-flux model. Numerical tests show a near-first-order convergence rate for the scheme, both in time and space, and confirm its stability.

### IV.1 Introduction

Dispersed two-phase flows and, in particular, bubbly flows are widely encountered in industrial applications as, for instance, nuclear safety studies, which are the context of the present work. Within the rather large panel of models dealing with such flows, the simplest is the so-called drift-flux model, which consists in balance equations for an equivalent continuum representing both the gaseous and the liquid phase. For isothermal flows, this approach leads to a system of three balance

equations, namely the overall mass, the gas mass and the momentum balance, which reads :

$$\left\{ \begin{array}{l} \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \\ \frac{\partial \rho y}{\partial t} + \nabla \cdot (\rho y u) = -\nabla \cdot (\rho y (1 - y) u_r) + \nabla \cdot (D \nabla y) \\ \frac{\partial \rho u}{\partial t} + \nabla \cdot (\rho u \otimes u) + \nabla p - \nabla \cdot \tau(u) = f_v \end{array} \right. \quad (\text{IV.1})$$

where  $t$  stands for the time,  $\rho$ ,  $u$  and  $p$  are the (average) density, velocity and pressure in the flow and  $y$  stands for the gas mass fraction. The diffusion coefficient  $D$  represents in most applications small scale perturbations of the flow due to the presence of the dispersed phase, sometimes called "diphasic turbulence" and  $u_r$  is the relative velocity between the liquid and the gaseous phase (the so-called drift velocity); for both these quantities, a phenomenologic relation must be supplied. The forcing term  $f_v$  may represent, for instance, the gravity forces. The tensor  $\tau$  is the viscous part of the stress tensor, given by the following expression :

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

For a constant viscosity, this relation yields :

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

and, in this case, this term is dissipative (*i.e.* for any regular velocity field  $u$  vanishing on the boundary, the integral of  $\nabla \cdot \tau(u) \cdot u$  over the computational domain is non-negative).

This system must be complemented by an equation of state, which takes the general form :

$$\rho = \varrho^{p,\alpha}(p, \alpha_g) = (1 - \alpha_g) \rho_\ell + \alpha_g \varrho_g(p) \quad (\text{IV.4})$$

where  $\alpha_g$  stands for the void fraction and  $\varrho_g(p)$  expresses the gas density as a function of the pressure; in the ideal gas approximation and for an isothermal flow,  $\varrho_g(\cdot)$  is simply a linear function :

$$\varrho_g(p) = \frac{p}{a^2} \quad (\text{IV.5})$$

where  $a$  is a constant characteristic of the gas, equal to the sound velocity in an isothermal (monophasic) flow. The density of the liquid phase  $\rho_\ell$  is assumed to be constant. Introducing the mass gas fraction  $y$  in (IV.4) by using the relation  $\alpha_g \varrho_g = \rho y$  leads to the following equation of state :

$$\rho = \varrho^{p,y}(p, y) = \frac{\varrho_g(p) \rho_\ell}{\rho_\ell y + (1 - y) \varrho_g(p)} \quad (\text{IV.6})$$

The problem is supposed to be posed over  $\Omega$ , an open bounded connected subset 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  $\rho$ ,  $u$  and  $y$ .

To design a numerical scheme for the solution of the system (IV.1), one is faced with several difficulties. First, since the fluid density  $\rho_\ell$  is supposed not to depend on the pressure, almost incompressible zones, *i.e.* zones where the void fraction is low, may coexist in the flow with compressible zones, *i.e.* zones where the void fraction remains significant. This feature makes the problem particularly difficult to solve from a numerical point of view, because the employed numerical scheme

will have to cope with a wide range of Mach numbers, starting from zero to, let us say, for low to moderate speed flows, a fraction of unity. Second, the gas mass fraction  $y$  can be expected, both for physical and mathematical reasons, to remain in the  $[0, 1]$  interval, and it appears strongly desirable that the numerical scheme reproduces this behaviour at the discrete level. Finally, it appears from numerical experiments that, in order to avoid numerical instabilities, the algorithm should preserve a constant pressure through moving interfaces between phases (*i.e.* contact discontinuities of the underlying hyperbolic system). To obtain a scheme stable in the low Mach number limit, the solution that we adopt here is to use an algorithm inspired from the incompressible flow numerics, namely from the class of finite element pressure correction methods, and which degenerates to a classical projection scheme when the fluid density is constant. The last two requirements are met thanks to an original pressure correction step in which the mass balance equation is solved simultaneously with a part of the gas mass balance. For technical reasons, the solution of this latter equation is itself split in two steps, the first step thus being incorporated to the pressure correction step and the second one being performed independently.

This work takes benefit of ideas developed in a wide literature, so we are only able to quote here some references, the choice of which will unfortunately probably appear somewhat arbitrary. For a description of projection schemes for incompressible flow, see *e.g.* [38, 53] and references herein. An extension to barotropic Navier-Stokes equations close to the scheme developed here can be found in [31], together with references to (a large number of) related works (see *e.g.* [42] for the seminal work and [75] for a comprehensive introduction). Extensions of pressure correction algorithms for multi-phase flows are scarcer, and seem to be restricted to iterative algorithms, often similar in spirit to the usual SIMPLE algorithm for incompressible flows [65, 55, 50]. The gas mass balance equation, *i.e.* the second equation of (IV.1), is a convection-diffusion equation which differs from the usual mass balance for chemical species in compressible multi-component flows studied by Larrouturou [51] by the addition of a non-linear term of the form  $\nabla \cdot \rho \varphi(y) u_r$ , where  $\varphi(\cdot)$  is a regular function such that  $\varphi(0) = \varphi(1) = 0$  (in the present case,  $\varphi(y) = y(1 - y)$ ). In [35], we propose a finite-volume scheme for the numerical approximation of this type of equation, and we prove the existence and uniqueness of the solution, together with the fact that it remains within physical bounds, *i.e.* within the interval  $[0, 1]$ . Here, the proof of the same results combines arguments from both [51] and [35].

Several theoretical issues concerning the proposed scheme are studied in this paper. First, the existence of a solution to the pressure correction step, which consists in an algebraic non-linear system, is obtained by a topological degree argument. Second, we address the stability of the scheme. At the continuous level, the existence of an entropy for the system when the drift velocity vanishes (*i.e.* the homogeneous model) is well-known. In addition, it is shown in [39], by a Chapman-Enskog expansion technique, that the two-fluid model can be reduced to the drift-flux model when a strong coupling of both phases is assumed, with a Darcy-like closure relation for the drift velocity, *i.e.* an expression of the form :

$$u_r = \frac{1}{\lambda} (1 - \alpha_g) \alpha_g \frac{\varrho_g(p) - \rho_\ell}{\rho} \nabla p \quad (\text{IV.7})$$

where  $\lambda$  is a positive phenomenological coefficient. The same relation can also be obtained by neglecting in the two-fluid model the difference of acceleration between both phases [63]. With such an expression for  $u_r$ , the drift term becomes a second order term, and it is shown in [39] that it is consistent with the entropy of the homogeneous model (*i.e.* that it generates a non-negative dissipation of the entropy). These results are proven here at the discrete level : up to a minor modification of the proposed scheme, which seems useless in practice, the entropy is conserved when  $u_r$  is equal to zero, and when the closure relation (IV.7) applies and with a specific discretization, the drift term generates a dissipation.

This paper is built as follows. The fractional step algorithm for the solution of the whole problem is first presented in section IV.2, together with some of its properties : the existence of a solution to each step of the algorithm, the fact that the unknowns are kept within their physical bounds and that the algorithm is able to preserve a constant pressure and a constant velocity through moving interfaces between phases. The proof of the existence of the solution to the pressure correction step is obtained as a consequence of a more general existence theory for some discrete problems associated to the drift-flux model, which is exposed in the appendix. Next two sections are devoted to the stability analysis of the scheme; after establishing estimates for the work of the pressure forces (section IV.3), we first address the case  $u_r = 0$  (section IV.4.1), then the case where  $u_r$  is given by the Darcy-like closure relation (IV.7) (section IV.4.2). Finally, numerical tests are reported in section IV.5; they include a problem exhibiting an analytical solution which allows to assess convergence properties of the discretization, a sloshing transient in a cavity, and the evolution of a bubble column.

For the sake of simplicity, we suppose for the presentation of the scheme and its analysis (sections IV.2, IV.3 and IV.4) that the velocity is prescribed to zero on the whole boundary  $\partial\Omega$  of the computational domain, and that the gas mass flux through  $\partial\Omega$  also vanishes, so that both  $u_r = 0$  and the normal component of  $\nabla y$  is zero on the boundary. Moreover, the analysis of the scheme assumes that pure liquid zones do not exist in the flow; with the proposed algorithm, this is a consequence of the fact that such zones are not present at the initial time (*i.e.*, at  $t = 0$ ,  $y \in (0, 1]$ ). Indeed, getting rid of this latter limitation at the theoretical level seems to be a difficult task. However, the numerical tests presented in section IV.5 are not restricted to these situations. In particular,  $y = 0$  in the liquid column in the sloshing problem, up to spurious phases mixing by the numerical diffusion near the free surface; it is also the case at the initial time in the bubble column simulation.

In the presentation of the scheme, the drift velocity is supposed to be known, *i.e.* to be given by a closure relation independent of the unknowns of the problem, and this still holds in numerical experiments. The case where  $u_r$  is given by (IV.7) is thus only treated from a theoretical point of view in section IV.4.2.

## IV.2 The numerical algorithm

We present in this section the numerical scheme considered in this paper. We begin by describing the proposed scheme in the time semi-discrete setting, then we introduce the spatial discretization spaces and we detail the discrete approximation and the properties for each step of the algorithm at hand.

### IV.2.1 Time semi-discrete formulation

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

1 - solve for  $\tilde{u}^{n+1}$

$$\frac{\rho^n \tilde{u}^{n+1} - \rho^{n-1} u^n}{\delta t} + \nabla \cdot (\rho^n u^n \otimes \tilde{u}^{n+1}) + \nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) = f_v^{n+1} \quad (\text{IV.8})$$

2 - solve for  $p^{n+1}$ ,  $u^{n+1}$ ,  $\rho^{n+1}$  and  $z^{n+1}$

$$\left\{ \begin{array}{l} \rho^n \frac{u^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - p^n) = 0 \\ \frac{\varrho^{p,z}(p^{n+1}, z^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho^{p,z}(p^{n+1}, z^{n+1}) u^{n+1}) = 0 \\ \frac{z^{n+1} - \rho^n y^n}{\delta t} + \nabla \cdot (z^{n+1} u^{n+1}) = 0 \\ \rho^{n+1} = \varrho^{p,z}(p^{n+1}, z^{n+1}) \end{array} \right. \quad (\text{IV.9})$$

3 - solve for  $y^{n+1}$

$$\frac{\rho^{n+1} y^{n+1} - z^{n+1}}{\delta t} + \nabla \cdot (\rho^{n+1} y^{n+1} (1 - y^{n+1}) u_r^{n+1}) = \nabla \cdot (D \nabla y^{n+1}) \quad (\text{IV.10})$$

The first step consists in a classical semi-implicit solution of the momentum balance equation to obtain a predicted velocity.

Step 2 is an original nonlinear pressure correction step, which couples the mass balance equation (second equation) with the transport terms of the gas mass balance equation (third equation). A new unknown is introduced in this step instead of the gas mass fraction, the partial gas density  $z$  given by  $z = \rho y$ . Thus, the equation of state must be reformulated to express the mixture density as a function of the partial gas density and of the pressure, which, from equation (IV.6), yields :

$$\rho = \varrho^{p,z}(p, z) = z \left( 1 - \frac{\rho_\ell a^2}{p} \right) + \rho_\ell \quad (\text{IV.11})$$

When the liquid and the gas densities are very different, this law presents much less steep variations than the relation linking the density and the mass fraction  $y$ , especially in the neighbourhood of  $y = 0$ ; this change of variable thus makes the resolution of this step much easier, and the overall algorithm more robust. In counterpart, it leads to split the gas mass balance equation : transport terms are dealt with in the present step, and the gas mass fraction is corrected in a next step (step 3) to take into account the drift terms. The pressure correction step would degenerate in the usual projection step as used in incompressible flows solvers if the density was constant (*i.e.*  $z = 0$ ). Taking (at the algebraic level, see section IV.2.4) the divergence of the first relation of (IV.9) and using the second one to eliminate the unknown velocity  $u^{n+1}$  yields a non-linear elliptic problem for the pressure. Solving at the same time this elliptic problem and the third equation by Newton's algorithm, we obtain the pressure and the gas mass fraction. Once the pressure is computed, the first relation yields the updated velocity and the fourth one gives the end-of-step density.

Finally, in the third step, the remaining terms of the gas mass balance are considered, and the end-of-step gas mass fraction is computed.

The motivations of this time discretization are the following ones : to keep the mass fraction  $y$  in the physical range  $[0, 1]$ , to allow the transport of phases interfaces without generating spurious pressure and velocity variations and to ensure the stability of (*i.e.* the conservation of the entropy by) the scheme. To show how this time splitting algorithm achieves these goals is the aim of the remainder of this paper.

### IV.2.2 Spatial discretization

Let  $\mathcal{M}$  be a decomposition of the domain  $\Omega$  either into convex quadrilaterals ( $d = 2$ ) or hexahedra ( $d = 3$ ) or in simplices. By  $\mathcal{E}$  and  $\mathcal{E}(K)$  we denote the set of all  $(d - 1)$ -edges  $\sigma$  of the

mesh and of the element  $K \in \mathcal{M}$  respectively. The set of edges included in the boundary of  $\Omega$  is denoted by  $\mathcal{E}_{\text{ext}}$  and the set of internal ones (*i.e.*  $\mathcal{E} \setminus \mathcal{E}_{\text{ext}}$ ) is denoted by  $\mathcal{E}_{\text{int}}$ . The decomposition  $\mathcal{M}$  is supposed to be regular in the usual sense of the finite element literature (e.g. [15]), and, in particular,  $\mathcal{M}$  satisfies the following properties :  $\Omega = \bigcup_{K \in \mathcal{M}} K$ ; if  $K, L \in \mathcal{M}$ , then  $\bar{K} \cap \bar{L}$  is reduced to the empty set, to a vertex or (if  $d = 3$ ) to a segment, or  $\bar{K} \cap \bar{L}$  is (the closure of) a common  $(d - 1)$ -edge of  $K$  and  $L$ , which is denoted by  $K|L$ . For each internal edge of the mesh  $\sigma = K|L$ ,  $n_{KL}$  stands for the normal vector of  $\sigma$ , oriented from  $K$  to  $L$ . By  $|K|$  and  $|\sigma|$  we denote the measure, respectively, of  $K$  and of the edge  $\sigma$ .

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 Babuska-Brezzi condition (*e.g.* [9]). 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 a coupling with a finite volume treatment of the other equations, as is proposed hereafter for the gas mass balance; this is the choice made here. The spatial discretization thus relies either on the so-called "rotated bilinear element"/ $P_0$  introduced by Rannacher and Turek [61] for quadrilateral or hexahedric meshes, or on the Crouzeix-Raviart element (see [18] for the seminal paper and, for instance, [24, p. 83–85] for a synthetic presentation) for simplicial meshes. The reference element  $\hat{K}$  for the rotated bilinear element is the unit  $d$ -cube (with edges parallel to the coordinate axes); the discrete functional space on  $\hat{K}$  is  $\tilde{Q}_1(\hat{K})^d$ , where  $\tilde{Q}_1(\hat{K})$  is defined as follows :

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

The reference element for the Crouzeix-Raviart is the unit  $d$ -simplex and the discrete functional space is the space  $P_1$  of affine polynomials. For both velocity elements used here, the degrees of freedom are determined by the following set of nodal functionals :

$$\{F_{\sigma,i}, \sigma \in \mathcal{E}(K), i = 1, \dots, d\}, \quad F_{\sigma,i}(v) = |\sigma|^{-1} \int_{\sigma} v_i \, d\gamma \quad (\text{IV.12})$$

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

$$W_h = \{ v_h \in L^2(\Omega) : v_h|_K \in W(K)^d, \forall K \in \mathcal{M}; \\ F_{\sigma,i}(v_h) \text{ continuous across each edge } \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d; \\ F_{\sigma,i}(v_h) = 0, \forall \sigma \in \mathcal{E}_{\text{ext}}, 1 \leq i \leq d \}$$

where  $W(K)$  is the space of functions on  $K$  generated by  $\tilde{Q}_1(\hat{K})$  through the  $Q_1$  mapping from  $\hat{K}$  to  $K$  for the Rannacher-Turek element and the space of affine functions on  $K$  for the Crouzeix-Raviart element. For both Rannacher-Turek and Crouzeix-Raviart discretizations, the pressure is approximated by the space  $L_h$  of piecewise constant functions :

$$L_h = \{q_h \in L^2(\Omega) : q_h|_K = \text{constant}, \forall K \in \mathcal{M}\}$$

Since only the continuity of the integral over each edge of the mesh is imposed, the velocities are 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  $c_i > 0$  independent of the mesh such that :

$$\forall p \in L_h, \quad \sup_{v \in W_h} \frac{\int_{\Omega, h} p \nabla \cdot v \, dx}{\|v\|_{1, b}} \geq c_i \|p - m(p)\|_{L^2(\Omega)}$$

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

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

From the definition (IV.12), each velocity degree of freedom can be univoquely 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  $i^{\text{th}}$  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 considered finite elements, reads :

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

where  $\varphi_\sigma$  is a scalar function.

Each degree of freedom for the pressure is associated to a mesh  $K$ , and the set of pressure degrees of freedom is denoted by  $\{p_K, K \in \mathcal{M}\}$ . As the pressure, the density  $\rho$ , the gas mass fraction  $y$  and the gas partial density  $z$  are approximated by piecewise constant functions over each element, and the associated sets of degrees of freedom are denoted by  $\{\rho_K, K \in \mathcal{M}\}$ ,  $\{y_K, K \in \mathcal{M}\}$  and  $\{z_K, K \in \mathcal{M}\}$  respectively.

### IV.2.3 Spatial discretization of the momentum balance equation

The main difficulty in the discretization of the momentum balance equation is to build a discrete convection operator which enjoys the discrete analogue of the kinetic energy relation, that is :

$$\int_{\Omega} \left[ \frac{\partial \rho u}{\partial t} + \nabla \cdot (\rho u \otimes u) \right] \cdot u \, dx = \frac{d}{dt} \int_{\Omega} \frac{1}{2} \rho |u|^2 \quad \text{provided that} \quad \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0$$

To this purpose, we follow an idea developped in [3], and already exploited for the same problem as here in [35]. The idea is to derive a finite-volume-like discretization of the convection operator, in order to apply the following result [31].

**Theorem IV.2.1 (Stability of a finite volume advection operator)**

Let  $(\rho_K^*)_{K \in \mathcal{M}}$  and  $(\rho_K)_{K \in \mathcal{M}}$  be two families of positive real numbers satisfying the following set of equations :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K - \rho_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} = 0 \quad (\text{IV.13})$$

where  $F_{\sigma,K}$  is a quantity associated to the edge  $\sigma$  and to the control volume  $K$  ; we suppose that, for any internal edge  $\sigma = K|L$ ,  $F_{\sigma,K} = -F_{\sigma,L}$ . Let  $(s_K^*)_{K \in \mathcal{M}}$  and  $(s_K)_{K \in \mathcal{M}}$  be two families of real numbers. The following stability property holds :

$$\sum_{K \in \mathcal{M}} z_K \left[ \frac{|K|}{\delta t} (\rho_K s_K - \rho_K^* s_K^*) + \sum_{\sigma=K|L} F_{\sigma,K} \frac{s_K + s_L}{2} \right] \geq \frac{1}{2} \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} [\rho_K s_K^2 - \rho_K^* s_K^{*2}]$$

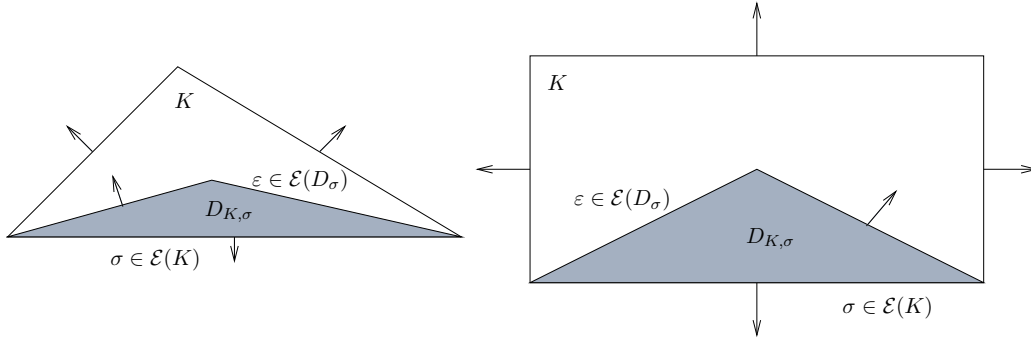


FIG. IV.1 – Diamond-cells for the Crouzeix-Raviart and Rannacher-Turek element.

To this purpose, we first define a control volume for each degree of freedom of the velocity, that is, in view of the discretization used here, around each barycenter of an internal edge. Let  $\sigma = K|L$  and  $D_{K,\sigma}$  be the conic volume having  $\sigma$  for basis and the mass center of  $K$  as additional vertex (see figure IV.1). The volume  $D_\sigma = D_{K,\sigma} \cup D_{L,\sigma}$  is referred to as the "diamond cell" associated to  $\sigma$  and  $D_{K,\sigma}$  is the half-diamond cell associated to  $\sigma$  and  $K$ . For the Crouzeix-Raviart element and, for the Rannacher-Turek element, when the mesh is a rectangle (in two dimensions) or a cuboid (in three dimensions), the integral of the shape function associated to the edge  $\sigma$  over the element  $K$  is the measure of the half-diamond cell  $D_{K,\sigma}$ . Thus, the application of the mass lumping to the terms of the form  $\rho u$  leads, in the equations associated to the velocity on the edge  $\sigma$ , to a discrete expression of the form  $\rho_\sigma u_\sigma$ , where  $\rho_\sigma$  results from an average of the values taken by the density in the two elements adjacent to  $\sigma$ , weighted by the measure of the half-diamonds :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad |D_\sigma| \rho_\sigma = |D_{K,\sigma}| \rho_K + |D_{L,\sigma}| \rho_L \quad (\text{IV.14})$$

where  $|D_\sigma|$  is the measure of the diamond cell  $D_\sigma$ ,  $|D_{K,\sigma}|$  and  $|D_{L,\sigma}|$  are the measure of the half-diamond cells associated respectively to  $\sigma$  and  $K$  and to  $\sigma$  and  $L$ . This lumped time derivative term naturally combines with a discretization of the advective term of the form :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad \text{advection term} \sim \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma} u_\varepsilon$$

where  $\mathcal{E}(D_\sigma)$  is the set of the edges of  $D_\sigma$ ,  $u_\varepsilon$  is a centered approximation of  $u$  on  $\varepsilon \in \mathcal{E}(D_\sigma)$  and  $F_{\varepsilon,\sigma}$  is a mass flux through  $\varepsilon$ . To proceed, we must now derive for this latter quantity an approximation which satisfies the compatibility condition (IV.13) of theorem IV.2.1 (in fact, the discrete mass balance over the diamond cells). Suppose that we are able to build, for any control volume  $K$ , a field  $\widetilde{\rho u}_K(x)$  such that  $\nabla \cdot \widetilde{\rho u}_K(x)$  remains constant inside the element  $K$  and that we take for the mass flux  $F_{\varepsilon,\sigma}$  through each diamond cell edge  $\varepsilon$  included in  $K$  :

$$F_{\varepsilon,\sigma} = \int_{\varepsilon} \widetilde{\rho u}_K(x) \cdot n_{\varepsilon,\sigma} \, d\gamma(x)$$

where  $n_{\varepsilon,\sigma}$  is the normal vector to  $\varepsilon$  outward  $D_\sigma$ . As the divergence of  $\widetilde{\rho u}_K$  is constant over  $K$ , it may be checked that, if the flux of  $\widetilde{\rho u}_K$  through each edge of  $K$  is the same as the mass flux used in a discrete mass balance over  $K$ , let say  $|\sigma| (\rho u)_\sigma \cdot n_\sigma$ , this mass balance is "carried over" the half-diamond cells  $D_{K,\sigma}$ , which, by summation over the two half-diamond cells, yields a compatibility condition of the desired form [35]. Such a field  $\widetilde{\rho u}_K(x)$  is derived for the Crouzeix-Raviart element by direct interpolation (*i.e.* using the standard expansion of the Crouzeix-Raviart elements) of the quantities  $((\rho u)_\sigma)_{\sigma \in \mathcal{E}(K)}$  :

$$\widetilde{\rho u}_K(x) = \sum_{\sigma \in \mathcal{E}(K)} \varphi_\sigma(x) (\rho u)_\sigma$$

For the Rannacher-Turek element, when the mesh is a rectangle or a cuboid, it is obtained by the following interpolation formula :

$$\widetilde{\rho u}_K(x) = \sum_{\sigma \in \mathcal{E}(K)} \alpha_\sigma(x \cdot n_\sigma) [(\rho u)_\sigma \cdot n_\sigma] n_\sigma$$

where the  $\alpha_\sigma(\cdot)$  are affine interpolation functions which are determined in such a way that the desired conditions hold, *i.e.* that the flux of  $\widetilde{\rho u}$  through each edge  $\sigma$  of  $K$  is  $|\sigma| (\rho u)_\sigma \cdot n_\sigma$ . Extension to more general grids is underway. Finally, since, in the proposed fractional step algorithm, the mass balance equation is considered only once the solution of the momentum balance has been computed, to obtain the desired compatibility condition (IV.13), we use the mass balance at the previous time step : the approximations of the density in the time derivative term are shifted of one time step and the quantities  $((\rho u)_\sigma)_{\sigma \in \mathcal{E}(K)}$  used to compute the mass fluxes  $F_{\varepsilon,\sigma}^n$  are chosen to be the mass fluxes obtained in the discrete mass balance at the previous time step. Since standard finite elements techniques are used to discretize the term  $\nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1})$ , this yields the following discrete momentum balance equation :

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

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

where the bilinear form  $a_d(\cdot, \cdot)$  represents the viscous term and is defined as follows :

$$\forall v \in W_h, \forall w \in W_h,$$

$$a_d(v, w) = \begin{cases} \mu \int_{\Omega,h} \left[ \nabla v : \nabla w + \frac{1}{3} \nabla \cdot v \nabla \cdot w \right] dx & \text{if (IV.3) holds (case of constant viscosity),} \\ \int_{\Omega,h} \tau(v) : \nabla w \, dx & \text{with } \tau \text{ given by (IV.2) otherwise.} \end{cases}$$

Note that, for Crouzeix-Raviart elements, a combined finite volume/finite element method similar to the technique employed here has already been analysed for a transient non-linear convection-diffusion equation by Feistauer and co-workers [1, 21, 29].

As a consequence of the stability of the convection operator, we have the following regularity result.

**Lemma IV.2.2 (Properties of the numerical scheme - velocity prediction)**

*Let us assume that the viscous term is dissipative (i.e.  $\forall v \in W_h$ ,  $a_d(v, v) \geq 0$ , which holds for the form of  $a_d(\cdot, \cdot)$  used in case of a constant viscosity); then the first step of the scheme, namely the velocity prediction step, has a unique solution.*

**Remark IV.2.3 (First time step)** *To ensure the compatibility condition (IV.13) at the first time step, a prediction step must be used to initialize the density :*

$$\frac{\rho^0 - \rho^{-1}}{\delta t} + \nabla \cdot (\rho^0 u^{-1}) = 0 \tag{IV.15}$$

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

**IV.2.4 Spatial discretization of the pressure correction step**

The discretization of the first equation of the pressure correction step (IV.9) is consistent with the momentum balance one, i.e. we use a mass lumping technique for the unsteady term and a standard finite element formulation for the gradient of the pressure increment :

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

Since the pressure is piecewise constant, the transposed of the discrete gradient operator takes the form of the finite volume standard discretization of the divergence based on the finite element mesh, thus the previous relation can be rewritten as follows :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \quad \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (u_\sigma^{n+1} - \tilde{u}_\sigma^{n+1}) + |\sigma| [(p_L^{n+1} - p_L^n) - (p_K^{n+1} - p_K^n)] n_{KL} = 0 \tag{IV.16}$$

Similarly, as the density is piecewise constant, the approximation of the time derivative of the density in the mass balance (second equation of (IV.9)) will also look as a finite volume term. This point suggests a finite volume discretization of this latter equation, which reads :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\varrho^{p,z}(p_K^{n+1}, z_K^{n+1}) - \rho_K^n) + \sum_{\sigma=K|L} F_{\sigma,K}^{n+1} = 0 \tag{IV.17}$$

To ensure the positivity of the density, we use an upwinding technique for the convection term, then the mass flux from  $K$  across  $\sigma = K|L$ ,  $F_{\sigma,K}^{n+1}$ , is expressed as follows :

$$F_{\sigma,K}^{n+1} = |\sigma| (\rho^{n+1} u^{n+1})|_\sigma \cdot n_\sigma = (v_{\sigma,K}^+)^{n+1} \varrho^{p,z}(p_K^{n+1}, z_K^{n+1}) - (v_{\sigma,K}^-)^{n+1} \varrho^{p,z}(p_L^{n+1}, z_L^{n+1})$$

where  $(v_{\sigma,K}^+)^{n+1}$  and  $(v_{\sigma,K}^-)^{n+1}$  stands respectively for  $\max(v_{\sigma,K}^{n+1}, 0)$  and  $-\min(v_{\sigma,K}^{n+1}, 0)$  with  $v_{\sigma,K}^{n+1} = |\sigma| u_\sigma^{n+1} \cdot n_{KL}$ .

Consistently with the mass balance equation, we use for the discretization of the third relation of (IV.9), *i.e.* the transport of the gas partial density  $z$ , a finite volume method with an upwind technique for the convection term  $\nabla \cdot (z u)$ . This yields the following discrete equation :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (z^{n+1} - \rho_K^n y_K^n) + \sum_{\sigma=K|L} (v_{\sigma,K}^+)^{n+1} z_K^{n+1} - (v_{\sigma,K}^-)^{n+1} z_L^{n+1} = 0 \quad (\text{IV.18})$$

In the following lemma, we state some properties of this pressure correction step which are obtained as a particular case of the existence theory presented in the appendix.

**Lemma IV.2.4 (Properties of the numerical scheme - pressure correction step)**

*Let the density of the liquid phase be constant and the gas phase obey the ideal gas law. Then, under the assumption that,  $\forall K \in \mathcal{M}$ ,  $\rho_K^n > 0$  and  $y_K^n \in (0, 1]$ , the system (IV.16)-(IV.18) has a solution, and any solution of this step is such that :*

$$\forall K \in \mathcal{M}, \quad \rho_K^{n+1} > 0, \quad p_K^{n+1} > 0, \quad z_K^{n+1} > 0 \quad \text{and} \quad \frac{z_K^{n+1}}{\rho_K^{n+1}} \in (0, 1]$$

Let us now turn to the practical solution of this pressure correction step. Keeping the same notation for the unknown functions and the vectors which gather their degrees of freedom, the algebraic formulation of this step reads :

$$\begin{cases} \frac{1}{\delta t} M_{\rho^n} (u^{n+1} - \tilde{u}^{n+1}) + B^t (p^{n+1} - p^n) = 0 \\ \frac{1}{\delta t} R (\varrho^{p,z}(p^{n+1}, z^{n+1}) - \rho^n) - B Q_{\rho^{n+1}}^{\text{up}} u^{n+1} = 0 \\ \frac{1}{\delta t} R (z^{n+1} - \rho^n y^n) - B Q_{z^{n+1}}^{\text{up}} u^{n+1} = 0 \end{cases} \quad (\text{IV.19})$$

In the first relation,  $M_{\rho^n}$  stands for the diagonal mass matrix weighted by the density at  $t^n$  (at edge center)  $\rho_\sigma^n$ , so the diagonal entry of  $M_{\rho^n}$  associated to the internal edge  $\sigma$  and the component  $i$  reads  $(M_{\rho^n})_{\sigma,i} = |D_\sigma| \rho_\sigma^n$ . The matrix  $B^t$  of  $\mathbb{R}^{N \times M}$ , where  $N = d \text{ card}(\mathcal{E}_{\text{int}})$  and  $M = \text{card}(\mathcal{M})$ , is associated to the gradient operator ; consequently, the matrix  $B$  is associated to the opposite of the divergence operator. In the second and in the third relation,  $Q_{w^{n+1}}^{\text{up}}$  (with  $w = \rho$  or  $w = z$ ) is a diagonal matrix, the entry of which corresponding to an edge  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , is obtained by simply taking  $w$  at  $t^{n+1}$  in the element located upstream of  $\sigma$  with respect to  $u^{n+1}$ , *i.e.* either  $w_K^{n+1}$  or  $w_L^{n+1}$ . The matrix  $R$  is diagonal and, for any  $K \in \mathcal{M}$ , its entry  $R_K$  is the measure of the element  $K$ .

The elliptic problem for the pressure is obtained by multiplying the first relation of (IV.19) by  $B Q_{\rho^{n+1}}^{\text{up}} (M_{\rho^{n+1}})^{-1}$  and using the second one. This equation reads :

$$L p^{n+1} + \frac{1}{\delta t^2} R \varrho^{p,z}(p^{n+1}, y^{n+1}) = L p^n + \frac{1}{\delta t^2} R \rho^n + \frac{1}{\delta t} B Q_{\rho^{n+1}}^{\text{up}} \tilde{u}^{n+1} \quad (\text{IV.20})$$

where, as seen in [31],  $L = B Q_{\rho^{n+1}}^{\text{up}} (M_{\rho^n})^{-1} B^t$  can be equivalently evaluated in the "finite volume way" by the following relation, valid for each element  $K$  :

$$(L p^{n+1})_K = \sum_{\sigma=K|L} \frac{\rho_{\text{up},\sigma}^{n+1}}{\rho_\sigma^n} \frac{|\sigma|^2}{|D_\sigma|} (p_K - p_L)$$

where  $\rho_{\text{up},\sigma}$  stands for the upwind density associated to the edge  $\sigma$ . One recognizes in this relation a usual finite volume diffusion operator, with a particular diffusion coefficient which, for instance, can be evaluated for rectangular parallelepipedic control volumes as  $d \rho_{\text{up},\sigma}^{n+1}/\rho_{\sigma}^n$ . The factor  $d$  should be suppressed to be consistent with what would be obtained by a finite volume discretization of this elliptic equation, if this latter was derived in the time semi-discrete setting : this fact is linked with the well-known non-consistency of the Rannacher-Turek or Crouzeix-Raviart discretization of the Darcy problem.

Then, equation (IV.20) is solved at the same time as the third equation of (IV.19) by a Newton algorithm. Denoting by the index  $k$  the Newton iterate, once  $p_{k+1}^{n+1}$  is known, the first relation of (IV.19) gives the updated value of the velocity :

$$u_{k+1}^{n+1} = \tilde{u}^{n+1} - \delta t (M_{\rho^n})^{-1} B^t (p_{k+1}^{n+1} - p^n) \quad (\text{IV.21})$$

Since, to preserve the positivity of the density, we need to use in the mass balance the value of the density upwinded with respect to  $u^{n+1}$ , equations (IV.20) and (IV.21) are not decoupled, by contrast with what happens in usual projection methods. They are thus solved sequentially, performing the upwinding with respect to  $u_k^{n+1}$  in (IV.20) and then updating the velocity by (IV.21), up to convergence.

#### IV.2.5 Spatial discretization of the correction step for $y$

So as to be consistent with the discretization of the first part of the gas mass balance, the correction step for  $y$  is discretized by the finite volume method, and the resulting discrete problem reads :

$$|K| \frac{\rho_K^{n+1} y_K^{n+1} - z_K^{n+1}}{\delta t} + \sum_{\sigma=K|L} (G_{\sigma,K}^{n+1})^+ g(y_K^{n+1}, y_L^{n+1}) - (G_{\sigma,K}^{n+1})^- g(y_L^{n+1}, y_K^{n+1}) + D \sum_{\sigma=K|L} \frac{|\sigma|}{d_{\sigma}} (y_K^{n+1} - y_L^{n+1}) = 0 \quad (\text{IV.22})$$

In this relation, for all edge  $\sigma = K|L$ ,  $d_{\sigma}$  is the Euclidean distance between two points  $x_K$  and  $x_L$  of the adjacent meshes  $K$  and  $L$ , supposed to be such that the segment  $[x_K, x_L]$  is perpendicular to  $K|L$ . These points may be defined as follows : if the control volume  $K$  is a rectangle or a cuboid,  $x_K$  is the barycenter of  $K$  ; if the control volume  $K$  is a simplex,  $x_K$  is the circumcenter of the vertices of  $K$ . Note that, in this latter case, the condition  $x_K \in K$  implies some geometrical constraints for  $K$ . Of course, in the cases where the diffusion coefficient  $D = 0$ , these limitations are useless.

The quantities  $(G_{\sigma,K}^{n+1})^+$  and  $(G_{\sigma,K}^{n+1})^-$  are defined as  $(G_{\sigma,K}^{n+1})^+ = \max(G_{\sigma,K}^{n+1}, 0)$  and  $(G_{\sigma,K}^{n+1})^- = -\min(G_{\sigma,K}^{n+1}, 0)$  respectively, with  $G_{\sigma,K}^{n+1}$  given by :

$$G_{\sigma,K}^{n+1} = \rho_{\sigma,\text{up}}^{n+1} \int_{\sigma} u_r^{n+1} \cdot n_{\sigma}$$

where  $\rho_{\sigma,\text{up}}^{n+1}$  stands for  $\rho_{\sigma,\text{up}}^{n+1} = \rho_K^{n+1}$  if  $u_{\sigma}^{n+1} \cdot n_{\sigma} \geq 0$  and  $\rho_{\sigma,\text{up}}^{n+1} = \rho_L^{n+1}$  otherwise. Note that this upwind choice with respect to  $u^{n+1}$  has no theoretical justification : in fact, the developments of this paper hold with any discretization for this density, and we use here the same discretization as in the mass balance simply to make the computer implementation easier. The function  $g(\cdot, \cdot)$  corresponds to an approximation of  $\varphi(y) = \max[y(1-y), 0]$  by a monotone numerical flux function. Let us recall the definition of this latter notion [26] :

**Definition IV.2.5 (Monotone numerical flux function)**

Let the function  $g(\cdot, \cdot) \in C(\mathbb{R}^2, \mathbb{R})$  satisfy the following assumptions :

1.  $g(a_1, a_2)$  is non-decreasing with respect to  $a_1$  and non-increasing with respect to  $a_2$ , for any real numbers  $a_1$  and  $a_2$ ,
2.  $g(\cdot, \cdot)$  is Lipschitz continuous with respect to both variables over  $\mathbb{R}$ ,
3.  $g(a_1, a_1) = \varphi(a_1)$ , for any  $a_1 \in \mathbb{R}$ .

Then  $g(\cdot, \cdot)$  is said to be a monotone numerical flux function for  $\varphi(\cdot)$ .

Several choices are possible for the numerical flux function  $g(\cdot, \cdot)$  and we refer to [26] for some examples and references. We adopt here the following simple flux-splitting formula :

$$g(a_1, a_2) = g_1(a_1) + g_2(a_2)$$

where  $g_1(a_1) = a_1$  if  $a_1 \in [0, 1]$  and  $g_1(a_1) = 0$  otherwise, and  $g_2(a_2) = -(a_2)^2$  if  $a_2 \in [0, 1]$  and  $g_2(a_2) = 0$  otherwise. Note that this choice does not exactly match the definition, as neither  $g_1(\cdot)$  nor  $g_2(\cdot)$  are continuous at  $a_1 = 1$ . However, this is unimportant, as one can prove, even in this case, that the solution  $y$  remains in the interval  $(0, 1]$ , as stated in the following lemma which is a weaker version of the result proven in [35, section 2].

**Lemma IV.2.6 (Existence and uniqueness for a discrete solution)**

Let us suppose that,  $\forall K \in \mathcal{M}$ ,  $\rho_K^{n+1} > 0$  and  $z_K^{n+1}/\rho_K^{n+1} \in (0, 1]$ . Then, there exists a unique solution to the considered discrete problem (IV.22), and this solution verifies  $y_K^{n+1} \in (0, 1]$ ,  $\forall K \in \mathcal{M}$ .

**IV.2.6 Some properties of the scheme**

The following theorem gathers some properties of the scheme, which are essentially straightforward consequences of lemmas IV.2.2, IV.2.4 and IV.2.6.

**Theorem IV.2.7 (Properties of the scheme)**

Let the density of the liquid phase be constant and the gas phase obey the ideal gas law. We suppose that the viscous term is dissipative (i.e.  $\forall v \in W_h$ ,  $a_d(v, v) \geq 0$ ). In addition, we assume that the initial density is positive and the initial gas mass fraction belongs to the interval  $(0, 1]$ . Then there exists a solution  $(u^n)_{1 \leq n \leq N}$ ,  $(p^n)_{1 \leq n \leq N}$ ,  $(\rho^n)_{1 \leq n \leq N}$ ,  $(z^n)_{1 \leq n \leq N}$  and  $(y^n)_{1 \leq n \leq N}$  to the scheme which enjoys the following properties, for all  $n \leq N$  :

– the unknowns lie in their physical range :

$$\forall K \in \mathcal{M}, \quad \rho_K^n > 0, \quad z_K^n > 0, \quad p_K^n > 0, \quad y_K^n \in (0, 1]$$

– the total mass, the gas mass and, if  $f_v = 0$ , the integral of the momentum are conserved :

$$\begin{aligned} \sum_{K \in \mathcal{M}} |K| \rho_K^n &= \sum_{K \in \mathcal{M}} |K| \rho_K^0 \\ \sum_{K \in \mathcal{M}} |K| z_K^n &= \sum_{K \in \mathcal{M}} |K| \rho_K^n y_K^n = \sum_{K \in \mathcal{M}} |K| \rho_K^0 y_K^0 \end{aligned}$$

We now turn to another feature of the scheme, which, from numerical experiments, seems to be crucial for the robustness of the algorithm. Let us suppose until the end of this section that the drift velocity  $u_r$ , the diffusive coefficient  $D$  and the forcing term  $f_v$  are set to zero. In addition, we momentarily forget the boundary condition, *i.e.* we reason as if the problem were posed in  $\mathbb{R}^n$ . Then the continuous problem enjoys the following property : if the initial velocity and the initial pressure are constant, let say  $u = u_0$  and  $p = p_0$  respectively, then they remain constant throughout the computational time interval, while  $\rho$  or  $z$  are transported by this (constant) velocity ; this solution corresponds to the transport of the contact discontinuity of the underlying hyperbolic system, the wave structure of which is quite similar to the Euler equations one [39]. The objective of the subsequent development is to prove that the numerical scheme considered in this paper presents the same behaviour : if, at the initial time,  $u_K^0 = u_0$  and  $p_K^0 = p_0$  for all  $K \in \mathcal{M}$ , then  $p_K^{n+1} = p_0$  and  $u_K^{n+1} = u_0$ , for all  $K \in \mathcal{M}$  and  $n < N$ .

Let us assume that, at time  $t = t_n$ , the velocity  $u^n$  and the pressure  $p^n$  take the constant value  $u_0$  and  $p_0$  respectively. We are now going to check that there exists a solution  $u^{n+1}$ ,  $p^{n+1}$ ,  $z^{n+1}$  and  $y^{n+1}$  to the scheme such that  $u^{n+1} = u_0$  and  $p^{n+1} = p_0$ . The discrete momentum balance equation reads, with a zero forcing term :

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

$$\frac{|D_\sigma|}{\delta t} (\rho_\sigma^n \tilde{u}_{\sigma,i}^{n+1} - \rho_\sigma^{n-1} u_{\sigma,i}^n) + \sum_{\substack{\varepsilon \in \mathcal{E}(D_\sigma), \\ \varepsilon = D_\sigma | D_{\sigma'}}} \frac{1}{2} F_{\varepsilon,\sigma}^n (\tilde{u}_{\sigma,i}^{n+1} + \tilde{u}_{\sigma',i}^{n+1}) - \int_{\Omega,h} p^n \nabla \cdot \varphi_\sigma^{(i)} dx = 0$$

Replacing  $u^n$  and  $p^n$  by  $u_0$  and  $p_0$  respectively and taking  $\tilde{u}_\sigma^{n+1} = u_0$  for all  $\sigma \in \mathcal{E}_{\text{int}}$ , this system becomes :

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L, \quad u_0 \left[ \frac{|D_\sigma|}{\delta t} (\rho_\sigma^n - \rho_\sigma^{n-1}) + \sum_{\varepsilon \in \mathcal{E}(D_\sigma)} F_{\varepsilon,\sigma}^n \right] = 0$$

which is verified thanks to the equivalence between mass balances over primal and dual meshes, as explained in section IV.2.3. We now turn to the pressure correction step, which we recall :

$$\left\{ \begin{array}{l} \frac{|D_\sigma|}{\delta t} \rho_\sigma^n (u_\sigma^{n+1} - \tilde{u}_\sigma^{n+1}) + |\sigma| [(p_L^{n+1} - p_L^n) - (p_K^{n+1} - p_K^n)] n_{KL} = 0, \quad \forall \sigma \in \mathcal{E}_{\text{int}}, \sigma = K|L \\ \frac{|K|}{\delta t} [\varrho^{p,z}(p_K^{n+1}, z_K^{n+1}) - \rho_K^n] \\ \quad + \sum_{\sigma \in \mathcal{E}(K)} \left[ (v_{\sigma,K}^+)^{n+1} \varrho^{p,z}(p_K^{n+1}, z_K^{n+1}) - (v_{\sigma,K}^-)^{n+1} \varrho^{p,z}(p_L^{n+1}, z_L^{n+1}) \right] = 0, \quad \forall K \in \mathcal{M} \\ \frac{|K|}{\delta t} (z_K^{n+1} - z_K^n) + \sum_{\sigma \in \mathcal{E}(K)} \left[ (v_{\sigma,K}^+)^{n+1} z_K^{n+1} - (v_{\sigma,K}^-)^{n+1} z_L^{n+1} \right] = 0, \quad \forall K \in \mathcal{M} \end{array} \right.$$

Taking  $u_\sigma^{n+1} = u_0$  for all  $\sigma \in \mathcal{E}_{\text{int}}$  and  $p_K^{n+1} = p_0$  for all  $K \in \mathcal{M}$ , the left hand side of the first equation of this system vanishes. Next, following [32], we remark that, at fixed pressure, the equation of state giving the density  $\rho$  as a function of  $z$  becomes an affine function :

$$\rho = \varrho^{p,z}(p_0, z) = z \left( 1 - \frac{\rho_\ell a^2}{p_0} \right) + \rho_\ell$$



Introducing this relation in the mass balance equation, we obtain :

$$\begin{aligned} & \frac{|K|}{\delta t} \left[ (z_K^{n+1} - z_K^n) \left( 1 - \frac{\rho_\ell a^2}{p_0} \right) \right] \\ & + \sum_{\sigma \in \mathcal{E}(K)} (v_{\sigma,K}^+)^{n+1} \left[ z_K^{n+1} \left( 1 - \frac{\rho_\ell a^2}{p_0} \right) + \rho_\ell \right] - (v_{\sigma,K}^-)^{n+1} \left[ z_L^{n+1} \left( 1 - \frac{\rho_\ell a^2}{p_0} \right) + \rho_\ell \right] = 0 \end{aligned}$$

which can be recast as :

$$\left( 1 - \frac{\rho_\ell a^2}{p_0} \right) \left[ \frac{|K|}{\delta t} (z_K^{n+1} - z_K^n) + \sum_{\sigma \in \mathcal{E}(K)} \left[ (v_{\sigma,K}^+)^{n+1} z_K^{n+1} - (v_{\sigma,K}^-)^{n+1} z_L^{n+1} \right] \right] + \rho_\ell \sum_{\sigma \in \mathcal{E}(K)} v_{\sigma,K}^{n+1} = 0$$

which, as the last term vanishes for  $u^{n+1} = u_0$ , is exactly the same equation as the gas mass balance. Thus,  $u^{n+1} = u_0$ ,  $p^{n+1} = p_0$ ,  $z^{n+1}$  given by this latter equation and  $y^{n+1}$  satisfying the correction step (which, for  $u_r = 0$  and  $D = 0$  becomes  $\rho^{n+1} y^{n+1} = z^{n+1}$ ) is a solution to the scheme. Consequently, provided that the solution is unique, the algorithm indeed preserves constant pressure and velocity through moving interfaces between phases, and transports this interface with this constant velocity.

**Remark IV.2.8 (Boundary conditions)** *The same property holds with a bounded computational domain when prescribing on the boundary either  $u = u_0$  or a Neumann condition compatible with  $u = u_0$  and  $p = p_0$ ; this fact has been confirmed by numerical experiments, although we leave its proof beyond the scope of this presentation, to avoid the technicalities of the description of these latter discrete boundary conditions.*

**Remark IV.2.9 (On the choice of coupling the mass balance and the gas mass balance equations)** *As in [35], one may be tempted, especially for computing efficiency reasons, to use a fully fractional step algorithm, i.e. to solve all the equations sequentially. The central argument of the preceding development is that, with a fixed pressure, the quantity  $py$  is affine with respect to  $\rho$ , and this fact originates from the particular form of the equation of state. Thus, for this argument to hold, it is mandatory for the density in the product  $py$  to be given by the equation of state  $\rho = \rho^{p,y}(y, p^*)$ , where only the pressure may be taken at the previous time step or at the previous stage of the algorithm (indeed, when checking as below that the interface is transported with a fixed pressure,  $p$  will be considered constant, in particular with respect to time). Hence, the transport terms in the gas mass balance should read, in the time semi-discrete setting :*

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

*But, in this case, the compatibility condition which yields a maximum principle for the advection operator, which here would read :*

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

*does not hold. So it seems that an algorithm keeping  $y$  within its physical bounds and transporting the interface at constant pressure and velocity necessarily couples the mixture and the gas mass balance.*

### IV.3 The stability induced by the pressure forces work

The aim of this section is to prove that the discretization at hand satisfies a stability bound which can be seen as the discrete analogue of the following equation :

$$- \int_{\Omega} p(x) \nabla \cdot u(x) \, dx = \frac{d}{dt} \int_{\Omega} f(x) \, dx \quad (\text{IV.23})$$

where  $f(x)$  stands for the volumetric free energy of the mixture. The role played by this estimate in the theory which is developed here is twofold. First, it provides an *a priori* bound for a class of discrete problems including the pressure correction step, which is the corner stone to prove the existence of a solution ; this development is presented in appendix. Second, it is crucial to derive stability results for the scheme.

Throughout this section, we suppose that both the drift velocity and the gas fraction diffusion vanishes, so the total mass and the gas mass balance equations simply read :

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) = 0 \quad (\text{IV.24})$$

$$\frac{\partial z}{\partial t} + \nabla \cdot (zu) = 0 \quad (\text{IV.25})$$

Of course, stability results for the complete problem will *in fine* depend on the fact that the neglected terms in (IV.25) are dissipative with respect to the free energy ; this point will be treated further.

This section is organized as follows. First, we prove this estimate in a general setting, *i.e.* without specifying the equation of state for the fluid. Then we explain as this theory applies to the case specifically adressed here, namely a constant density fluid and a gaseous phase obeying the ideal gas law.

#### IV.3.1 Abstract estimates

The formal computation which allows to derive estimate (IV.23) in the continuous setting is the following. The first assumption is that, through the (system of) equation(s) of state, the specific free energy can be expressed as a function of the mixture density and the gas partial density, which we write  $f = f(\rho, z)$ . Then multiplying the mass balance equation by the derivative of  $f$  with respect to  $\rho$ , the gas mass balance equation by the derivative of  $f(\cdot, \cdot)$  with respect to  $z$  and finally summing these relations, we obtain :

$$\frac{\partial f}{\partial \rho} \left[ \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho u) \right] + \frac{\partial f}{\partial z} \left[ \frac{\partial z}{\partial t} + \nabla \cdot (zu) \right] = 0$$

which yields :

$$\frac{\partial}{\partial t} f(\rho(x, t), z(x, t)) + \frac{\partial f}{\partial \rho} \nabla \cdot (\rho u) + \frac{\partial f}{\partial z} \nabla \cdot (zu) = 0$$

Developping the divergence terms, we get :

$$\frac{\partial}{\partial t} f(\rho(x, t), z(x, t)) + u \cdot \left[ \frac{\partial f}{\partial \rho} \nabla \rho + \frac{\partial f}{\partial z} \nabla z \right] + \nabla \cdot u \left[ \rho \frac{\partial f}{\partial \rho} + z \frac{\partial f}{\partial z} \right] = 0 \quad (\text{IV.26})$$

The second term of this relation is equal to  $u \cdot \nabla f(\rho(x, t), z(x, t))$ . Adding and subtracting  $f \nabla \cdot u$ , we thus have :

$$\frac{\partial}{\partial t} f(\rho(x, t), z(x, t)) + \nabla \cdot (f(\rho(x, t), z(x, t)) u) + \nabla \cdot u \left[ \rho \frac{\partial f}{\partial \rho} + z \frac{\partial f}{\partial z} - f \right] = 0 \quad (\text{IV.27})$$

Since the integral of  $\nabla \cdot (f(\rho(x, t), z(x, t)) u)$  over the computational domain is zero thanks to the fact that the velocity is supposed to vanish at the boundary, this equation is the relation we are seeking, provided that the free energy is such that the following relation holds :

$$\rho \frac{\partial f}{\partial \rho} + z \frac{\partial f}{\partial z} - f = p$$

We are going now to reproduce this computation at the discrete level.

**Theorem IV.3.10 (Stability due to the pressure work)**

Let  $\mathcal{C}$  be an open convex subset of  $\mathbb{R}^2$  and  $f(\cdot, \cdot)$  be a convex continuously differentiable function from  $\mathcal{C}$  to  $\mathbb{R}$ . We suppose that  $(\rho_K)_{K \in \mathcal{M}}$ ,  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(z_K)_{K \in \mathcal{M}}$  and  $(z_K^*)_{K \in \mathcal{M}}$  are four families of real numbers such that,  $\forall K \in \mathcal{M}$ ,  $(\rho_K, z_K) \in \mathcal{C}$ ,  $(\rho_K^*, z_K^*) \in \mathcal{C}$  and the following relations hold :

$$\begin{cases} \frac{|K|}{\delta t} (\rho_K - \rho_K^*) + \sum_{\sigma=K|L} v_{\sigma,K} \rho_\sigma = 0 \\ \frac{|K|}{\delta t} (z_K - z_K^*) + \sum_{\sigma=K|L} v_{\sigma,K} z_\sigma = 0 \end{cases} \quad (\text{IV.28})$$

where  $\rho_\sigma$  and  $z_\sigma$  are given by  $\rho_\sigma = \rho_K$  and  $z_\sigma = z_K$  if  $v_{\sigma,K} \geq 0$ ,  $\rho_\sigma = \rho_L$  and  $z_\sigma = z_L$  otherwise. Then the following estimate holds :

$$\sum_{K \in \mathcal{M}} -p_K \left[ \sum_{\sigma=K|L} v_{\sigma,K} \right] \geq \sum_{K \in \mathcal{M}} |K| \frac{f(\rho_K, z_K) - f(\rho_K^*, z_K^*)}{\delta t}$$

where the family of real numbers  $(p_K)_{K \in \mathcal{M}}$  is given by :

$$\forall K \in \mathcal{M}, \quad p_K = \rho_K \frac{\partial f}{\partial \rho}(\rho_K, z_K) + z_K \frac{\partial f}{\partial z}(\rho_K, z_K) - f(\rho_K, z_K)$$

**Proof.**

Let us multiply the first relation of (IV.28) by the derivative with respect to  $\rho$  of  $f(\cdot, \cdot)$ , the second one by the derivative with respect to  $z$  of  $f(\cdot, \cdot)$ , both being evaluated at  $(\rho_K, z_K)$ , and sum :

$$\begin{aligned} & \underbrace{\frac{|K|}{\delta t} \left[ (\rho_K - \rho_K^*) \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} + (z_K - z_K^*) \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} \right]}_{T_{\partial/\partial t}} \\ & + \underbrace{\left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} \sum_{\sigma=K|L} v_{\sigma,K} \rho_\sigma + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} \sum_{\sigma=K|L} v_{\sigma,K} z_\sigma}_{T_{div,K}} = 0 \end{aligned} \quad (\text{IV.29})$$

The second term of the previous relation,  $T_{div,K}$ , can be recast as :

$$T_{div,K} = \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} \left[ \sum_{\sigma=K|L} v_{\sigma,K} (\rho_\sigma - \rho_K) + \rho_K \sum_{\sigma=K|L} v_{\sigma,K} \right] + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} \left[ \sum_{\sigma=K|L} v_{\sigma,K} (z_\sigma - z_K) + z_K \sum_{\sigma=K|L} v_{\sigma,K} \right] \quad (IV.30)$$

This relation is the discrete equivalent to equation (IV.26) : up to the multiplication by  $1/|K|$ , the first summations in the first term and the second term at the right hand side are the analogues of  $u \cdot \nabla \rho$  and  $u \cdot \nabla z$  respectively, while the second summations are the analogues of  $\rho \nabla \cdot u$  and  $z \nabla \cdot u$  respectively. Adding and subtracting  $f(\rho_K, z_K)$ , we obtain a discrete equivalent of relation (IV.27) :

$$T_{div,K} = \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} \sum_{\sigma=K|L} v_{\sigma,K} (\rho_\sigma - \rho_K) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} \sum_{\sigma=K|L} v_{\sigma,K} (z_\sigma - z_K) + f(\rho_K, z_K) \sum_{\sigma=K|L} v_{\sigma,K} + \left[ \rho_K \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} + z_K \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} - f(\rho_K, z_K) \right] \sum_{\sigma=K|L} v_{\sigma,K}$$

In the last term, we recognize, as in the continuous setting,  $\rho_K \sum_{\sigma=K|L} v_{\sigma,K}$ . The process will be completed if we put the first three terms of the right hand side in the divergence form. To this end, let us sum up the term  $T_{div,K}$  over  $K \in \mathcal{M}$  and reorder the summation :

$$\sum_{K \in \mathcal{M}} T_{div,K} = \sum_{K \in \mathcal{M}} \rho_K \left[ \sum_{\sigma=K|L} v_{\sigma,K} \right] + \sum_{\sigma \in \mathcal{E}_{int}} T_{div,\sigma} \quad (IV.31)$$

where, if  $\sigma = K|L$  :

$$T_{div,\sigma} = v_{\sigma,K} \left[ \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} (\rho_\sigma - \rho_K) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} (z_\sigma - z_K) + f(\rho_K, z_K) - \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_L, z_L)} (\rho_\sigma - \rho_L) - \left( \frac{\partial f}{\partial z} \right)_{(\rho_L, z_L)} (z_\sigma - z_L) - f(\rho_L, z_L) \right]$$

In this relation, there are two possible choices for the orientation of  $\sigma$ , *i.e.*  $K|L$  or  $L|K$ ; we choose this orientation in order to have  $v_{\sigma,K} \geq 0$ . The function  $(\rho, z) \mapsto f(\rho, z)$  is by assumption continuously differentiable and convex on the convex set  $\mathcal{C}$  containing both  $(\rho_K, z_K)$  and  $(\rho_L, z_L)$ , so the technical lemma IV.3.11 hereafter applies and there exists  $(\bar{\rho}_\sigma, \bar{z}_\sigma)$  in the segment  $[(\rho_K, z_K), (\rho_L, z_L)]$  (itself included in  $\mathcal{C}$ ) such that :

$$\left. \begin{array}{l} \text{if } (\rho_K, z_K) \neq (\rho_L, z_L) : \\ \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} (\bar{\rho}_\sigma - \rho_K) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} (\bar{z}_\sigma - z_K) + f(\rho_K, z_K) \\ = \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_L, z_L)} (\bar{\rho}_\sigma - \rho_L) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_L, z_L)} (\bar{z}_\sigma - z_L) + f(\rho_L, z_L) \\ \text{otherwise : } \quad (\bar{\rho}_\sigma, \bar{z}_\sigma) = (\rho_K, z_K) = (\rho_L, z_L) \end{array} \right\} \quad (IV.32)$$

By definition, the choice  $(\rho_\sigma, z_\sigma) = (\bar{\rho}_\sigma, \bar{z}_\sigma)$  is such that the term  $T_{div,\sigma}$  vanishes, which means that the first three terms at the right hand side of equation (IV.30) are a conservative approximation of the quantity  $\nabla \cdot (fu)$  appearing in equation (IV.27), with the following expression for the flux :

$$\begin{aligned} F_{\sigma,K} &= f_\sigma v_{\sigma,K}, \quad \text{with :} \\ f_\sigma &= \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} (\bar{\rho}_\sigma - \rho_K) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} (\bar{z}_\sigma - z_K) + f(\rho_K, z_K) \\ &= \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_L, z_L)} (\bar{\rho}_\sigma - \rho_L) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_L, z_L)} (\bar{z}_\sigma - z_L) + f(\rho_L, z_L) \end{aligned}$$

Then the term  $T_{div,\sigma}$  can be rewritten as :

$$\begin{aligned} T_{div,\sigma} &= v_{\sigma,K} (\rho_\sigma - \bar{\rho}_\sigma) \left[ \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} - \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_L, z_L)} \right] \\ &\quad + v_{\sigma,K} (z_\sigma - \bar{z}_\sigma) \left[ \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} - \left( \frac{\partial f}{\partial z} \right)_{(\rho_L, z_L)} \right] \end{aligned}$$

With the orientation taken for  $\sigma$ , an upwind choice yields :

$$\begin{aligned} T_{div,\sigma} &= v_{\sigma,K} (\rho_K - \bar{\rho}_\sigma) \left[ \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_K, z_K)} - \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_L, z_L)} \right] \\ &\quad + v_{\sigma,K} (z_K - \bar{z}_\sigma) \left[ \left( \frac{\partial f}{\partial z} \right)_{(\rho_K, z_K)} - \left( \frac{\partial f}{\partial z} \right)_{(\rho_L, z_L)} \right] \end{aligned}$$

and, by the inequality of lemma IV.3.11 hereafter,  $T_{div,\sigma}$  can be seen to be non-negative. Let us now turn to  $T_{\partial/\partial t}$ . As the function  $(\rho, z) \mapsto f(\rho, z)$  is convex on the convex set  $\mathcal{C}$  and both  $(\rho_K, z_K)$  and  $(\rho_K^*, z_K^*)$  belong to  $\mathcal{C}$ , we have :

$$T_{\partial/\partial t} \geq |K| \frac{f(\rho_K, z_K) - f(\rho_K^*, z_K^*)}{\delta t} \quad (\text{IV.33})$$

Then, summing for  $K \in \mathcal{M}$  and using relations (IV.29), (IV.31) and (IV.33) concludes the proof.  $\blacksquare$

In the course of the preceding proof, we used the following technical lemma.

#### Lemma IV.3.11

Let  $\mathcal{C}$  be an open convex subset of  $\mathbb{R}^2$ ,  $f(\cdot, \cdot)$  be a convex continuously differentiable function from  $\mathcal{C}$  to  $\mathbb{R}$  and  $(\rho_1, z_1)$  and  $(\rho_2, z_2)$  be two distinct elements of  $\mathcal{C}$ . Then there exists  $\zeta \in [0, 1]$  such that  $(\bar{\rho}, \bar{z}) = (1 - \zeta) (\rho_1, z_1) + \zeta (\rho_2, z_2)$  satisfies the following relation :

$$\begin{aligned} f(\rho_1, z_1) + \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_1, z_1)} (\bar{\rho} - \rho_1) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_1, z_1)} (\bar{z} - z_1) = \\ f(\rho_2, z_2) + \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_2, z_2)} (\bar{\rho} - \rho_2) + \left( \frac{\partial f}{\partial z} \right)_{(\rho_2, z_2)} (\bar{z} - z_2) \end{aligned} \quad (\text{IV.34})$$

In addition, the following inequality holds :

$$T = (\rho_1 - \bar{\rho}) \left[ \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_1, z_1)} - \left( \frac{\partial f}{\partial \rho} \right)_{(\rho_2, z_2)} \right] + (z_1 - \bar{z}) \left[ \left( \frac{\partial f}{\partial z} \right)_{(\rho_1, z_1)} - \left( \frac{\partial f}{\partial z} \right)_{(\rho_2, z_2)} \right] \geq 0$$

**Proof.**

Let us consider the function  $g(\cdot)$  defined by :

$$\zeta \mapsto f((1 - \zeta) (\rho_1, z_1) + \zeta (\rho_2, z_2))$$

By assumption, the function  $g(\cdot)$  is defined over  $[0, 1]$ , convex and continuously differentiable. Moreover, it may be checked that equation (IV.34) equivalently reads :

$$g(0) + g'(0) \zeta = g(1) + g'(1) (\zeta - 1)$$

or, reordering terms :

$$[g'(1) - g'(0)] \zeta = g(0) - (g(1) - g'(1))$$

Since  $g(\cdot)$  is convex, if  $g'(1) = g'(0)$ , the function  $g(\cdot)$  is affine and  $g(0) - (g(1) - g'(1))$  vanishes, so the preceding relation is satisfied with any value of  $\zeta$ . Otherwise, the preceding relation allows to compute  $\zeta$  and, still by convexity of  $g(\cdot)$ , both  $g'(1) - g'(0)$  and  $g(0) - (g(1) - g'(1))$  are positive, and so is  $\zeta$ . Still in this second case, this relation equivalently reads :

$$[g'(1) - g'(0)] (\zeta - 1) = g(0) + g'(0) - g(1)$$

which, as  $g(0) + g'(0) - g(1)$  is negative, shows that  $\zeta \leq 1$ . Finally, the quantity  $T$  simply reads  $\zeta [g'(1) - g'(0)]$ , and is thus non-negative. ■

**Remark IV.3.12 (Discretization of the convective terms and conservation of the entropy)** *From the above computation, it appears that the choice of  $\bar{\rho}_\sigma$  and  $\bar{z}_\sigma$  defined by equation (IV.32), for the convective terms in the mass and the gas mass balance equations, is a convenient one to obtain an exact discrete counterpart of the continuous identity (IV.23), and thus, in fine, to build a scheme exactly conserving the entropy. The upwind choice yields a dissipation, and nothing can be said for the centered one.*

### IV.3.2 The case of a constant density liquid and an ideal gas

Let us suppose that  $\rho_\ell$  is constant and  $\rho_g$  is linearly increasing with the pressure :

$$\rho_g = \frac{1}{a^2} p$$

where  $a$  is a positive real number (from a physical point of view, it is the sound velocity in a pure gaseous isothermal flow). For any positive  $\rho$  and  $z$  such that  $z - \rho + \rho_\ell > 0$ , the relation (IV.6) giving the mixture density as a function of the gas mass fraction and the phasic densities may be recast under the following form :

$$\frac{1}{a^2} p = \varrho_g^{\rho, z}(\rho, z) = \frac{z \rho_\ell}{z + \rho_\ell - \rho} \tag{IV.35}$$

Let us define the volumetric free energy of the mixture by :

$$f(\rho, z) = a^2 z \log(\varrho_g^{\rho, z}(\rho, z)) \tag{IV.36}$$

This function is continuously differentiable over the convex subset of  $\mathbb{R}^2$  :

$$\mathcal{C} = \{(\rho, z) \in \mathbb{R}^2 \text{ s.t. } \rho > 0, z > 0, z - \rho + \rho_\ell > 0\} \tag{IV.37}$$

We are now going to show that it verifies the other two assumptions of theorem IV.3.10, namely that  $f(\cdot)$  is convex and satisfies the identity :

$$T_p = \rho \frac{\partial f}{\partial \rho} + z \frac{\partial f}{\partial z} - f = p$$

This latter relation can be proven without referring to the specific form of  $f(\cdot, \cdot)$ , thanks to of the following property, which would be verified by the volumetric free energy function associated to any mixture composed of a constant density liquid phase and a barotropic gaseous phase :

$$f(\rho, z) = z f_g(\varrho_g^{\rho, z}(\rho, z)) \quad \text{with : } f'_g(s) = \frac{\wp(s)}{s^2}$$

where  $\wp(\cdot)$  is the function giving the pressure as a function of the gas density (thus, in particular,  $f'(\rho_g) = p/\rho_g^2$ ) and  $f_g(\cdot)$  stands for the specific free energy of the gaseous phase. Developing the derivatives and using the definition of  $f_g(\cdot)$ , we get :

$$T_p = \rho z f'_g(\rho_g) \frac{\partial \varrho_g^{\rho, z}}{\partial \rho} + z^2 f'_g(\rho_g) \frac{\partial \varrho_g^{\rho, z}}{\partial z} + z f_g(\rho_g) - z f_g(\rho_g) = z \frac{p}{\rho_g^2} \left[ \rho \frac{\partial \varrho_g^{\rho, z}}{\partial \rho} + z \frac{\partial \varrho_g^{\rho, z}}{\partial z} \right] \quad (\text{IV.38})$$

From the expression (IV.35), we have :

$$\frac{\partial \varrho_g^{\rho, z}}{\partial \rho} = \frac{\rho_g^2}{\rho_\ell z} \quad \text{and} \quad \frac{\partial \varrho_g^{\rho, z}}{\partial z} = \frac{\rho_g^2(\rho_\ell - \rho)}{\rho_\ell z^2} \quad (\text{IV.39})$$

Substituting in (IV.38) leads to :

$$T_p = \rho_g^2 f'_g(\rho_g) = p$$

The convexity of  $f(\cdot, \cdot)$  is obtained from its explicit form :

$$f(\rho, z) = a^2 z \log \left( \frac{z \rho_\ell}{z + \rho_\ell - \rho} \right)$$

Differentiating twice this expression, we get :

$$\frac{\partial^2 f}{\partial \rho^2} = a^2 \frac{z}{(z + \rho_\ell + \rho)^2}, \quad \frac{\partial^2 f}{\partial z^2} = a^2 \frac{(\rho_\ell - \rho)^2}{z(z + \rho_\ell + \rho)^2}, \quad \frac{\partial^2 f}{\partial \rho \partial z} = \frac{\partial^2 f}{\partial z \partial \rho} = a^2 \frac{\rho_\ell - \rho}{(z + \rho_\ell + \rho)^2}$$

It is thus easy to check that the determinant of the Hessian matrix  $A$  of  $f(\cdot, \cdot)$  is zero while its trace is positive. One eigenvalue of  $A$  is thus zero and the second one is positive, and  $f(\cdot, \cdot)$  is convex.

## IV.4 Stability analysis

The aim of this section is to provide some results concerning the stability of (*i.e.* the conservation of the entropy by) the scheme considered in this paper. First (section IV.4.1), in the case where both the drift velocity  $u_r$  and the diffusion coefficient for the mass fraction of the dispersed phase  $D$  vanish (*i.e.* for the homogeneous model), we prove that the entropy (*i.e.* the usual entropy associated to the homogeneous model) is conserved by the scheme, up to a step of renormalization of the pressure which is precisely stated. Note that this step, which was implemented for monophasic flows in [31], could be added in the present scheme ; however, we have chosen not to consider it further than in this theoretical section, as, in practice, its beneficial effects were not clear. Second (section IV.4.2), we show that, as in the continuous case, if the drift velocity is proportional to the

gradient of the pressure, the drift term is dissipative with respect to the same entropy; for this property to hold, a particular discretization of the drift term has to be implemented.

In this section, we use the following discrete norm and semi-norm :

$$\begin{aligned} \forall v \in W_h, \quad \|v\|_{h,\rho}^2 &= \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma| \rho_\sigma |v_\sigma|^2 \\ \forall q \in L_h, \quad |q|_{h,\rho}^2 &= \sum_{\sigma \in \mathcal{E}_{\text{int}}, \sigma=K|L} \frac{1}{\rho_\sigma} \frac{|\sigma|^2}{|D_\sigma|} (q_K - q_L)^2 \end{aligned} \quad (\text{IV.40})$$

where  $\rho = (\rho_\sigma)_{\sigma \in \mathcal{E}_{\text{int}}}$  is a family of positive real numbers. The function  $\|\cdot\|_{h,\rho}^2$  defines a norm over  $W_h$ , and  $|\cdot|_{h,\rho}$  can be seen as a weighted version of the  $H^1$  semi-norm classical in the finite volume context [26]. The following relation links this latter semi-norm to the problem at hand :

$$\forall q \in L_h, \quad (\text{B M}_\rho^{-1} \text{B}^t q, q) = |q|_{h,\rho}^2 \quad (\text{IV.41})$$

where  $\text{B}^t$ ,  $\text{B}$  and  $\text{M}_\rho$  are the discrete gradient, (opposite of the) divergence and mass matrix defined in section IV.2.4. A proof of this equality can be found in [31, section 3.4].

#### IV.4.1 First case : $u_r = 0$ , $D = 0$

With a zero drift velocity and a zero diffusion coefficient, the numerical scheme at hand reads, in the time semi-discrete setting :

1 - solve for  $\tilde{u}^{n+1}$

$$\frac{\rho^n \tilde{u}^{n+1} - \rho^{n-1} u^n}{\delta t} + \nabla \cdot (\rho^n u^n \otimes \tilde{u}^{n+1}) + \nabla p^n - \nabla \cdot \tau(\tilde{u}^{n+1}) = f_v^{n+1} \quad (\text{IV.42})$$

2 - solve for  $p^{n+1}$ ,  $u^{n+1}$ ,  $\rho^{n+1}$  and  $z^{n+1}$

$$\left\{ \begin{aligned} \rho^n \frac{u^{n+1} - \tilde{u}^{n+1}}{\delta t} + \nabla(p^{n+1} - p^n) &= 0 \\ \frac{\varrho^{p,z}(p^{n+1}, z^{n+1}) - \rho^n}{\delta t} + \nabla \cdot (\varrho^{p,z}(p^{n+1}, z^{n+1}) u^{n+1}) &= 0 \\ \frac{z^{n+1} - \rho^n y^n}{\delta t} + \nabla \cdot (z^{n+1} u^n) &= 0 \\ \rho^{n+1} &= \varrho^{p,z}(p^{n+1}, z^{n+1}) \end{aligned} \right. \quad (\text{IV.43})$$

3 - solve for  $y^{n+1}$

$$\rho^{n+1} y^{n+1} = z^{n+1} \quad (\text{IV.44})$$



**Proposition IV.4.13 (A partial stability result)**

Let the density of the liquid phase be constant, the gas phase obey the ideal gas law and  $f(\rho, z)$  be the corresponding volumetric free energy of the mixture, defined by (IV.36). We suppose that the viscous term is dissipative (i.e.  $\forall v \in W_h, a_d(v, v) \geq 0$ ). In addition, we assume that the density  $\rho^n$  is positive and the gas mass fraction  $y^n$  belongs to the interval  $(0, 1]$ . Let  $\tilde{u}^{n+1}, u^{n+1}, p^{n+1}, z^{n+1}$  and  $\rho^{n+1}$  be a solution to equations (IV.42)-(IV.43), with a zero forcing term. Then the following bound holds :

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^n}^2 + \int_{\Omega} f(\rho^{n+1}, z^{n+1}) dx + \delta t a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) + \frac{\delta t^2}{2} |p^{n+1}|_{h, \rho^n}^2 \\ \leq \frac{1}{2} \|u^n\|_{h, \rho^{n-1}}^2 + \int_{\Omega} f(\rho^n, \rho^n y^n) dx + \frac{\delta t^2}{2} |p^n|_{h, \rho^n}^2 \end{aligned}$$

**Proof.**

Multiplying each equation of the first step of the scheme (IV.42) by the corresponding unknown (i.e the corresponding component of the velocity  $\tilde{u}^{n+1}$  on the corresponding edge  $\sigma$ ) and summing over the edges and the components yields, by virtue of theorem IV.2.1 :

$$\frac{1}{2\delta t} \|\tilde{u}^{n+1}\|_{h, \rho^n}^2 - \frac{1}{2\delta t} \|u^n\|_{h, \rho^{n-1}}^2 + a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) - \int_{\Omega, h} p^n \nabla \cdot \tilde{u}^{n+1} dx \leq 0 \quad (\text{IV.45})$$

On the other hand, the first relation of system equation (IV.43) reads, in the algebraic setting :

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

Reordering this relation and multiplying by  $M_{\rho^n}^{-1/2}$  (recall that  $M_{\rho^n}$  is diagonal), we obtain :

$$\frac{1}{\delta t} M_{\rho^n}^{1/2} u^{n+1} + M_{\rho^n}^{-1/2} B^t p^{n+1} = \frac{1}{\delta t} M_{\rho^n}^{1/2} \tilde{u}^{n+1} + M_{\rho^n}^{-1/2} B^t p^n$$

Squaring this relation gives :

$$\begin{aligned} \left( \frac{1}{\delta t} M_{\rho^n}^{1/2} u^{n+1} + M_{\rho^n}^{-1/2} B^t p^{n+1}, \frac{1}{\delta t} M_{\rho^n}^{1/2} u^{n+1} + M_{\rho^n}^{-1/2} B^t p^{n+1} \right) = \\ \left( \frac{1}{\delta t} M_{\rho^n}^{1/2} \tilde{u}^{n+1} + M_{\rho^n}^{-1/2} B^t p^n, \frac{1}{\delta t} M_{\rho^n}^{1/2} \tilde{u}^{n+1} + M_{\rho^n}^{-1/2} B^t p^n \right) \end{aligned}$$

which reads :

$$\begin{aligned} \frac{1}{\delta t^2} (M_{\rho^n} u^{n+1}, u^{n+1}) + (M_{\rho^n}^{-1} B^t p^{n+1}, B^t p^{n+1}) + \frac{2}{\delta t} (u^{n+1}, B^t p^{n+1}) = \\ \frac{1}{\delta t^2} (M_{\rho^n} \tilde{u}^{n+1}, \tilde{u}^{n+1}) + (M_{\rho^n}^{-1} B^t p^n, B^t p^n) + \frac{2}{\delta t} (\tilde{u}^{n+1}, p^n) \end{aligned}$$

Multiplying by  $\delta t/2$ , remarking that,  $\forall v \in W_h, (M_{\rho^n} v, v) = \|v\|_{h, \rho^n}^2$  and that, thanks to relation (IV.41),  $\forall q \in L_h, (M_{\rho^n}^{-1} B^t q, B^t q) = (B M_{\rho^n}^{-1} B^t q, q) = |q|_{h, \rho^n}^2$ , we get :

$$\begin{aligned} \frac{1}{2\delta t} \|u^{n+1}\|_{h, \rho^n}^2 + \frac{\delta t}{2} |p^{n+1}|_{h, \rho^n}^2 + (u^{n+1}, B^t p^{n+1}) \\ - \frac{1}{2\delta t} \|\tilde{u}^{n+1}\|_{h, \rho^n}^2 - \frac{\delta t}{2} |p^n|_{h, \rho^n}^2 - (\tilde{u}^{n+1}, B^t p^n) = 0 \end{aligned} \quad (\text{IV.46})$$

The quantity  $-(\tilde{u}^{n+1}, \mathbf{B}^t p^n)$  is nothing more than the opposite of the term  $\int_{\Omega, h} p^n \nabla \cdot \tilde{u}^{n+1} dx$  appearing in (IV.45), so summing (IV.45) and (IV.46) makes these terms disappear, leading to :

$$\begin{aligned} \frac{1}{2\delta t} \|u^{n+1}\|_{h, \rho^n}^2 - \frac{1}{2\delta t} \|u^n\|_{h, \rho^{n-1}}^2 + a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) \\ + \frac{\delta t}{2} |p^{n+1}|_{h, \rho^n}^2 - \frac{\delta t}{2} |p^n|_{h, \rho^n}^2 + (u^{n+1}, \mathbf{B}^t p^{n+1}) \leq 0 \end{aligned}$$

Finally,  $(u^{n+1}, \mathbf{B}^t p^{n+1})$  is precisely the pressure work which is likely to be bounded by the time derivative of the volumetric free energy of the mixture. We know from theorem IV.7.22 that any solution to the system (IV.43) satisfies  $\rho^{n+1} > 0$ ,  $z^{n+1} > 0$  and  $p^{n+1} > 0$ . In view of the different forms of the equation of state gathered in (IV.51), this implies that this solution belongs to the convex set  $\mathcal{C}$  defined by (IV.37), inside which the free energy is well defined, regular and convex. Hence, with this solution, theorem IV.3.10 indeed applies and we get :

$$\begin{aligned} \frac{1}{2\delta t} \|u^{n+1}\|_{h, \rho^n}^2 + a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) + \frac{\delta t}{2} |p^{n+1}|_{h, \rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} f(\rho^{n+1}, z^{n+1}) dx \\ \leq \frac{1}{2\delta t} \|u^n\|_{h, \rho^{n-1}}^2 + \frac{\delta t}{2} |p^n|_{h, \rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} f(\rho^n, \rho^n y^n) dx \end{aligned}$$

which concludes the proof. ■

**Theorem IV.4.14 (Stability of the scheme, case  $u_r = D = 0$ )**

*Let the density of the liquid phase be constant, the gas phase obey the ideal gas law and  $f(\rho, z)$  be the corresponding volumetric free energy of the mixture, defined by (IV.36). We suppose that the viscous term is dissipative (i.e.  $\forall v \in W_h, a_d(v, v) \geq 0$ ). In addition, we assume that the initial density is positive and the initial gas mass fraction belongs to the interval  $(0, 1]$ .*

*We now add to the scheme (IV.42)-(IV.44) the following renormalization step of the pressure, to be performed at the very beginning of the time step, before the velocity prediction step :*

$$\text{Solve for } \tilde{p}^{n+1} : \quad -\nabla \cdot \left( \frac{1}{\rho^n} \nabla \tilde{p}^{n+1} \right) = -\nabla \cdot \left( \frac{1}{\sqrt{\rho^n \rho^{n-1}}} \nabla p^n \right)$$

*or, in the algebraic setting :*

$$\mathbf{B} \mathbf{M}_{\rho^n}^{-1} \mathbf{B}^t \tilde{p}^{n+1} = \mathbf{B} \mathbf{M}_{\sqrt{\rho^n \rho^{n-1}}}^{-1} \mathbf{B}^t p^n$$

*Accordingly, the pressure used in the velocity prediction step must be changed to  $\tilde{p}^{n+1}$ .*

*Let  $(\tilde{u}^n)_{0 \leq n \leq N}$ ,  $(u^n)_{0 \leq n \leq N}$ ,  $(p^n)_{0 \leq n \leq N}$ ,  $(z^n)_{0 \leq n \leq N}$  and  $(\rho^n)_{0 \leq n \leq N}$  be the solution to this scheme, with a zero forcing term. Then the following entropy conservation result holds for  $0 \leq n < N$  :*

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^n}^2 + \int_{\Omega} f(\rho^{n+1}, z^{n+1}) dx + \delta t \sum_{k=1}^{n+1} a_d(\tilde{u}^k, \tilde{u}^k) + \frac{\delta t^2}{2} |p^{n+1}|_{h, \rho^n}^2 \\ \leq \frac{1}{2} \|u^0\|_{h, \rho^0}^2 + \int_{\Omega} z^0 f_g(\rho^0, z^0) dx + \frac{\delta t^2}{2} |p^0|_{h, \rho^0}^2 \end{aligned}$$

**Proof.**

By the same proof as for the scheme without the pressure renormalization step, we get :

$$\begin{aligned} \frac{1}{2\delta t} \|u^{n+1}\|_{h,\rho^n}^2 + a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) + \frac{\delta t}{2} |p^{n+1}|_{h,\rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} f(\rho^{n+1}, z^{n+1}) dx \\ \leq \frac{1}{2\delta t} \|u^n\|_{h,\rho^{n-1}}^2 + \frac{\delta t}{2} |\tilde{p}^{n+1}|_{h,\rho^n}^2 + \frac{1}{\delta t} \int_{\Omega} f(\rho^n, \rho^n y^n) dx \end{aligned}$$

and the conclusion follows by summing over the time steps, remarking that  $z^{n+1} = \rho^{n+1} y^{n+1}$  and, thanks to the renormalization step (see [31] for a detailed computation) :

$$|\tilde{p}^{n+1}|_{h,\rho^n}^2 \leq |p^n|_{h,\rho^{n-1}}^2$$

■

Note that a similar pressure renormalization step has already been introduced for variable density incompressible flows [37].

#### IV.4.2 Dissipativity of the drift term

We address in this section the case where the drift velocity is given by the Darcy-like closure relation (IV.7) :

$$u_r = \frac{1}{\lambda} (1 - \alpha_g) \alpha_g \frac{\varrho_g(p) - \rho_\ell}{\rho} \nabla p$$

In this relation,  $\lambda$  is a positive phenomenological coefficient and  $\alpha_g$  is the void fraction, which can be expressed as a function of the unknowns used in the scheme as  $\alpha_g = z/\varrho_g(p)$ . We recall the spatial discretization of the drift term in the correction step for the gas mass fraction  $y$ , namely  $\nabla \cdot (\rho y (1 - y) u_r)$ , given in section IV.2.5 :

$$\sum_{\sigma=K|L} G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K)$$

where the function  $g(\cdot, \cdot)$  corresponds to an approximation of  $\varphi(y) = \max[y(1 - y), 0]$  by a monotone numerical flux function,  $G_{\sigma,K}^+ = \max(G_{\sigma,K}, 0)$ ,  $G_{\sigma,K}^- = -\min(G_{\sigma,K}, 0)$  and  $G_{\sigma,K}$  is an approximation for the flux of  $\rho u_r$  through the edge  $\sigma = K|L$ . With the closure relation (IV.7) for  $u_r$ , a natural discretization for this quantity reads :

$$G_{\sigma,K} = |\sigma| \rho_{\sigma,\text{up}} \left[ \frac{\alpha_g (1 - \alpha_g)}{\lambda \rho} (\rho_g - \rho_\ell) \right]_{\sigma} (p_K - p_L) \quad (\text{IV.47})$$

where  $\rho_{\sigma,\text{up}}$  is a density on  $\sigma$ , for which, for practical implementation reasons, we choose an upwind approximation with respect to the mean velocity  $u$ . The goal of this section is to show that it is possible to approximate :

$$\left[ \frac{\alpha_g (1 - \alpha_g)}{\lambda \rho} (\rho_g - \rho_\ell) \right]_{\sigma}$$

in such a way that this drift term is dissipative with respect to the entropy of the system.

We begin this section by stating a consequence of the equation of state for the mixture which is central to the present development.

**Lemma IV.4.15**

Let the density of the liquid phase be constant, the gas phase obey the ideal gas law,  $f(\rho, z)$  be the corresponding volumetric free energy of the mixture, defined by (IV.36), and  $h(\rho, z)$  be the partial derivative of  $f(\cdot, \cdot)$  with respect to the second variable  $z$ . Then the following results hold :

1.  $h(\cdot, \cdot)$  only depends on the pressure, i.e. there exists a function  $h_p(\cdot)$  such that, for  $\rho$  and  $z$  in the convex set  $\mathcal{C}$  defined by (IV.37),  $h(\rho, z) = h_p(\wp(\rho, z))$ , where  $\wp(\cdot, \cdot)$  is the function giving the pressure as a function of  $\rho$  and  $z$  :

$$p = \wp(\rho, z) = a^2 \varrho_g^{\rho, z}(\rho, z) = a^2 \frac{z \rho_\ell}{z + \rho_\ell - \rho}$$

2. the derivative of  $h_p(\cdot)$  is given by :

$$h'_p(p) = \frac{\rho_\ell - \rho_g(p)}{\rho_\ell \rho_g(p)}$$

3. for any positive real numbers  $p_1$  and  $p_2$  such that  $p_1 < p_2$ , there exists  $p_{1,2} \in [p_1, p_2]$  such that :

$$h'_p(p_{1,2}) \frac{h_p(p_1) - h_p(p_2)}{p_1 - p_2} \geq 0$$

**Proof.**

As  $(\rho, z) \in \mathcal{C}$ , the pressure or, equivalently, the gas density  $\rho_g$  can be expressed as a function of  $(\rho, z)$  by  $\rho_g = \varrho_g^{\rho, z}(\rho, z)$ . By the definition of  $f(\cdot, \cdot)$ , we thus have :

$$h(\rho, z) = \frac{\partial f}{\partial z} = f_g(\rho_g) + z \frac{\partial f_g}{\partial z} = a^2 \log\left(\frac{p}{a^2}\right) + z \frac{\partial f_g}{\partial \rho_g} \frac{\partial \varrho_g^{\rho, z}}{\partial z}$$

Then using the expression (IV.39) of the derivative of  $\varrho_g^{\rho, z}(\cdot, \cdot)$  with respect to the second variable, we get :

$$h(\rho, z) = a^2 \log\left(\frac{p}{a^2}\right) + p \frac{\rho_\ell - \rho}{\rho_\ell z}$$

Using the fact that  $\rho = (1 - \alpha_g)\rho_\ell + \alpha_g \rho_g$  and thus  $\rho_\ell - \rho = \alpha_g (\rho_\ell - \rho_g) = \frac{z}{\rho_g} (\rho_\ell - \rho_g)$ , we have :

$$h(\rho, z) = a^2 \log\left(\frac{p}{a^2}\right) + p \frac{\rho_\ell - \rho_g}{\rho_\ell \rho_g}$$

By definition of  $\rho_g$ , i.e.  $\rho_g = p/a^2$ , we thus get :

$$h(\rho, z) = a^2 \left[ \log\left(\frac{p}{a^2}\right) + \frac{\rho_\ell - p/a^2}{\rho_\ell} \right] = h_p(p)$$

Taking the derivative of this relation yields the desired expression for  $h'_p(\cdot)$  and, as  $h_p(\cdot)$  is continuously differentiable in  $[p_1, p_2]$ , the existence of  $p_{1,2}$  follows by Lagrange's theorem. ■

We are now in position to state and prove the following stability result.

**Proposition IV.4.16**

Let the density of the liquid phase be constant, the gas phase obey the ideal gas law and  $f(\rho, z)$  be the corresponding volumetric free energy of the mixture, defined by (IV.36). Let  $\mathcal{C}$  be the convex set defined by (IV.37) and  $(\rho_K)_{K \in \mathcal{M}}$ ,  $(y_K)_{K \in \mathcal{M}}$  and  $(z_K)_{K \in \mathcal{M}}$  be such that,  $\forall K \in \mathcal{M}$ ,  $(\rho_K, \rho_K y_K) \in \mathcal{C}$ ,  $(\rho_K, z_K) \in \mathcal{C}$ , and the following relation is satisfied :

$$\frac{|K|}{\delta t} (\rho_K y_K - z_K) + \sum_{\sigma=K|L} G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K) = 0 \quad (\text{IV.48})$$

where  $g(\cdot, \cdot)$  corresponds to an approximation of  $\varphi(y) = \max[y(1-y), 0]$  by a monotone numerical flux function,  $G_{\sigma,K}^+ = \max(G_{\sigma,K}, 0)$ ,  $G_{\sigma,K}^- = -\min(G_{\sigma,K}, 0)$  and  $G_{\sigma,K}$  is given by the relation (IV.47). Then, if  $g(y_K, y_L) \geq 0$  for all  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , there exists a discretization for the term :

$$\left[ \frac{\alpha_g (1 - \alpha_g)}{\lambda \rho} (\rho_g - \rho_\ell) \right]_\sigma$$

in (IV.47) such that the following stability estimate holds :

$$\frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [f(\rho_K, \rho_K y_K) - f(\rho_K, z_K)] \leq 0$$

which means that the drift term is dissipative with respect to the entropy of the system.

**Proof.**

We multiply equation (IV.48) by the partial derivative of  $f(\cdot, \cdot)$  with respect to the second variable, taken at the point  $(\rho_K, \rho_K y_K)$ , and sum up over the control volumes of the mesh :

$$\sum_{K \in \mathcal{M}} h(\rho_K, \rho_K y_K) \left[ \frac{|K|}{\delta t} (\rho_K y_K - z_K) + \sum_{\sigma=K|L} G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K) \right] = T_1 + T_2 = 0$$

where  $T_1$  and  $T_2$  reads :

$$T_1 = \sum_{K \in \mathcal{M}} \frac{|K|}{\delta t} h(\rho_K, \rho_K y_K) [\rho_K y_K - z_K]$$

$$T_2 = \sum_{K \in \mathcal{M}} h(\rho_K, \rho_K y_K) \left[ \sum_{\sigma=K|L} G_{\sigma,K}^+ g(y_K, y_L) - G_{\sigma,K}^- g(y_L, y_K) \right]$$

As the fonction  $f(\cdot, \cdot)$  is convex, we have :

$$T_1 \geq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [f(\rho_K, \rho_K y_K) - f(\rho_K, z_K)] \quad (\text{IV.49})$$

Let us turn to  $T_2$ . Reordering the sum, we get :

$$T_2 = \sum_{\sigma \in \mathcal{E}_{\text{int}}} |\sigma| \rho_{\sigma, \text{up}} g_{\text{up}}(y_K, y_L, u_r) \left[ \frac{\alpha_g (1 - \alpha_g)}{\lambda \rho} (\rho_g - \rho_\ell) \right]_\sigma (p_K - p_L) [h(\rho_K, \rho_K y_K) - h(\rho_L, \rho_L y_L)]$$

where  $g_{\text{up}}(y_K, y_L, u_r) = g(y_K, y_L)$  if  $u_r \cdot n_{KL} \geq 0$  and  $g_{\text{up}}(y_K, y_L, u_r) = g(y_L, y_K)$  if  $u_r \leq 0$ ; in any case, we have, by assumption,  $g_{\text{up}}(y_K, y_L, u_r) \geq 0$ . We now choose, for the approximation of the quantity defined on  $\sigma$  in the preceding relation, an expression of the form :

$$\left[ \frac{\alpha_g (1 - \alpha_g)}{\lambda \rho} (\rho_g - \rho_\ell) \right]_\sigma \equiv \frac{(\alpha_g)_\sigma [1 - (\alpha_g)_\sigma]}{\lambda \rho_\sigma} (\rho_g(p_\sigma) - \rho_\ell)$$

where  $(\alpha_g)_\sigma$  and  $\rho_\sigma$  stand for approximations of the void fraction and the density on  $\sigma$ , respectively, which only need here to be supposed non-negative. Applying lemma IV.4.15,  $T_2$  reads :

$$T_2 = \sum_{\sigma \in \mathcal{E}_{\text{int}}} |\sigma| \rho_{\sigma, \text{up}} g_{\text{up}}(y_K, y_L, u_r) \frac{(\alpha_g)_\sigma (1 - \alpha_g)_\sigma}{\lambda \rho_\sigma} \rho_\ell \rho_g(p_\sigma) h'_p(p_\sigma) (p_K - p_L) [h_p(p_K) - h_p(p_L)]$$

If  $p_K = p_L$ , the term associated to  $K|L$  in this sum vanishes. Otherwise, from the third assertion of lemma IV.4.15, there exists  $p_\sigma \in [\min(p_K, p_L), \max(p_K, p_L)]$  such that the product  $h'_p(p_\sigma) (p_K - p_L) [h_p(p_K) - h_p(p_L)]$  is positive. Since we choose  $(\alpha_g)_\sigma$  such that  $(\alpha_g)_\sigma \geq 0$ , all the other quantities are positive, and this concludes the proof.  $\blacksquare$

The following proposition extends the stability result of the preceding section to the case  $u_r \neq 0$ .

**Proposition IV.4.17 (Stability of the scheme, case  $u_r \neq 0$ )**

*Let the density of the liquid phase be constant, the gas phase obey the ideal gas law and  $f(\rho, z)$  be the corresponding volumetric free energy of the mixture, defined by (IV.36). We suppose that the viscous term is dissipative (i.e.  $\forall v \in W_h, a_d(v, v) \geq 0$ ). In addition, we assume that the density  $\rho^n$  is positive and the gas mass fraction  $y^n$  belongs to the interval  $(0, 1]$ . Let  $\tilde{u}^{n+1}, u^{n+1}, p^{n+1}, z^{n+1}, \rho^{n+1}$  and  $y^{n+1}$  be a solution to the equations of one time step of the scheme, with a zero forcing term. We suppose that the drift velocity is given by a Darcy-like relation (IV.7) and that the discretization of the correction step for the gas mass fraction  $y^{n+1}$  is such that the stability result of proposition IV.4.16 applies. Then the following bound holds :*

$$\begin{aligned} \frac{1}{2} \|u^{n+1}\|_{h, \rho^n}^2 + \int_{\Omega} f(\rho^{n+1}, \rho^{n+1} y^{n+1}) dx + \delta t a_d(\tilde{u}^{n+1}, \tilde{u}^{n+1}) + \frac{\delta t^2}{2} |p^{n+1}|_{h, \rho^n}^2 \\ \leq \frac{1}{2} \|u^n\|_{h, \rho^{n-1}}^2 + \int_{\Omega} f(\rho^n, \rho^n y^n) dx + \frac{\delta t^2}{2} |p^n|_{h, \rho^n}^2 \end{aligned}$$

**Proof.**

Proposition IV.4.16 yields :

$$\frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| [f(\rho_K^{n+1}, \rho_K^{n+1} y_K^{n+1}) - f(\rho_K^n, z_K^n)] \leq 0$$

The conclusion thus follows by summing this relation with the estimate of proposition IV.4.13.  $\blacksquare$

Finally, note that, as in the preceding section, this partial stability result yields the same entropy decrease estimate for the whole scheme as in the preceding section if a renormalization step for the pressure is added to the scheme.

**Remark IV.4.18 (On the choice of the monotone numerical flux function)** *As stated in section IV.2, we have adopted for the numerical tests presented hereafter the following flux-splitting formula :*

$$g(a_1, a_2) = g_1(a_1) + g_2(a_2)$$

where  $g_1(a_1) = a_1$  if  $a_1 \in [0, 1]$  and zero otherwise and  $g_2(a_2) = -(a_2)^2$  if  $a_2 \in [0, 1]$  and zero otherwise. This numerical monotone flux does not satisfy the hypothesis of proposition IV.4.17, as it is not always non-negative. However, several other choices are possible for the numerical flux function  $g(\cdot, \cdot)$  (e.g.[26]), and some of them solve this problem. Thanks to the fact that  $\varphi(s) = s(1 - s)$  is positive  $\forall s \in [0, 1]$ , it is the case, for example, for the flux obtained with a one-dimensional Godunov scheme for each interface :

$$g(a_1, a_2) = \begin{cases} \max\{\varphi(s), a_2 \leq s \leq a_1\} & \text{if } a_2 \leq a_1 \\ \min\{\varphi(s), a_1 \leq s \leq a_2\} & \text{if } a_1 \leq a_2 \end{cases}$$

## IV.5 Numerical results

This section is devoted to numerical tests of the proposed scheme. We first address a problem built in such a way that it admits an analytical solution, to assess the convergence properties of the scheme. Then several additional tests are performed, to check the stability of the algorithm and the quality of the results.

### IV.5.1 Assessing the convergence against an analytic solution

We address here a problem built by the so-called technique of manufactured solutions : the computational domain and the solution are chosen *a priori* and the initial conditions, the boundary conditions and the forcing terms are adjusted consequently. Let thus the computational domain be  $\Omega = (0, 1) \times (-1/2, 1/2)$ , and the density and the momentum take the following expressions :

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

The pressure and the partial gas density are linked to the density by the equation of state (IV.11), where the liquid density  $\rho_\ell$  is set at  $\rho_\ell = 5$  and the quantity  $a^2$  in the equation of state of the gas (IV.5) is given by  $a^2 = 1$  (so  $\rho_g = p$ ). We choose the following expression for the unknowns  $y$  and  $z$  :

$$y = \frac{2.5 - 0.5 \rho}{4.5 \rho} \quad z = \rho y = \frac{2.5 - 0.5 \rho}{4.5}$$

The relative velocity is constant and given by  $u_r = (0, 1)^t$  and the diffusive coefficient  $D$  is set to  $D = 0.1$ . The analytical expression for the pressure is obtained from the equation of state (*i.e.* relation (IV.35)). These functions satisfy the mass balance equation ; for the gas mass fraction and momentum balance, we add the corresponding right-hand side. In this latter equation, we suppose that the divergence of the stress tensor is given by :

$$\nabla \cdot \tau(u) = \mu \Delta u + \frac{\mu}{3} \nabla \nabla \cdot u, \quad \mu = 10^{-2}$$

and we use for the viscous term the corresponding form for the bilinear form  $a_d(\cdot, \cdot)$  (see section IV.2.3).

Errors for the velocity, pressure and gas mass fraction obtained at  $t = 0.5$ , as a function of the time step and for various meshings, are drawn on figure IV.2, figure IV.3 and figure IV.4, 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. Computations are made with  $20 \times 20$ ,  $40 \times 40$  and  $80 \times 80$  uniform meshes (so with square cells and the Rannacher-Turek element). 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

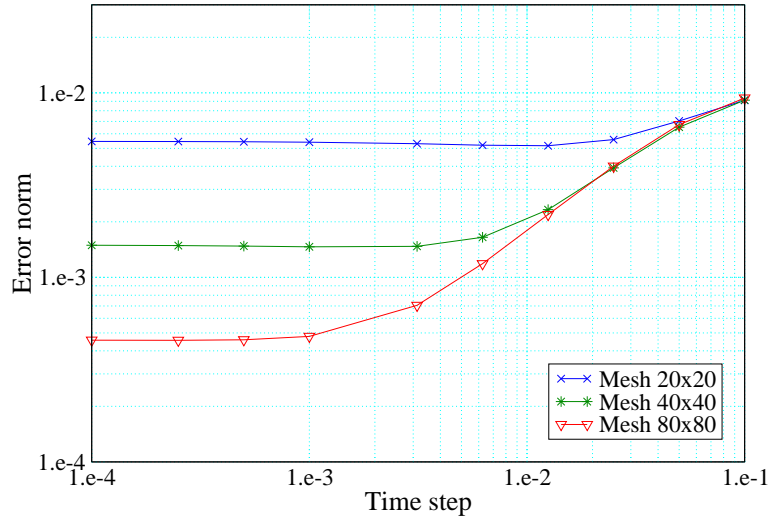


FIG. IV.2 – Error for the velocity at  $t = 0.5$ , as a function of the time step ( $L^2$  norm).

spatial discretization error. The value of the errors on this plateau then show a spatial convergence order close to one, which is consistent with the choice of an upwind discretization for the advection terms in the mass and gas mass fraction balance equations.

#### IV.5.2 The glass of water problem

We now consider as physical domain a square with unitary length edges. This domain is filled at the initial time with water in the bottom half ( $y = 0$  for  $x_2 < 0.5$ ) and with air in the top half ( $y = 1$  for  $x_2 > 0.5$ ). At the initial time, the fluid is at rest and the pressure is set to the uniform value  $p = 10^5 Pa$ . The density of the water is  $\rho_\ell = 1000 kg.m^{-3}$  and, for the air,  $\rho_g = p/a^2$ , where  $a^2$  is such that  $\rho_g = 1.2 kg.m^{-3}$  at  $p = 10^5 Pa$ . The diffusion coefficient  $D$  and the drift velocity are set to zero, and the viscosity is  $\mu = 10^{-2}$ .

Under the action of the gravity (with  $g = 9.81 m.s^{-2}$ ), the water is supposed to stay at rest, and the air moves up to obtain a final state where  $u = 0$  and the pressure obeys the ordinary differential equation  $dp/dx_2 = \rho_g(p)$ , its mean value being fixed by the conservation of the air mass and volume. The air thus compresses in the bottom of the air column and decompresses in the top, so it first moves downward; then, due to inertial effects, some oscillations occur, which are damped by the viscosity.

The first goal of this test is to check that the numerical scheme does not generate any spurious velocity field, due to what is usually called in the litterature "a bad balancing of the forcing term". In fact, this problem indeed appears for the finite element discretization used in this paper, and is cured here by using a numerical integration of the forcing term of the momentum balance (here the gravity) specially designed for this purpose. The guideline to build this discretization is to ensure that the discretization of a gradient is indeed a discrete gradient (*i.e.* if there exists a function  $\psi$  such that the forcing term  $f$  can be recast under the form  $f = \nabla\psi$ , the discrete right-hand side of the momentum balance belongs to the range of the discrete gradient  $B^t$ ). It has been extensively tested in the framework of the development at IRSN of the ISIS code, aimed at the description of low Mach number turbulent reactive flows as generated by fires, for the Rannacher-Turek element and rectangle or cuboid control volumes; its extension to Crouzeix-Raviart discretizations is underway. This test confirms the good behaviour of this discretization.



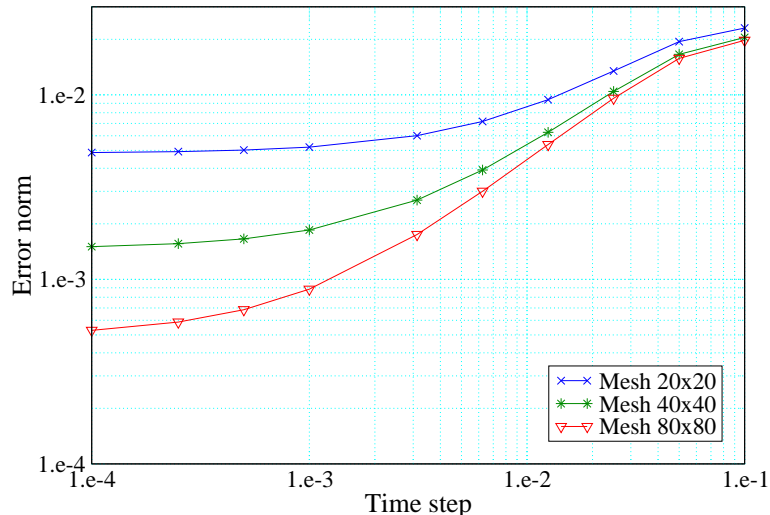


FIG. IV.3 – Error for the pressure at  $t = 0.5$ , as a function of the time step (discrete  $L^2$  norm).

Second, from our experience, an instability is likely to occur with fractional step algorithms, which should be termed "a sloshing instability". It seems to originate from the decoupling of the computation of the pressure and the free surface position. To explain its mechanism, let us imagine that the mesh is composed of only two vertical column of control volumes, which we refer to in the following as the left and right ones respectively. If the free surface in the left column of cells is higher than in the right one, the pressure, which, if we start from a situation where the fluid is approximately at rest, is almost equal to the hydrostatic one, is greater in the left column. This generates (for instance, through the velocity update in the pressure correction step) a velocity field from the left to the right. When updating the free surface (*i.e.*, here, the distribution of the gas mass fraction  $y$ ), this motion in turn induces a level increase in the right column of cell. If the free surface is now higher in this latter, the opposite motion occurs, and one can easily imagine that, in some cases, this phenomenon can develop up to lead to a blow up of the computation. This has been observed with the completely decoupled algorithm presented in [35] when the difference between the liquid and gas density is large, and could not be cured by a (reasonable) decrease of the time step. By coupling the computation of the pressure and the gas mass fraction, the present algorithm is likely to cure this problem : from the present test case and the sloshing transient presented hereafter, it seems to be indeed the case.

We document here a calculation which is performed with a  $40 \times 40$  uniform mesh (Rannacher-Turek element), using  $\delta t = 10^{-2} s$  as time step. The qualitative behaviour of the computed flow corresponds to what is expected ; in particular, a steady state with a fluid at rest is obtained after a short transient, where the pressure distribution is hydrostatic. The pressure along the line  $x_1 = 0.5$  is shown on figure IV.5.

### IV.5.3 Transport of interfaces

The transport of interfaces between phases is a standard test case for multiphase flow solvers. This test is presented here in the case of a one-dimensional, two-dimensional and three-dimensional geometry. In each case, the densities of the two phases, namely water and air, are  $\rho_\ell = 1000 kg.m^{-3}$  and  $\rho_g = p/a^2$ , where  $a^2$  is such that  $\rho_g = 1.2 kg.m^{-3}$  at  $p = 10^5 Pa$ . The pressure is initially constant, and takes the value  $p = 10^5 Pa$ . The relative velocity and the diffusive coefficient are set

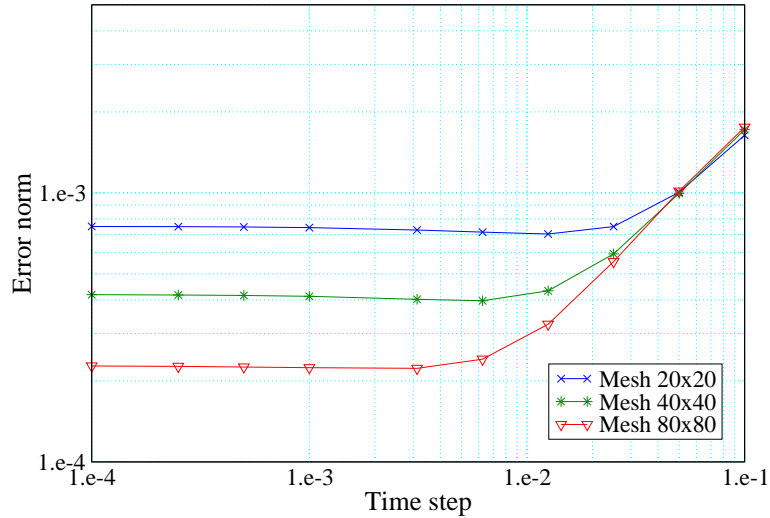


FIG. IV.4 – Error for the gas mass fraction at  $t = 0.5$ , as a function of the time step (discrete  $L^2$  norm).

to zero, while the viscosity is given by  $\mu = 1 Pa.s$ .

The one dimensional case consists in the advection of a material interface initially located at  $x = 0.5 m$ , at the mid-plane of a domain of length  $L = 1 m$ , with water on the left and air on the right. The initial velocity is set to  $u = 1 m.s^{-1}$ . We use a regular mesh composed of rectangular cells (with the Rannacher-Turek element); since this problem is one dimensional, only one grid cell is used in the direction perpendicular to the flow, and 100 cells in the parallel one. Results obtained with the time step  $dt = 1.2510^{-2} s$  (so with a CFL number equal to 1.25) at  $t = 0.2 s$  are displayed on figure IV.6.

The second test and the third one are respectively two-dimensional and three-dimensional cases, and consist in a skew-to-the mesh advection of a disc ( $d = 2$ ) or a bubble ( $d = 3$ ), initially of unit radius and located at the center of the domain  $\Omega = (0, 5 m)^d$ . The void fraction is set to one inside the disc or the bubble, and to zero elsewhere. Each component of the initial velocity is set to one. Computations are made with triangular and tetrahedral control volumes, obtained by cutting in two triangles each cell of a  $40 \times 40$  uniform mesh ( $d = 2$ ) and in six tetrahedra each cell of  $20 \times 20 \times 20$  uniform mesh ( $d = 3$ ), using the Crouzeix-Raviart elements. The time step is equal to  $dt = 10^{-1} s$ , so the CFL number is approximately equal to 0.5 ( $d = 2$ ) and 1 ( $d = 3$ ). The results are drawn on figure IV.7 ( $d = 2$ ) and on figure IV.8 ( $d = 3$ ).

In both cases, the velocity and the pressure remain constant up to the accuracy of the linear solver or the Newton algorithm, and the gas mass fraction is transported by this constant velocity. Specially in the multidimensional cases, the results show the diffusive behaviour of the numerical scheme, due to the upwind discretization. Less diffusive discretizations should be developed, using, for example, MUSCL-like techniques. Note however that the drift has for effect to sharpen an interface perpendicular to the drift velocity, which explains for instance the thinness of the mixed zone in the bubble column test presented hereafter.

#### IV.5.4 Two-dimensional sloshing in cavity

Two layers of non-miscible fluids (air and water) are superimposed with the lighter one on top of the heavier one. The gravity (with  $g = 9.81 m.s^{-2}$ ) is acting in the vertical downward

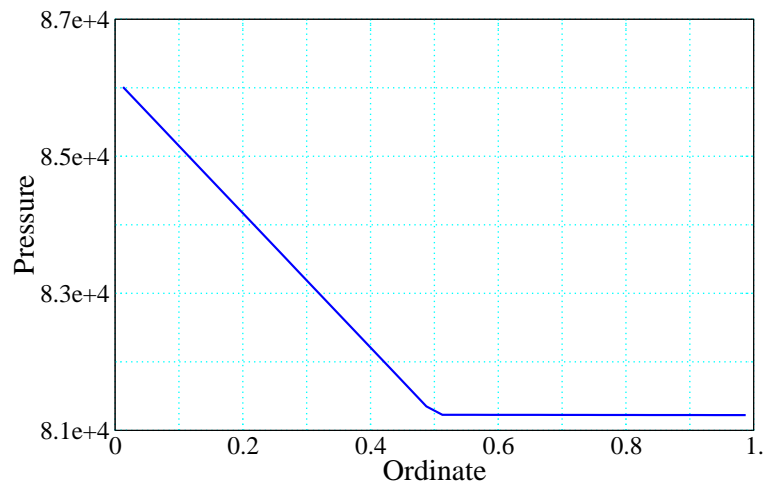


FIG. IV.5 – The glass of water problem : pressure at the end of the transient.

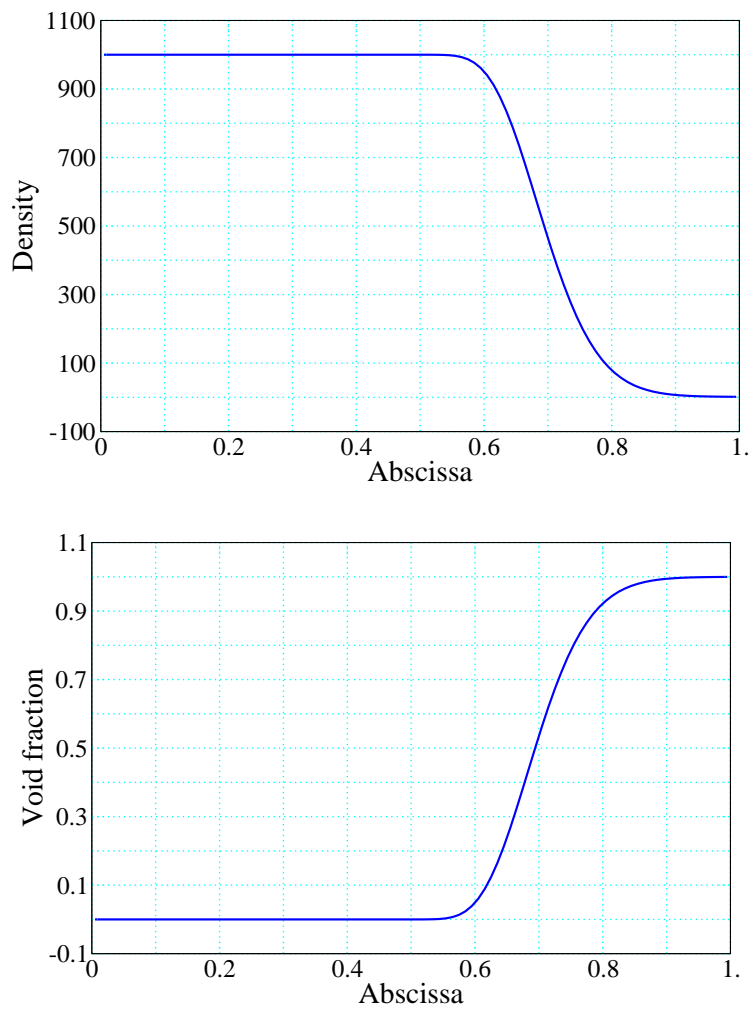


FIG. IV.6 – Transport of interfaces : density and void fraction at time 0.2 s.

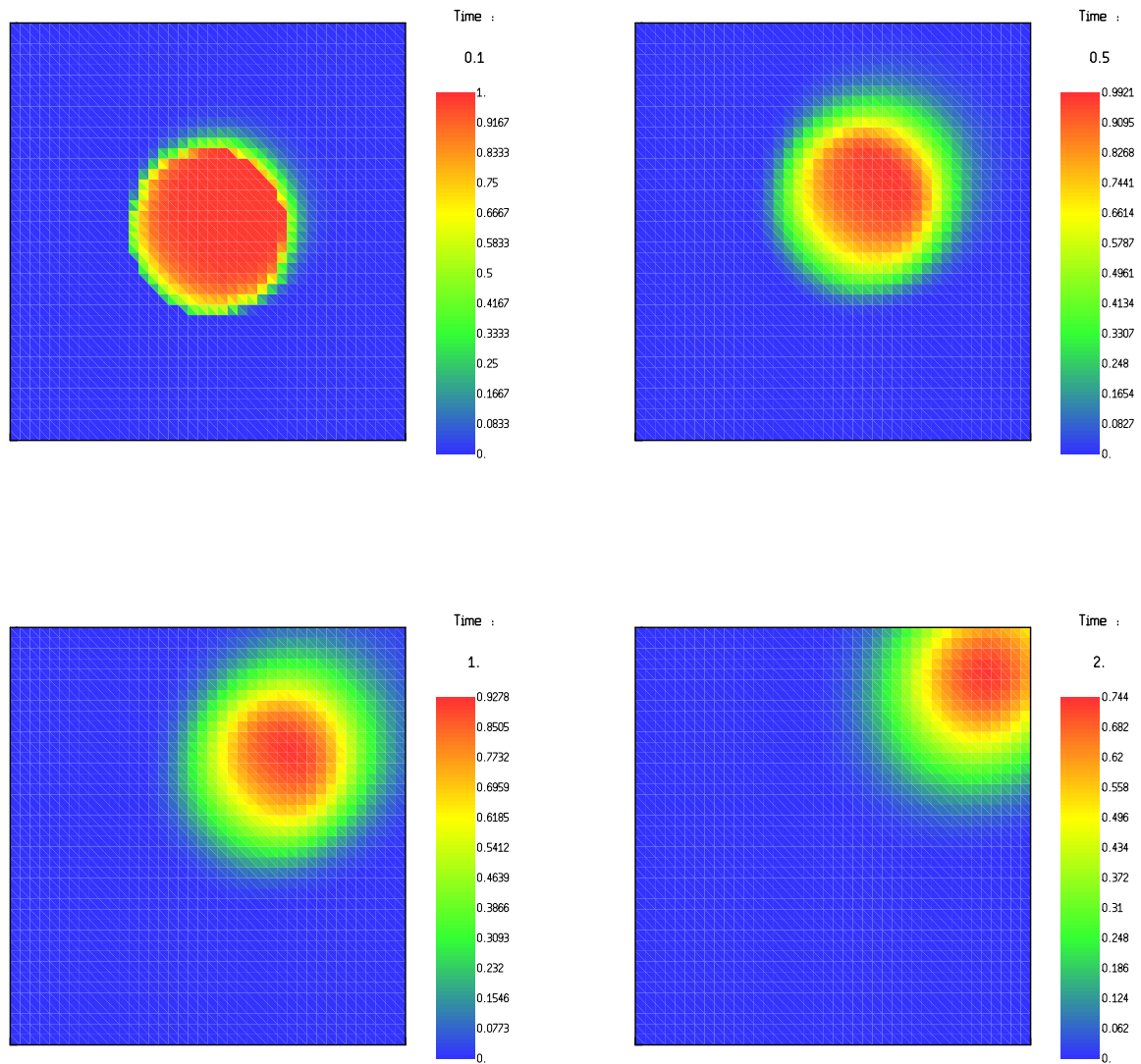


FIG. IV.7 – Transport of a circular bubble : void fraction at time 0.1 s, 0.5 s, 1 s and 2 s.

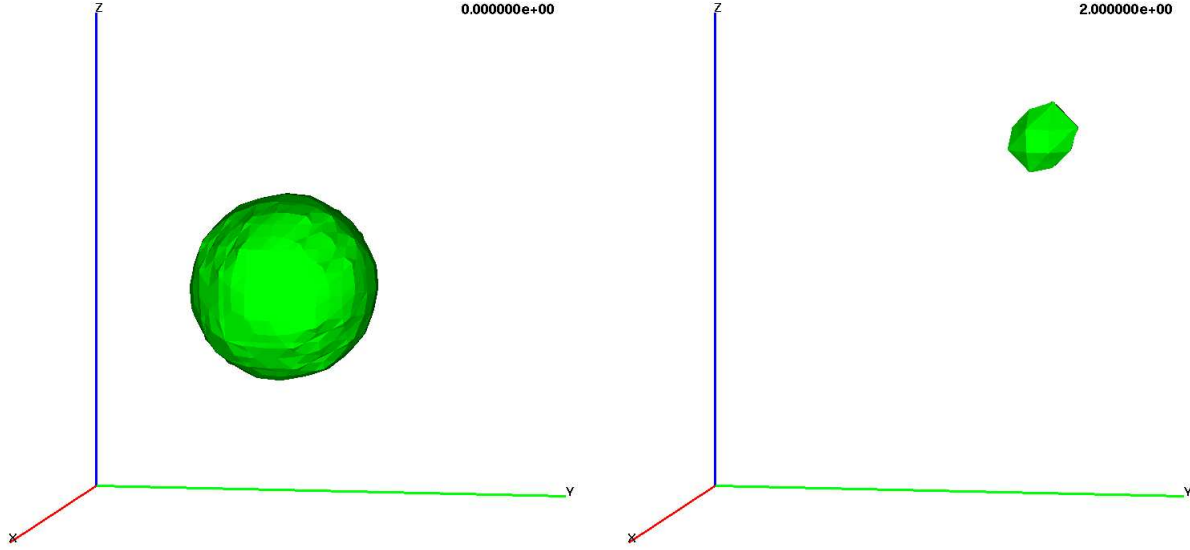


FIG. IV.8 – Transport of a spherical bubble : void fraction at time 0 s and 2 s (contour of the isosurface  $\alpha_g = 0.5$ ).

direction. The length of the rectangular cavity is  $L = 1\text{ m}$ , the height of each layer is respectively  $h_\ell = 1\text{ m}$  and  $h_g = 1.25\text{ m}$ , so the total height of the box is  $2.25\text{ m}$ . The water and air densities are respectively  $\rho_\ell = 1000\text{ kg.m}^{-3}$  and  $\rho_g = p/a^2$  where  $a^2$  is such that  $\rho_g = 1.2\text{ kg.m}^{-3}$  at  $p = 10^5\text{ Pa}$ . The diffusion coefficient  $D$  and the drift velocity are set to zero. A perfect slip condition is imposed on the whole boundary. At initial time, both fluids are at rest, then the cavity is submitted to an horizontal acceleration given by  $a_0 = 0.1\text{ m.s}^{-2}$ .

In the case where both fluids are supposed incompressible and the convection and diffusion terms may be neglected, an analytical solution for the flow in a rectangular cavity is provided in [12]. In particular, the shape of the interface is given by the following relation :

$$\xi = \frac{a_0}{g} \left[ x - \frac{L}{2} + \sum_{n \geq 0} \frac{4}{L k_{2n+1}^2} \cos(\omega_{2n+1} t) \cos(k_{2n+1} t) \right]$$

where the wave number  $k_n$  is defined by :

$$k_n = \frac{2\pi n}{L}$$

and  $\omega_n$  is given by :

$$\omega_n^2 = \frac{g k_n (\rho_\ell - \rho_g)}{\rho_g \coth(k_n h_g) + \rho_\ell \coth(k_n h_\ell)}$$

In practice, to compute this analytical solution, we perform the summation up to  $n = 200$ .

So as to remain in the domain of validity of the solution, the amplitude of the fluid oscillations must be very small ; hence a very fine mesh is necessary near the free surface, to capture its motion. The mesh is thus made of about 41 000 rectangular cells (with the Rannacher-Turek element) and, in the vertical direction, the space step is adapted in such a way that it is smaller near the interface between the two phases and equal to  $\delta x_2 = 0.0005\text{ m}$ , and increases when moving away the free

surface, up to  $\delta x_2 = 0.05 m$  at the top and bottom sections. In the horizontal direction, the mesh is uniform with step size  $\delta x_1 = 1/70 m$ . Calculations with different viscosities have been performed, these latter being supposed to vary with the mixture density :  $\mu = \rho/100$ ,  $\mu = \rho/1000$ ,  $\mu = \rho/10000$ .

The numerical results are reported on figure IV.9 ( $\mu = \rho/100$ ), figure IV.10 ( $\mu = \rho/1000$ ), and figure IV.11 ( $\mu = \rho/10000$ ) respectively. Comparing the obtained shape for the interface with the analytical solution, we observe that the numerical solution is closer to the analytical one with  $\mu = \rho/1000$  than with  $\mu = \rho/100$ , certainly because the fluid is too viscous in this latter case. More surprisingly, when reducing the viscosity to  $\mu = \rho/10000$ , the numerical solution also becomes less accurate. Our explanation is that, to obtain a good solution, it is necessary to respect a balance between approaching the physical problem (which, in this case, would suggest  $\mu = 0$ ) and keeping sufficient coercivity to ensure a reasonable convergence of the numerical approximation (which, on the contrary, requires a high value for the viscosity). With a more refined mesh, viscosity thus probably could be decreased, and the solution be closer to the analytical one. However, with this mesh already, results seem to be rather more accurate than those available in the literature [12].

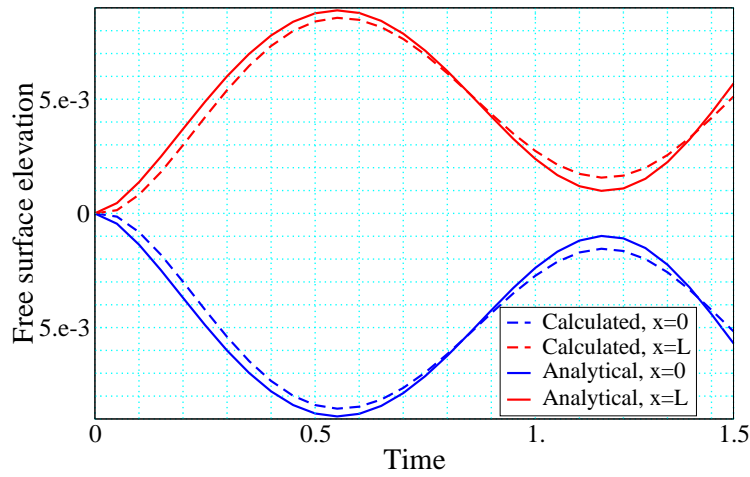


FIG. IV.9 – Sloshing in cavity : analytical solution and numerical solution with  $\mu = \rho/100$ .

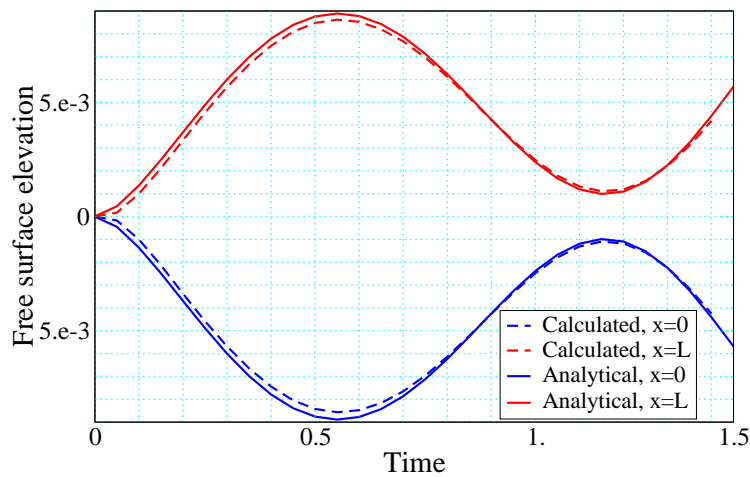


FIG. IV.10 – Sloshing in cavity : analytical solution and numerical solution with  $\mu = \rho/1000$ .

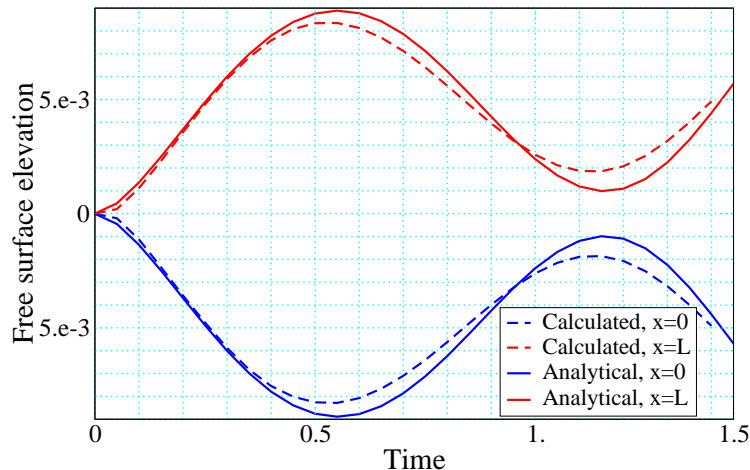


FIG. IV.11 – Sloshing in cavity : analytical solution and numerical solution with  $\mu = \rho/10000$ .

#### IV.5.5 Bubble column

We address in this section a classical benchmark for diphasic flow solvers, namely the flow in a pseudo two dimensional bubble column investigated experimentally by Becker *et al.*[5]. The apparatus has a rectangular cross section with the following dimensions : its width is  $L = 50\text{ cm}$ , its depth is  $8\text{ cm}$  and it is  $H = 200\text{ cm}$  high (see figure IV.12). It is filled with water up to the height  $h = 150\text{ cm}$ . A gas sparger, positioned  $15\text{ cm}$  from the left wall, is used to introduce an air flow of  $q = 8\text{ l.min}$  into the system. The circular sparger has a diameter of  $40\text{ mm}$  and a pore size of  $40\text{ }\mu\text{m}$ . Several liquid circulation cells can be observed in the column, the location and size of which continuously change. The bubble swarm is influenced by these vortices and therefore rises in a meander-like way. The direction of its lower part is stable and directed towards the nearest sidewall ; its upper part changes its shape and location in a quasiperiodic way, according to transient liquid circulations [62].

To simulate this experiment, we choose the following data. The boundary conditions are defined at the inlet as follows :

$$u_{imp} = \frac{q}{S \alpha_{g,imp}}$$

where  $S$  is the gas inlet area and  $\alpha_{g,imp} = 1$  is the void fraction imposed at the inlet. Along the walls and at the outlet of the column, homogeneous Dirichlet conditions are used for the velocity. Initial conditions are set to  $u = 0\text{ m.s}^{-1}$  and  $p = p_0$  where  $p_0 = 10^5\text{ Pa}$  is the ambient pressure. The density of the liquid is  $\rho_\ell = 1000\text{ kg.m}^{-3}$  ; the gas obeys an ideal gas equation of state  $\rho_g = p/a^2$ , where  $a^2$  is such that  $\rho_g = 1.2\text{ kg.m}^{-3}$  at  $p = 10^5\text{ Pa}$ . The diffusion coefficient  $D$  is set to zero, the drift velocity is constant and given by  $u_r = (0, 0.2)^t\text{ m.s}^{-1}$ .

For this test case, we use a regular meshing composed of rectangular cells (with the Rannacher-Turek element) with 76 cells in the horizontal direction, out of which 4 are for the gas inlet, and 300 in the vertical one. Calculations with time steps up to  $\delta t = 10^{-1}\text{ s}$  have been performed, observing that smaller time steps yield a thinner free surface.

The viscosity is a parameter which is difficult to adjust, since, in this simulation which is based on the system of equations governing a laminar flow, it must represent in some way the turbulent diffusion, *i.e.* the effects of fluctuations of the flow at microscopic scales, which may originate from the usual turbulence phenomena (sometimes termed "monophasic turbulence") and from the

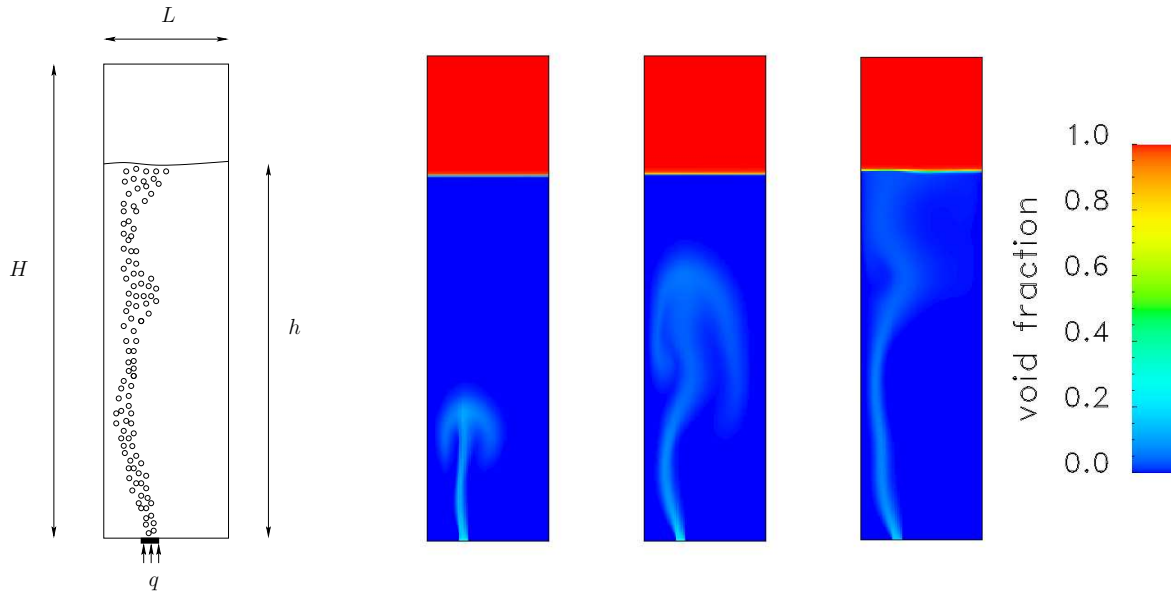


FIG. IV.12 – Bubble column : geometry of the problem and void fraction at times  $2s$ ,  $4s$  and  $40s$ .

perturbation of the velocity field due to the motion of the bubbles (sometimes termed "diphasic turbulence"). Calculations with a viscosity ranging from  $\mu = 10^{-3} Pa.s$  to  $\mu = 10^2 Pa.s$  have been performed. With smaller viscosities, we observe more oscillations of the free surface, the bubble swarm reaches the free surface faster and is farther from the sidewall.

Finally, the numerical results obtained with  $\delta t = 10^{-2} s$  and a viscosity of  $\mu = 1 Pa.s$  are reported on figure IV.12. With this value of the viscosity and these mesh and time steps, numerical convergence seems to be reached, at least visually. One can observe the stability and the thinness of the free surface. Results qualitatively reproduce the expected behaviour, which is the best we can hope with the rather crude modelling of turbulence which we adopted.

## IV.6 Conclusion

In this paper, we address the drift-flux model, which, for isothermal flows, consists in a system of three balance equations, namely the overall mass, the gas mass and the momentum balance complemented by an equation of state and a phenomenologic relation for the drift velocity.

For this problem, we develop a pressure correction scheme combining finite element and finite volume discretizations, which enjoys the following properties. First, the existence of a solution to each step of the algorithm is proven. Then, essential stability features of the continuous problem still hold at the discrete level : the unknowns are kept within their physical bounds (in particular, the gas mass fraction remains in the  $[0, 1]$  interval); in the homogeneous case (*i.e.* when the drift velocity vanishes), the discrete entropy of the system decreases; in addition, when using for the drift velocity the Darcy-like relation suggested in [39], the drift term becomes dissipative. Since, when the density is constant, this fractional step algorithm degenerates to an usual incremental projection method based on an *inf-sup* stable approximation, stability can be expected in the zero Mach number limit. Finally, the present algorithm preserves a constant pressure and a constant velocity through moving interfaces between phases (*i.e.* contact discontinuities of the underlying



hyperbolic system). To achieve this latter goal, the key ingredient is to couple the mass balance and the transport terms of the gas mass balance in an original pressure correction step.

We chose in this paper to only consider the case of a constant density liquid phase and of a gaseous phase obeying the ideal gas law. Dealing with a more general barotropic gas phase is certainly the simplest generalization, but the present theory also seems to extend to the case of a compressible fluid with minor modifications : for the stability study, essentially, the expression for the volumetric free energy of the mixture should be replaced by the usual expression applying when both phases are compressible, see for instance [39]; the existence theory would probably be simpler, since an upper bound for the density would provide in this case an estimate for the pressure. Returning to the case of an incompressible fluid, extending the present theory to deal with pure liquid zones appears on the contrary to be a difficult task, since the role played by the pressure in such a system seems to deserve some clarifications.

Numerical tests show a near-first-order convergence in space and time, consistent with the implemented discretization : first order backward Euler method in time and standard upwinding of the convection terms in the mass and gas mass fraction balance equations. With respect to this latter point, using more accurate space discretization (typically, MUSCL-like techniques) should certainly be desirable.

To assess the robustness of this algorithm, various numerical tests have been performed. They show in particular that free surface flows are computed without any instability, keeping a rather sharp interface throughout the computation. In addition, pure monophasic liquid zones are supported, although, as already mentioned, this case remains beyond the scope of the theory developed here. This scheme is now implemented in the ISIS code developed at IRSN and daily used for industrial applications.

## Appendix

### IV.7 Existence of a solution to a class of discrete diphasic problems

We address in this section the following abstract discrete problem :

$$\left\{ \begin{array}{l} a(u, \varphi_\sigma^{(i)}) - \int_{\Omega, h} p \nabla \cdot \varphi_\sigma^{(i)} dx = \int_{\Omega} f \cdot \varphi_\sigma^{(i)} dx, \quad \forall \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d \\ \frac{|K|}{\delta t} [\varrho^{p,z}(p_K, z_K) - \varrho^{p,z}(p_K^*, z_K^*)] \\ \quad + \sum_{\sigma=K|L} v_{\sigma,K}^+ \varrho^{p,z}(p_K, z_K) - v_{\sigma,K}^- \varrho^{p,z}(p_L, z_L) = 0, \quad \forall K \in \mathcal{M} \\ \frac{|K|}{\delta t} (z_K - z_K^*) + \sum_{\sigma=K|L} v_{\sigma,K}^+ z_K - v_{\sigma,K}^- z_L = 0, \quad \forall K \in \mathcal{M} \end{array} \right. \quad (\text{IV.50})$$

This problem is supposed to be obtained from (part of) a continuous problem by a space discretization combining Rannacher-Turek or Crouzeix-Raviart finite elements and finite volumes ; notations related to discrete quantities are given in section IV.2.2 and are not recalled here. The bilinear form  $a(\cdot, \cdot)$  is only assumed to be such that  $\|u\|_a = [a(u, u)]^{1/2}$  defines a norm over the discrete space  $W_h$ . The quantities  $(v_{\sigma,K}^+)^{n+1}$  and  $(v_{\sigma,K}^-)^{n+1}$  stands respectively for  $\max(v_{\sigma,K}^{n+1}, 0)$  and  $-\min(v_{\sigma,K}^{n+1}, 0)$  with  $v_{\sigma,K}^{n+1} = |\sigma| u_\sigma^{n+1} \cdot n_{KL}$ . Note that, in the last two equations, the flux summation excludes the external edges, which implicitly expresses the fact that the velocity is supposed to vanish on the boundary.

This system must be completed by three equations of state. The first two ones giving the liquid density  $\rho_\ell$  and the gas density  $\rho_g$  as a function of the pressure : we suppose here that the density of the liquid is constant and that the gas obeys the equation of state of ideal gases, which, for the sake of conciseness, we suppose here to be simply  $\rho_g = p$ . The last equation relates the mixture density  $\rho$  with the gas mass fraction  $y$  or the gas partial density  $z = \rho y$  and the phases density, and may take the three following forms :

$$\begin{aligned} p = \varphi(\rho, z) &= \frac{z \rho_\ell}{z + \rho_\ell - \rho} & \rho = \varrho^{p,z}(p, z) &= \rho_\ell + z \left(1 - \frac{\rho_\ell}{p}\right) \\ \rho = \varrho^{p,y}(p, y) &= \frac{1}{\frac{1-y}{\rho_\ell} + \frac{y}{p}} \end{aligned} \quad (\text{IV.51})$$

These three relations are equivalent as soon as the following assumptions for the unknowns of this system are satisfied :

$$\rho > 0, \quad p > 0, \quad z > 0 \quad \text{and} \quad 0 < y \leq 1 \quad (\text{IV.52})$$

These assumptions are natural, except for the hypothesis that  $y$  or  $z$  does not vanish, which excludes the existence of purely liquid zones. This latter assumption is assumed to hold for the initial quantities, *i.e.* we suppose that :

$$\forall K \in \mathcal{M}, \quad y_K^* = \frac{z_K^*}{\rho_K^*} \in (0, 1] \quad (\text{IV.53})$$

where  $\rho_K^* = \varrho^{p,z}(p_K^*, z_K^*)$ .

Our aim in this section is to prove that there exists a solution to system (IV.50) complemented with one of the relations of (IV.51), under the assumption (IV.53), and that any such solution satisfies the inequalities (IV.52).

We begin this section with two preliminary lemmas.

**Lemma IV.7.19**

Let  $(x_K^*)_{K \in \mathcal{M}}$  and  $(x_K)_{K \in \mathcal{M}}$  be two families of real numbers satisfying the following set of equations :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (x_K - x_K^*) + \sum_{\sigma=K|L} v_{\sigma,K}^+ x_K - v_{\sigma,K}^- x_L = 0$$

We suppose that,  $\forall K \in \mathcal{M}$ ,  $x_K^* > 0$ . Let  $\|\nabla_h \cdot u\|_\infty$  be defined by :

$$\|\nabla_h \cdot u\|_\infty = \max_{K \in \mathcal{M}} \left[ 0, \frac{1}{|K|} \sum_{\sigma=K|L} v_{\sigma,K} \right]$$

Then,  $\forall K \in \mathcal{M}$ ,  $x_K$  satisfies :

$$\frac{\min_{K \in \mathcal{M}} x_K^*}{1 + \delta t \|\nabla_h \cdot u\|_\infty} \leq x_K \leq \frac{1}{\min_{K \in \mathcal{M}} |K|} \sum_{K \in \mathcal{M}} |K| x_K^*$$

**Proof.**

The first inequality follows from an application of the discrete maximum principle lemma which can be found in [31] (lemma 2.5, section 2.3). The second one then follows from the fact that, by conservativity,  $\sum_{K \in \mathcal{M}} x_K = \sum_{K \in \mathcal{M}} x_K^*$ , remarking that, by the preceding relation, the values  $x_K$ , for  $K \in \mathcal{M}$ , are all positive. ■

The proof of the following result can be found in [51].

**Lemma IV.7.20**

Let  $(\rho_K^*)_{K \in \mathcal{M}}$ ,  $(x_K^*)_{K \in \mathcal{M}}$ ,  $(\rho_K)_{K \in \mathcal{M}}$  and  $(x_K)_{K \in \mathcal{M}}$  be four families of real numbers satisfying the following set of equations :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K x_K - \rho_K^* x_K^*) + \sum_{\sigma=K|L} v_{\sigma,K}^+ \rho_K x_K - v_{\sigma,K}^- \rho_L x_L = 0$$

We suppose that,  $\forall K \in \mathcal{M}$ ,  $\rho_K^* > 0$ ,  $\rho_K > 0$  and :

$$\forall K \in \mathcal{M}, \quad \frac{|K|}{\delta t} (\rho_K - \rho_K^*) + \sum_{\sigma=K|L} v_{\sigma,K}^+ \rho_K - v_{\sigma,K}^- \rho_L = 0$$

Then the following discrete maximum principle holds :

$$\forall K \in \mathcal{M}, \quad \min_{L \in \mathcal{M}} x_L^* \leq x_K \leq \max_{L \in \mathcal{M}} x_L^*$$

We now state the abstract theorem which will be used hereafter ; this result follows from standard arguments of the topological degree theory (see [19] for an overview of the theory and e.g. [25, 31] for other uses in the same objective as here, namely the proof of existence of a solution to a numerical scheme).

**Theorem IV.7.21 (A result from the topological degree theory)**

Let  $N$  and  $M$  be two positive integers and  $V$  be defined as follows :

$$V = \{(x, y, z) \in \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M \text{ such that } y > 0 \text{ and } z > 0\}$$

where, for any real number  $c$  and vector  $y$ , the notation  $y > c$  means that each component of  $y$  is greater than  $c$ . Let  $b \in \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M$  and  $f(\cdot)$  and  $F(\cdot, \cdot)$  be two continuous functions respectively from  $V$  and  $V \times [0, 1]$  to  $\mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M$  satisfying :

(i)  $F(\cdot, 1) = f(\cdot)$  ;

(ii)  $\forall \theta \in [0, 1]$ , if an element  $v$  of  $\bar{\mathcal{O}}$  (the closure of  $\mathcal{O}$ ) is such that  $F(v, \theta) = b$ , then  $v \in \mathcal{O}$ , where  $\mathcal{O}$  is defined as follows :

$$\mathcal{O} = \{(x, y, z) \in \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M \text{ s.t. } \|x\| < M \text{ and } \epsilon < y < M \text{ and } \epsilon < z < M\}$$

with  $M$  and  $\epsilon$  two positive constants and  $\|\cdot\|$  a norm defined over  $\mathbb{R}^N$  ;

(iii) the topological degree of  $F(\cdot, 0)$  with respect to  $b$  and  $\mathcal{O}$  is equal to  $d_0 \neq 0$ .

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

We are now in position to prove the existence of a solution to the considered discrete system.

**Theorem IV.7.22 (Existence of a solution)**

Under the assumption (IV.53), the nonlinear system (IV.50) complemented with the relation (IV.51) admits at least one solution, and any possible solution is such that :

$$\forall K \in \mathcal{M}, \quad \rho_K > 0, \quad z_K > 0, \quad 0 < y_K = \frac{z_K}{\rho_K} \leq 1, \quad p_K > 0$$

**Proof.**

This proof makes use of theorem IV.7.21 twice, by linking the initial problem (IV.50) to a linear one through two successive homotopies. Let  $N = d \operatorname{card}(\mathcal{E}_{\text{int}})$  and  $M = \operatorname{card}(\mathcal{M})$ ; we identify the finite element space of discrete velocity with  $\mathbb{R}^N$  and the finite volume space of pressure and partial density with  $\mathbb{R}^M$ . Let  $V$  be defined by  $V = \{(u, p, z) \in \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M \text{ such that } p > 0 \text{ and } z > 0\}$ .

Step 1 : first homotopy

We consider the function  $F : V \times [0, 1] \rightarrow \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M$  given by :

$$F(u, p, z, \theta) = \begin{cases} v_{\sigma,i}, & \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d \\ q_K, & K \in \mathcal{M} \\ s_K, & K \in \mathcal{M} \end{cases}$$

with :

$$\begin{cases} v_{\sigma,i} = a(u, \varphi_{\sigma}^{(i)}) - \int_{\Omega,h} p \nabla \cdot \varphi_{\sigma}^{(i)} dx - \int_{\Omega} f \cdot \varphi_{\sigma}^{(i)} dx \\ q_K = \frac{|K|}{\delta t} [\varrho_{\theta}^{p,z}(p_K, z_K) - \varrho_{\theta}^{p,z}(p_K^*, z_K^*)] \\ \quad + \sum_{\sigma=K|L} v_{\sigma,K}^+ \varrho_{\theta}^{p,z}(p_K, z_K) - v_{\sigma,K}^- \varrho_{\theta}^{p,z}(p_L, z_L), \\ s_K = \frac{|K|}{\delta t} [z_K - \varrho_{\theta}^{p,z}(p_K^*, z_K^*) y_K^*] + \sum_{\sigma=K|L} v_{\sigma,K}^+ z_K - v_{\sigma,K}^- z_L, \end{cases}$$

where the function  $\varrho_{\theta}^{p,z}(\cdot, \cdot)$  is implicitly defined by the following relation :

$$\varrho_{\theta}^{p,z}(p, z) = \varrho_{\theta}^{p,y}(p, y) = \frac{1}{\frac{1-y}{\varrho_{\ell,\theta}(p)} + \frac{y}{p}} \quad \text{with} \quad \varrho_{\ell,\theta}(p) = \frac{1}{\frac{\theta}{\rho_{\ell}} + \frac{1-\theta}{p}} \quad \text{and} \quad z = \rho y$$

Note that this definition makes sense (*i.e.* using  $z = \rho y$ , the function  $\varrho^{p,z}(\cdot, \cdot)$  can be explicitly computed from the expression of  $\varrho^{p,y}(\cdot, \cdot)$ ) as soon as  $p > 0$ , and thus for any  $(u, p, z) \in V$ .

Problem  $F(u, p, z, 1) = 0$  is exactly the same as system (IV.50).

Let  $\epsilon$  and  $M$  be two positive real numbers, and  $\mathcal{O}$  be defined by :

$$\mathcal{O} = \{(u, p, z) \in \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M \text{ s.t. } \|u\|_a < M, \epsilon < p < M \text{ and } \epsilon < z < M\}$$

We now suppose that  $(u, p, z) \in \bar{\mathcal{O}}$  (and thus, in particular,  $p \geq \epsilon$ ) and that  $F(u, p, z, \theta) = 0$  and provide estimates for  $(u, p, z)$ .

We begin by the following elementary bound, which is useful throughout the proof. From the definition of  $\varrho_{\ell,\theta}(p)$ , we observe that  $\min(\rho_{\ell}, p) \leq \varrho_{\ell,\theta}(p) \leq \max(\rho_{\ell}, p)$ . In the same way, provided

that  $y \in [0, 1]$ ,  $\min(\varrho_{\ell, \theta}(p), p) \leq \varrho_{\theta}^{p, z}(p, z) \leq \max(\varrho_{\ell, \theta}(p), p)$ . Hence,  $\min(\rho_{\ell}, p) \leq \varrho_{\theta}^{p, z}(p, z) \leq \max(\rho_{\ell}, p)$  and, thanks to assumption (IV.53) :

$$\forall \theta \in [0, 1], \forall K \in \mathcal{M}, \quad \varrho_{\theta}^{p, z}(p_K^*, z_K^*) \leq \bar{\rho}^* \quad \text{with} \quad \bar{\rho}^* = \max \left[ \left( \max_{K \in \mathcal{M}} p_K^* \right), \rho_{\ell} \right]$$

Step 1.1 :  $\|\cdot\|_a$  estimate for the velocity.

Let us first recast the equation of state of the mixture under a more convenient form. Substituting its definition for  $\varrho_{\ell, \theta}(p)$  in  $\varrho_{\theta}^{p, y}(p, y)$ , we get :

$$\rho = \frac{1}{\frac{\theta(1-y)}{\varrho_{\ell}} + \frac{y + (1-\theta)(1-y)}{p}} = \frac{1}{\frac{1-y'}{\varrho_{\ell}} + \frac{y'}{p}} \quad (\text{IV.54})$$

with  $y'(y, \theta) = y + (1-\theta)(1-y)$ . Then, taking  $y = z/\varrho_{\theta}^{p, z}(p, z)$  as unknown in the third equation of  $F(u, p, z, \theta) = 0$ , we get, for any  $K \in \mathcal{M}$  :

$$\frac{|K|}{\delta t} [\varrho_{\theta}^{p, z}(p_K, z_K) y_K - \varrho_{\theta}^{p, z}(p_K^*, z_K^*) y_K^*] + \sum_{\sigma=K|L} v_{\sigma, K}^+ \varrho_{\theta}^{p, z}(p_K, z_K) y_K - v_{\sigma, K}^- \varrho_{\theta}^{p, z}(p_L, z_L) y_L = 0$$

As, by the second equation of  $F(u, p, z, \theta) = 0$ , this relation vanishes for the constant function  $y_K = 1$ ,  $\forall K \in \mathcal{M}$ , we also obtain :

$$\begin{aligned} \frac{|K|}{\delta t} [\varrho_{\theta}^{p, z}(p_K, z_K) y'_K - \varrho_{\theta}^{p, z}(p_K^*, z_K^*) (y')_K^*] \\ + \sum_{\sigma=K|L} v_{\sigma, K}^+ \varrho_{\theta}^{p, z}(p_K, z_K) y'_K - v_{\sigma, K}^- \varrho_{\theta}^{p, z}(p_L, z_L) y'_L = 0 \end{aligned} \quad (\text{IV.55})$$

where, by assumption (IV.53),  $\forall K \in \mathcal{M}$ ,  $(y')_K^* = y_K^* + (1-\theta)(1-y_K^*) \in (1-\theta + \theta \underline{y}^*, 1]$ , with  $\underline{y}^* = \min_{K \in \mathcal{M}} y_K^*$ . We thus obtain a new problem, which keeps the structure of system (IV.50), with the same equation of state (*i.e.* relation (IV.54)) and just a modified initial value for  $z$  (*i.e.*  $z_K^*$  changed to  $\varrho_{\theta}^{p, z}(p_K^*, z_K^*) (y')_K^*$ ). The unknown  $p$  is still an unknown of this new problem, and we thus have by assumption  $p \geq 0$ . In addition, by lemma IV.7.20, any solution of this new problem is such that the gas mass fraction verifies  $1-\theta + \theta \underline{y}^* < y \leq 1$ , and thus the density and the gas partial density are positive. The unknowns thus belong to the domain where the free energy is correctly defined, and theorem IV.3.10 applies. Multiplying the first equation of  $F(u, p, z, \theta) = 0$  by  $u_{\sigma, i}$ , summing over  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $1 \leq i \leq d$  and using Young's inequality thus yields :

$$\frac{1}{2} \|u\|_a^2 + \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| \varrho_{\theta}^{p, z}(p_K, z_K) y'_K \log(p_K) \leq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| \varrho_{\theta}^{p, z}(p_K^*, z_K^*) (y')_K^* \log(p_K^*) + \frac{1}{2} \|f\|_{-a}^2$$

where  $\|\cdot\|_{-a}$  stands for the dual norm of  $\|\cdot\|_a$  with respect to the  $L^2$  inner product. The summation at the right hand side of this relation is bounded by  $(1/\delta t) |\Omega| \bar{\rho}^* \log(\bar{\rho}^*)$  where  $\bar{\rho}^* = \max_{K \in \mathcal{M}} p_K^*$ . By conservativity of equation (IV.55),  $\sum_{K \in \mathcal{M}} |K| \varrho_{\theta}^{p, z}(p_K^*, z_K^*) (y')_K^* \leq |\Omega| \bar{\rho}^*$ . Since, by assumption,  $p \geq \epsilon$ , we thus get :

$$\|u\|_a^2 \leq \frac{2}{\delta t} |\Omega| \bar{\rho}^* |\log(\epsilon)| + \frac{2}{\delta t} |\Omega| \bar{\rho}^* \log(\bar{\rho}^*) + \|f\|_{-a}^2$$

For  $\epsilon > 0$  small enough, we thus have :

$$\|u\|_a \leq c_1 |\log(\epsilon)|^{1/2} \quad (\text{IV.56})$$

where, in this relation and throughout the proof, we denote by  $c_i$  a real number only depending on the data of the problem, *i.e.*  $\Omega$ ,  $\bar{\rho}^*$ ,  $\bar{p}^*$ ,  $f$ ,  $a(\cdot, \cdot)$ ,  $\delta t$  and the mesh, and the expression "ε small enough" stands for  $\epsilon < c'_1$  where  $c'_1$  is a positive real number itself only depending on the data.

Step 1.2 :  $L^\infty$  estimates for  $z$ .

By equivalence of the norms over finite dimensional spaces, inequality (IV.56) also yields a bound for  $u$  in the  $L^\infty$  norm and, finally, for  $\|\nabla_h \cdot u\|_\infty$  :

$$\|\nabla_h \cdot u\|_\infty \leq c_2 |\log(\epsilon)|^{1/2}$$

By lemma IV.7.19, we thus get from the third relation of the system  $F(u, p, z, \theta) = 0$ , still for  $\epsilon$  small enough :

$$z \geq c_3 |\log(\epsilon)|^{-1/2} \quad (\text{IV.57})$$

On the other hand, we get from the same relation by conservativity :

$$z \leq c_4 \quad (\text{IV.58})$$

Step 1.3 :  $L^\infty$  estimates for  $p$ .

From the first relation of (IV.51), using the bounds for  $z$ , we get :

$$p \geq c_5 |\log(\epsilon)|^{-1/2} \quad (\text{IV.59})$$

To obtain an upper bound for  $p$ , we first remark that, as the considered spatial discretization satisfies a discrete *inf-sup* condition, a bound for  $u$  provides a bound for  $p - m(p)$  where  $m(p)$  stands for the mean value of  $p$ . By equivalence of norms on finite dimensional spaces, we can choose to express this bound in the seminorm defined by  $\forall q \in L_h$ ,  $|q|_{1,1,h} = \sum_{\sigma \in \mathcal{E}_{\text{int}}} (\sigma=K|L) |q_K - q_L|$ . With this semi-norm, the mean value of  $p$  disappears, and we get for  $\epsilon$  small enough :

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} (\sigma=K|L) |p_K - p_L| \leq c_6 |\log(\epsilon)|^{1/2} \quad (\text{IV.60})$$

An upper bound for  $p$  in one cell of the mesh, say  $K_0$ , would then provide an upper bound for  $p$ , since, for any  $K \in \mathcal{M}$ , it is possible to build a path from  $K_0$  to  $K$  crossing each internal edge at most once. To obtain such an estimate, we follow the following idea. If the pressure is somewhere lower than  $\rho_\ell$ , we are done; otherwise, with the chosen equation of state  $\rho_\theta^{p,z}(\cdot, \cdot)$ , when  $\theta$  varies, the liquid is everywhere denser than for  $\theta = 1$  and we are going to show that, even if its total mass also increases, the volume that it occupies is lower than for  $\theta = 1$ . Hence, the remaining volume for the gas is bounded away from zero, and, by conservation of the gas mass, the pressure cannot blow up everywhere. First, we need to introduce the phase volumetric fractions. The equation of state (IV.51) can be written as :

$$\frac{\rho}{p} + \frac{\rho - z}{\rho_\ell} = 1$$

and, as  $\rho - z = \rho(1 - y)$  and  $y \leq 1$ , both fractions at the left hand side of this relation are non-negative. We may thus define  $\alpha_g \in [0, 1]$  and  $\alpha_\ell \in [0, 1]$ , referred to as the gas and liquid volume fraction respectively, by :

$$\alpha_g = \frac{\rho}{p} \quad \alpha_\ell = \frac{\rho - z}{\rho_\ell}$$

Note that  $\alpha_g + \alpha_\ell = 1$ . Combining the second and the third relation of the system  $F(u, p, z, \theta) = 0$ , summing over the control volumes of the mesh and remarking that the fluxes cancel by conservativity, we get :

$$\sum_{K \in \mathcal{M}} |K| (\alpha_\ell)_K = \sum_{K \in \mathcal{M}} |K| \frac{(1 - y_K^*) \varrho_\theta^{p,z}(p_K^*, z_K^*)}{\rho_\ell} \quad (\text{IV.61})$$

Let us denote by  $(\alpha_\ell^*)_{K,1}$  the liquid void fraction with the equation of state corresponding to  $\theta = 1$  :

$$(\alpha_\ell^*)_{K,1} = \frac{1 - y_K^*}{\rho_\ell \left[ \frac{y_K^*}{p_K^*} + \frac{1 - y_K^*}{\rho_\ell} \right]}$$

Exploiting the form (IV.54) of the equation of state for  $\theta \neq 0$ , we obtain from relation (IV.61) :

$$\sum_{K \in \mathcal{M}} |K| (\alpha_\ell)_K = \sum_{K \in \mathcal{M}} |K| (\alpha_\ell^*)_{K,1} \frac{\frac{y_K^*}{p_K^*} + \frac{1 - y_K^*}{\rho_\ell}}{\frac{(y')_K^*}{p_K^*} + \frac{1 - (y')_K^*}{\rho_\ell}} \quad \text{with} \quad (y')_K^* = y_K^* + \theta(1 - y_K^*)$$

If we suppose that  $p_K \geq \rho_\ell$ , the fraction in the above equation is bounded by 1 : indeed, both the numerator and the denominator are harmonic averages of  $p_K^*$  and  $\rho_\ell$ , the weight associated to  $p_K^*$  being larger in the denominator, since  $(y')_K^*$  is closer to 1 than  $y_K^*$ . We thus get :

$$\sum_{K \in \mathcal{M}} |K| (\alpha_g)_K \geq c_7 = |\Omega| - \sum_{K \in \mathcal{M}} |K| (\alpha_\ell)_{K,1}$$

where  $c_7$  is positive by assumption, since  $\forall K \in \mathcal{M}, y_K^* > 0$ . Thus there exists  $K_0 \in \mathcal{M}$  such that  $(\alpha_g)_{K_0} \geq c_8 = c_7/|\Omega|$ . On the other hand, we have, still by conservativity :

$$\sum_{K \in \mathcal{M}} |K| (\alpha_g)_K p_K = \sum_{K \in \mathcal{M}} |K| z_K = \sum_{K \in \mathcal{M}} |K| z_K^* = \sum_{K \in \mathcal{M}} |K| \varrho_\theta^{p,z}(p_K^*, z_K^*) y_K^* \leq |\Omega| \bar{\rho}^*$$

We thus get, since all the  $(\alpha_g)_K$  and  $p_k$  are non-negative :

$$(\alpha_g)_{K_0} p_{K_0} \leq \frac{|\Omega|}{|K_0|} \bar{\rho}^*$$

and thus, as  $(\alpha_g)_{K_0}$  is bounded by below, the pressure is bounded by a quantity only depending on the data. As a consequence, for  $\epsilon$  small enough :

$$p \leq c_9 |\log(\epsilon)|^{1/2} \quad (\text{IV.62})$$



### Step 2 : second homotopy

We consider the function  $F : V \times [0, 1] \rightarrow \mathbb{R}^N \times \mathbb{R}^M \times \mathbb{R}^M$  given by :

$$F(u, p, z, \theta) = \begin{cases} v_{\sigma,i} = a(u, \varphi_{\sigma}^{(i)}) - \theta \int_{\Omega,h} p \nabla \cdot \varphi_{\sigma}^{(i)} dx - \int_{\Omega} f \cdot \varphi_{\sigma}^{(i)} dx, & \sigma \in \mathcal{E}_{\text{int}}, 1 \leq i \leq d \\ q_K = \frac{|K|}{\delta t} (p_K - p_K^*) + \theta \sum_{\sigma=K|L} v_{\sigma,K}^+ p_K - v_{\sigma,K}^- p_L, & K \in \mathcal{M} \\ s_K = \frac{|K|}{\delta t} (z_K - \frac{p_K^*}{\rho_K^*} z_K^*) + \theta \sum_{\sigma=K|L} v_{\sigma,K}^+ z_K - v_{\sigma,K}^- z_L, & K \in \mathcal{M} \end{cases}$$

The system  $F(u, p, z, 1) = 0$  is the same as the system obtained at the end of the preceding homotopy for  $\theta = 0$ , and the system  $F(u, p, z, 0) = 0$  is linear and clearly regular (by stability of the bilinear form  $a(\cdot, \cdot)$ ).

In addition, the third equation is now decoupled from the first two ones, and these latter have the structure of a monophasic compressible problem as studied in [31]. From this theory, an estimate similar to the first one in the preceding step is available and reads :

$$\frac{1}{2} \|u\|_a^2 + \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| p_K \log(p_K) \leq \frac{1}{\delta t} \sum_{K \in \mathcal{M}} |K| p_K^* \log(p_K^*) + \frac{1}{2} \|f\|_{-a}^2$$

Since the function  $s \mapsto s \log(s)$  is bounded by below on  $(0, +\infty)$ , this latter relation yields :

$$\|u\|_a \leq c_{10} \quad (\text{IV.63})$$

By lemma IV.7.19, we thus directly get :

$$c_{11} \leq p \leq c_{12}, \quad c_{13} \leq z \leq c_{14} \quad (\text{IV.64})$$

### Conclusion

We choose  $\epsilon$  small enough for the relations (IV.56), (IV.57), (IV.60) and (IV.62) to hold,  $\epsilon < \min(c_{11}, c_{13})$  and, in addition :

$$\epsilon < \max(c_3, c_5) |\log(\epsilon)|^{-1/2}$$

which is possible because the function  $s \mapsto s \log s$  tends to zero when  $s$  tends to zero. Let now  $M$  be such that :

$$M > \max \left[ \max(c_1, c_9) |\log(\epsilon)|^{1/2}, c_4, c_{10}, c_{12}, c_{14} \right]$$

Then, from inequalities (IV.56), (IV.57), (IV.58), (IV.59), (IV.62), (IV.63) and (IV.64), we get that throughout both homotopies, the unknown  $(u, p, z)$  remains in  $\mathcal{O}$ . As the last linear system is regular and admits a solution in  $\mathcal{O}$ , the topological degree of  $F(\cdot, \cdot, \cdot, \theta)$  with respect to  $\mathcal{O}$  and zero remains different of zero all along both homotopies, which proves the existence of a solution in  $\mathcal{O}$ .

We now turn to the proof of the *a priori* estimates  $\rho > 0$ ,  $z > 0$ ,  $0 < y^* \leq 1$  and  $p > 0$ . The fact that, if  $\rho^* > 0$  and  $z^* > 0$ , then  $\rho > 0$  and  $z > 0$  is a direct consequence of lemma IV.7.19 applied to the second and third relation of problem (IV.50). In addition, as both  $\rho^* > 0$  and  $\rho > 0$ , lemma IV.7.20 applies and thus, as  $0 < y^* \leq 1$ , we have  $0 < y \leq 1$ . If  $p \geq \rho_{\ell}$ , the fact that  $p > 0$  is evident. In the other case, by the equation of state written as a function of  $p$  and  $y$  (third form of (IV.51)), we get first that :

$$p \leq \rho < \rho_{\ell} \quad (\text{IV.65})$$

and, second, that, since  $\rho > 0$ , the pressure does not vanish. Thus the second form of this same relation (IV.51) can be written :

$$\rho = \frac{z}{p} p + \left(1 - \frac{z}{p}\right) \rho_\ell = \alpha_g p + (1 - \alpha_g) \rho_\ell$$

As  $p \neq \rho_\ell$ , the void fraction  $\alpha_g$  thus reads :

$$\alpha_g = \frac{\rho - \rho_\ell}{\rho_\ell - p}$$

which, by inequalities (IV.65), yields  $\alpha_g > 0$  and, finally, since  $z > 0$ ,  $p > 0$ . ■

This existence result applies directly to the pressure correction step used in the algorithm presented in this paper, with a particular expression for the bilinear form  $a(\cdot, \cdot)$ , which reads, dropping for short the time exponents :

$$a(u, v) = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma|}{\delta t} \rho_\sigma u_\sigma \cdot v_\sigma$$

Note that the analysis is performed here with a very simple equation of state for the gas ( $p = \rho$ ), but would be readily extended to general barotropic laws  $p = \wp(\rho)$ , under the mild assumptions that the corresponding free energy exists and is convex and the function  $\wp(\cdot)$  is increasing and one to one from  $(0, +\infty)$  to  $(0, +\infty)$ .

Let us now turn to the discretization of a stationary diphasic problem. As happens in the monophasic case, [34], it is likely that, in the case where the velocity is prescribed on the whole boundary, this problem needs to be completely determined the data of the total mixture mass (say  $M_m$ ) and of the total gas mass (say  $M_g$ ) present in the computational domain. A natural way to impose these two conditions is to add to this problem two regularizing terms in the mass balance and the gas mass balance :

$$\left\{ \begin{array}{l} c(h) |K| \left[ \varrho^{p,z}(p_K, z_K) - \frac{M_m}{|\Omega|} \right] + \sum_{\sigma=K|L} v_{\sigma,K}^+ \varrho^{p,z}(p_K, z_K) - v_{\sigma,K}^- \varrho^{p,z}(p_L, z_L) = 0 \quad \forall K \in \mathcal{M} \\ c(h) |K| \left[ z_K - \frac{M_g}{|\Omega|} \right] + \sum_{\sigma=K|L} v_{\sigma,K}^+ z_K - v_{\sigma,K}^- z_L = 0 \quad \forall K \in \mathcal{M} \end{array} \right.$$

where  $c(h)$  is a regularization parameter tending to zero with the size of the mesh. In this case, the present existence theory directly applies, provided that the momentum balance equation remains linear with respect to the velocity. Of course, under the same restriction, this is true also for an implicit discretization of a time-dependent problem.

In view of the stability results provided for the advection operator, adding such a term to the first relation of the problem (*i.e.* the momentum balance) should lead to a rather straightforward extension of the present existence result ; the advection term would be multiplied by the homotopy parameter and the stability (*i.e.* an analogue to estimate (IV.56)) would stem from the diffusion term. Note that, in this case, to keep the stability of the advection term, a regularization term consistent with the mass balance one should also be introduced in the momentum balance equation.

We have shown in this paper that, with Darcy's law for the drift velocity and a particular discretization for this term, the drift term is dissipative. So this term does not prevent obtaining stability estimates as (IV.56) ; this suggests that the existence theory developped here may perhaps be extended to the complete drift flux model.

Finally, we have not dealt in this study with the case where liquid monophasic zones ( $z = 0$ ) exist in the flow. In such zones, the pressure changes of mathematical nature : it is no more a parameter entering the equation of state and determined by the local density, but a Lagrange multiplier for the incompressibility constraint. Note that this fact is already underlying in the present study : indeed, the incompressibility of the liquid prevents to derive  $L^\infty$  estimates for the pressure from  $L^\infty$  estimates for the density (which are readily obtained using a conservation argument), and we must invoke to this purpose the stability of the discrete gradient (*i.e.* the discrete *inf-sup* condition), that is typically the argument allowing to control the pressure in incompressible flow problems. However, obtaining *a priori* estimates when  $z$  may vanish in the flow seems a difficult task, which should deserve more efforts. On the contrary, obtaining existence results for two barotropic phases seems to be rather simpler than the analysis performed here.



# Bibliographie

- [1] Ph. Angot, V. Dolejší, M. Feistauer, and J. Felcman. Analysis of a combined barycentric finite volume-nonconforming finite element method for nonlinear convection-diffusion problems. *Applications of Mathematics*, 4 :263–310, 1998.
- [2] 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.
- [3] 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*.
- [4] G.K. Batchelor and J.T. Green. The determination of the bulk-stress in a suspension of spherical particles to order  $c^2$ . *Journal of Fluid Mechanics*, 56 :401–427, 1972.
- [5] S. Becker, A. Sokolichin, and G. Eigenberger. Gas-liquid flow in bubble columns and loop reactors : Part II. comparison of detailed experiments and flow simulations. *Chemical Engineering Science*, 49(24B) :5747–5762, 1994.
- [6] M. Bern, D. Eppstein, and J. Gilbert. Provably good mesh generation. *Journal of Computer and System Sciences*, 48 :384–409, 1994.
- [7] H. Bijl and P. Wesseling. A unified method for computing incompressible and compressible flows in boundary-fitted coordinates. *Journal of Computational Physics*, 141 :153–173, 1998.
- [8] S. Brenner. Korn’s inequalities for piecewise  $H^1$  vector fields. *Mathematics of Computation*, 73(247) :1067–1087, 2003.
- [9] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, 1991.
- [10] M.O. Bristeau, R. Glowinski, L. Dutto, J. Périaux, and G. Rogé. Compressible viscous flow calculations using compatible finite element approximations. *International Journal for Numerical Methods in Fluids*, 11 :719–749, 1990.
- [11] V. Casulli and D. Greenspan. Pressure method for the numerical solution of transient, compressible fluid flows. *International Journal for Numerical Methods in Fluids*, 4 :1001–1012, 1984.
- [12] G. Chanteperdrix. *Modélisation et simulation numérique d’écoulements diphasiques à interface libre. Application à l’étude des mouvements de liquides dans les réservoirs de véhicules spatiaux*. Energétique et dynamique des fluides, Ecole Nationale Supérieure de l’Aéronautique et de l’Espace, 2004.

- 
- [13] G. Chen, C. Levermore, and T. Liu. Hyperbolic conservation laws with stiff relaxation terms and entropy. *Communications on Pure and Applied Mathematics*, XLVII :787–830, 1994.
- [14] A.J. Chorin. Numerical solution of the Navier-Stokes equations. *Mathematics of Computation*, 22 :745–762, 1968.
- [15] P. G. Ciarlet. Handbook of numerical analysis volume II : Finite elements methods – Basic error estimates for elliptic problems. In P. Ciarlet and J.L. Lions, editors, *Handbook of Numerical Analysis, Volume II*, pages 17–351. North Holland, 1991.
- [16] P. Colella and K. Pao. A projection method for low speed flows. *Journal of Computational Physics*, 149 :245–269, 1999.
- [17] F. Coquel, K. El Amine, E. Godlewski, B. Perthame, and P. Rasle. A numerical method using upwind schemes for the resolution of two-phase flows. *Journal of Computational Physics*, 136 :272–288, 1997.
- [18] 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.
- [19] K. Deimling. *Nonlinear Functional Analysis*. Springer, New-York, 1980.
- [20] I. Demirdžić, Ž. Lilek, and M. Perić. A collocated finite volume method for predicting flows at all speed. *International Journal for Numerical Methods in Fluids*, 16 :1029–1050, 1993.
- [21] V. Dolejší, M. Feistauer, J. Felcman, and A. Kliková. Error estimates for barycentric finite volumes combined with nonconforming finite elements applied to nonlinear convection-diffusion problems. *Applications of Mathematics*, 47 :301–340, 2002.
- [22] P. Drábek and J. Milota. *Methods of nonlinear analysis*. Birkhäuser Advanced Texts. 2007.
- [23] Donald A. Drew and Stephen L. Passman. *Theory of Multicomponent Fluids*, volume 135 of *Applied Mathematical Sciences*. Springer, 1999.
- [24] A. Ern. *Aide Mémoire Éléments finis*. Dunod, Paris, 2005.
- [25] 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.
- [26] 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.
- [27] R. Eymard and R. Herbin. Entropy estimate for the approximation of the compressible barotropic Navier-Stokes equations using a collocated finite volume scheme. *in preparation*, 2007.
- [28] E. Feireisl. Dynamics of viscous compressible flows. volume 26 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, 2004.
- [29] H. Feistauer, J. Felcman, and I. Straškraba. *Mathematical and computational methods for compressible flows*. Oxford Science Publications. Clarendon Press, 2003.

- [30] M. Fortin, H. Manouzi, and A. Soulaïmani. On finite element approximation and stabilization methods for compressible viscous flows. *International Journal for Numerical Methods in Fluids*, 17 :477–499, 1993.
- [31] T. Gallouët, L. Gastaldo, R. Herbin, and J.C. Latché. An unconditionally stable pressure correction scheme for compressible barotropic Navier-Stokes equations. *submitted to Mathematical Modelling and Numerical Analysis*.
- [32] T. Gallouët, J.-M. Hérard, and N. Seguin. On the use of symmetrizing variables for vacuums. *Calcolo*, 40 :163–194, 2003.
- [33] 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.
- [34] T. Gallouët, R. Herbin, and J.C. Latché. A convergent finite-element/finite-volume scheme for the compressible Stokes problem – part I : the isothermal case. *submitted to Mathematics of Computation*.
- [35] L. Gastaldo, R. Herbin, and J.C. Latché. On a discretization of phases mass balance in segregated algorithms for the drift-flux model. *submitted to IMA Journal of Numerical Analysis*.
- [36] D. Gilbarg and N. Trudinger. *Elliptic Partial Differential Equations of Second Order*. Classics in Mathematics. Springer, third edition, 2001.
- [37] J.-L. Guermond and L. Quartapelle. A projection FEM for variable density incompressible flows. *Journal of Computational Physics*, 165 :167–188, 2000.
- [38] 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.
- [39] H. Guillard and F. Duval. A Darcy law for the drift velocity in a two-phase flow model. *Journal of Computational Physics*, 224 :288–313, 2007.
- [40] H. Guillard and A. Murrone. On the behavior of upwind schemes in the low Mach number limit : II. Godunov type schemes. *Computer & Fluids*, 33(4) :655–675, may 2004.
- [41] F.H. Harlow and A.A. Amsden. Numerical calculation of almost incompressible flow. *Journal of Computational Physics*, 3 :80–93, 1968.
- [42] F.H. Harlow and A.A. Amsden. A numerical fluid dynamics calculation method for all flow speeds. *Journal of Computational Physics*, 8 :197–213, 1971.
- [43] M. Ishii. *Thermo-Fluid Dynamic Theory of Two-Phase Flow*. Collection de la Direction des Etudes et Recherches d’Electricité de France. Eyrolles, Paris, 1975.
- [44] R.I. Issa. Solution of the implicitly discretised fluid flow equations by operator splitting. *Journal of Computational Physics*, 62 :40–65, 1985.
- [45] R.I. Issa, A.D. Gosman, and A.P. Watkins. The computation of compressible and incompressible recirculating flows by a non-iterative implicit scheme. *Journal of Computational Physics*, 62 :66–82, 1986.
- [46] R.I. Issa and M.H. Javareshkian. Pressure-based compressible calculation method utilizing total variation diminishing schemes. *AIAA Journal*, 36 :1652–1657, 1998.

- 
- [47] J.B. Joshi, V.S. Vitankar, A.A. Kulkarni, M.T. Dhotre, and K. Ekambara. Coherent flow structures in bubble column reactors. *Chemical Engineering Science*, 57 :3157–3183, 2002.
- [48] K.C. Karki and S.V. Patankar. Pressure based calculation procedure for viscous flows at all speeds in arbitrary configurations. *AIAA Journal*, 27 :1167–1174, 1989.
- [49] M.H. Kobayashi and J.C.F. Pereira. Characteristic-based pressure correction at all speeds. *AIAA Journal*, 34 :272–280, 1996.
- [50] D. Kuzmin and S. Turek. Numerical simulation of turbulent bubbly flows. In *3rd International Symposium on Two-Phase Flow Modelling and Experimentation, Pisa, 22-24 September*, 2004.
- [51] B. Larrouturou. How to preserve the mass fractions positivity when computing compressible multi-component flows. *Journal of Computational Physics*, 95 :59–84, 1991.
- [52] P.-L. Lions. Mathematical topics in fluid mechanics – volume 2 – compressible models. volume 10 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, 1998.
- [53] 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.
- [54] F. Moukalled and M. Darwish. A high-resolution pressure-based algorithm for fluid flow at all speeds. *Journal of Computational Physics*, 168 :101–133, 2001.
- [55] F. Moukalled, M. Darwish, and B. Sekar. A pressure-based algorithm for multi-phase flow at all speeds. *Journal of Computational Physics*, 190 :550–571, 2003.
- [56] P. Nithiarasu, R. Codina, and O.C. Zienkiewicz. The Characteristic-Based Split (CBS) scheme – a unified approach to fluid dynamics. *International Journal for Numerical Methods in Engineering*, 66 :1514–1546, 2006.
- [57] A. Novotný and I. Straškraba. Introduction to the mathematical theory of compressible flow. volume 27 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, 2004.
- [58] F. Panescu. *Modélisation Eulerienne d'écoulements diphasiques à phase dispersée et simulation numérique par une méthode volumes finis - éléments finis*. Mathématiques appliquées, Université de Nice Sophia Antipolis, 2006.
- [59] G. Patnaik, R.H. Guirguis, J.P. Boris, and E.S. Oran. A barely implicit correction for flux-corrected transport. *Journal of Computational Physics*, 71 :1–20, 1987.
- [60] E.S. Politis and K.C. Giannakoglou. A pressure-based algorithm for high-speed turbomachinery flows. *International Journal for Numerical Methods in Fluids*, 25 :63–80, 1997.
- [61] R. Rannacher and S. Turek. Simple nonconforming quadrilateral Stokes element. *Numerical Methods for Partial Differential Equations*, 8 :97–111, 1992.
- [62] A. Sokolichin and G. Eigenberger. Applicability of the standard k- $\epsilon$  turbulence model to the dynamic simulation of bubble columns : Part I. detailed numerical simulations. *Chemical Engineering Science*, 54(13-14) :2273–2284, 1999.
- [63] A. Sokolichin, G. Eigenberger, and A. Lapin. Simulation of buoyancy driven bubbly flow : Established simplifications and open questions. *AIChE Journal*, 50(1) :24–45, 2004.



- [64] A. Sokolichin, G. Eigenberger, A. Lapin, and A. Lubbert. Dynamic Numerical Simulation of Gas-Liquid Two-Phase Flows Euler/Euler versus Euler/Lagrange. *Chemical Engineering Science*, 52(4) :611–627, 1997.
- [65] B. Spalding. Numerical computation of multiphase flow and heat transfer. In *Recent Advances in Numerical Methods in Fluids – Volume 1*, pages 139–168, 1980.
- [66] H. Staedtke, G. Franchello, B. Worth, U. Graf, P. Romstedt, A. Kumbaro, J. García-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.
- [67] H. Bruce Stewart and Burton Wendroff. Two-phase flow : Models and methods. *Journal of Computational Physics*, 56 :363–409, 1984. Review Article.
- [68] R. Temam. Sur l’approximation de la solution des équations de Navier-Stokes par la méthode des pas fractionnaires II. *Arch. Rat. Mech. Anal.*, 33 :377–385, 1969.
- [69] D.R. van der Heul, C. Vuik, and P. Wesseling. Stability analysis of segregated solution methods for compressible flow. *Applied Numerical Mathematics*, 38 :257–274, 2001.
- [70] D.R. van der Heul, C. Vuik, and P. Wesseling. A conservative pressure-correction method for flow at all speeds. *Computer & Fluids*, 32 :1113–1132, 2003.
- [71] J.P. Van Dormaal, G.D. Raithby, and B.H. McDonald. The segregated approach to predicting viscous compressible fluid flows. *Transactions of the ASME*, 109 :268–277, 1987.
- [72] D. Vidović, A. Segal, and P. Wesseling. A superlinearly convergent Mach-uniform finite volume method for the Euler equations on staggered unstructured grids. *Journal of Computational Physics*, 217 :277–294, 2006.
- [73] C. Wall, C.D. Pierce, and P. Moin. A semi-implicit method for resolution of acoustic waves in low Mach number flows. *Journal of Computational Physics*, 181 :545–563, 2002.
- [74] I. Wenneker, A. Segal, and P. Wesseling. A Mach-uniform unstructured staggered grid method. *International Journal for Numerical Methods in Fluids*, 40 :1209–1235, 2002.
- [75] P. Wesseling. Principles of computational fluid dynamics. volume 29 of *Springer Series in Computational Mathematics*. Springer, 2001.
- [76] O.C. Zienkiewicz and R. Codina. A general algorithm for compressible and incompressible flow – Part I. The split characteristic-based scheme. *International Journal for Numerical Methods in Fluids*, 20 :869–885, 1995.

