Jürgen Fuhrmann
Mario Ohlberger
Christian Rohde *Editors*

# Finite Volumes for Complex Applications VII - Elliptic, Parabolic and Hyperbolic Problems

FVCA 7, Berlin, June 2014

Springer

# Springer Proceedings in Mathematics & Statistics

Volume 78

**Springer Proceedings in Mathematics & Statistics**

This book series features volumes composed of selected contributions from workshops and conferences in all areas of current research in mathematics and statistics, including OR and optimization. In addition to an overall evaluation of the interest, scientific quality, and timeliness of each proposal at the hands of the publisher, individual contributions are all refereed to the high quality standards of leading journals in the field. Thus, this series provides the research community with well-edited, authoritative reports on developments in the most exciting areas of mathematical and statistical research today.

Jürgen Fuhrmann · Mario Ohlberger
Christian Rohde
Editors

# Finite Volumes for Complex Applications VII - Elliptic, Parabolic and Hyperbolic Problems

FVCA 7, Berlin, June 2014

Springer

*Editors*
Jürgen Fuhrmann
Weierstrass Institute for Applied
   Analysis and Stochastics
Berlin
Germany

Mario Ohlberger
Institute for Computational and
   Applied Mathematics and Center
   for Nonlinear Sciences (CeNoS)
University of Münster
Münster
Germany

Christian Rohde
Institute of Applied Analysis
   and Numerical Simulation
University of Stuttgart
Stuttgart
Germany

# Preface

The finite volume method in its various forms is a space discretization technique for partial differential equations based on the fundamental physical principle of conservation. It has been used successfully in many applications, including fluid dynamics, magnetohydrodynamics, structural analysis, nuclear physics, and semiconductor theory. Recent decades have brought significant success to the theoretical understanding of the method. Many finite volume methods preserve further qualitative or asymptotic properties, including maximum principles, dissipativity, monotone decay of the free energy, and asymptotic stability.

Due to these properties, finite volume methods belong to the wider class of compatible discretization methods, which preserve qualitative properties of continuous problems at the discrete level. This structural approach to the discretization of partial differential equations becomes particularly important for multiphysics and multiscale applications.

The triennial series of conferences "International Symposium on Finite Volumes for Complex Applications—Problems and Perspectives (FVCA)" brings together mathematicians, physicists, and engineers interested in this kind of physically motivated discretization. Contributions to the further advancement of the theoretical understanding of suitable finite volume, finite element, discontinuous Galerkin and other discretization schemes, and the exploration of new application fields have been welcomed.

Previous conferences on this series have been held in Rouen (1996), Duisburg (1999), Porquerolles (2002), Marrakech (2005), Aussois (2008), and Prague (2011).

The present volumes contain the invited and contributed papers presented as posters or talks at the Seventh International Symposium on Finite Volumes for Complex Applications held in Berlin, June 15–20, 2014.

The contributions in the first volume deal with the theoretical aspects of the method. They focus on topics such as preservation of physical properties on the discrete level, convergence, stability and error analysis, physically consistent coupling between discretizations for different processes, connections to other discretization methods, the relationship between grids and discretization schemes, complex geometries and adaptivity shock waves and other flow discontinuities, new and existing schemes and their limitations, and bottlenecks in the solution of large-scale problems.

As described, finite volume and related methods are of great practical value, as is demonstrated by the contributions to the second volume of the proceedings. Fields of application include atmospheric and ocean modeling, chemical engineering and combustion energy generation and storage, electro-reaction-diffusion systems, and porous media.

The volume editors thank the authors for their high quality contributions, the members of the program committee for supporting the organization of the review process, and all reviewers for their thorough work on the evaluation of each of the contributions.

The production of the proceedings was continuously supported by the Editor's team at Springer Verlag.

Without the financial contributions of the Deutsche Forschungsgemeinschaft (DFG), the Weierstrass Institute for Applied Analysis and Stochastics, the DFG Priority Program 1276 "Metström," the Westfälische Universität Münster, the Stuttgart Research Centre for Simulation Technology (Simtech) and the Czech Technical University of Prague, the organization of the conference and the production of the proceedings would not have been possible.

The Berlin Brandenburgische Akademie der Wissenschaften provided an impressive conference venue in the center of Berlin.

Finally, we thank the local organizers and the staff at the Weierstrass Institute for Applied Analysis and Stochastics for carrying the main organizational burden and for providing a friendly atmosphere at the conference.

March 2014                                                                   Jürgen Fuhrmann
                                                                            Mario Ohlberger
                                                                            Christian Rohde

# Organization Committees

## Organizing Committee

Peter Bastian
Robert Eymard
Jürgen Fuhrmann
Jiří Fürst
Annegret Glitzky
Volker John
Rupert Klein
Alexander Linke
Mario Ohlberger
Christian Rohde
Jörn Sesterhenn

## Proceedings Committee

Remi Abgrall
Brahim Amaziane
Boris Andreianov
Peter Bastian
Fayssal Benkhaldoun
Franck Boyer
Yves Coudière
Andreas Dedner
Vit Dolejsi
Jerome Droniou
Denis Dutykh
Alexandre Ern
Robert Eymard
Jürgen Fuhrmann
Jiří Fürst

Jan Giesselmann
Annegret Glitzky
Khaled Hassouni
Christiane Helzel
Jean-Marc Hérard
Danielle Hilhorst
Florence Hubert
Volker John
Rupert Klein
Robert Kloefkorn
Peter Knabner
Alexander Linke
Konstantin Lipnikov
Andreas Meister
Mario Ohlberger
Christian Rohde
Martin Rumpf
Jörn Sesterhenn
Martin Vohralik
Petra Wittbold

# Contents

**Part III   Elliptic and Parabolic Problems**

# Part III
# Elliptic and Parabolic Problems

# Asymptotic-Preserving Methods for an Anisotropic Model of Electrical Potential in a Tokamak

**Philippe Angot, Thomas Auphan and Olivier Guès**

**Abstract** A 2D nonlinear model for the electrical potential in the edge plasma in a tokamak generates a stiff problem due to the low resistivity in the direction parallel to the magnetic field lines. An asymptotic-preserving method based on a micro-macro decomposition is studied in order to have a well-posed problem, even when the parallel resistivity goes to 0. Numerical tests with a finite difference scheme show a bounded condition number for the linearised discrete problem solved at each time step, which confirms the theoretical analysis on the continuous problem.

**MSC2010**: 00B25, 41A60, 65M30

## 1 Introduction

The fusion reaction can be performed using a tokamak, a machine whose shape is toroidal. The plasma is confined and warmed in the core of the tokamak to produce the fusion reaction. This technique is expected to maintain the fusion reaction during a long time (more than five minutes, for the ITER project).

One of the main challenges for this objective is to control the wall-plasma interactions. Indeed, the magnetic confinement is not perfect and the plasma is in contact with the wall. In a tokamak such as TORE SUPRA, an obstacle called the limiter, is settled at the bottom of the machine. Due to the strong magnetic confinement, the plasma transport essentially occurs along the magnetic field lines. Thus, the parallel

P. Angot · T. Auphan (✉) · O. Guès
CNRS, Aix Marseille Université, Centrale Marseille, I2M, UMR 7373, 13453 Marseille, France
e-mail: thomas.auphan@univ-amu.fr; tauphan@cmi.univ-mrs.fr

P. Angot
e-mail: philippe.angot@univ-amu.fr

O. Guès
e-mail: olivier.gues@univ-amu.fr

resistivity $\eta$ is very small (typically, $\eta = 10^{-6}$), generating a strong anisotropy in the model. The area where the magnetic lines are interrupted by the limiter is called the scrape-off layer. The numerical simulation of the edge plasma transport allows us to better understand the interactions with the wall.

## 2 Anisotropic Model of the Electrical Potential

In this paper, we focus on a 2D model of the electrical potential of the edge plasma $\phi_\eta$ in a tokamak with a limiter configuration. A schematic representation of the domain is given in Fig. 1. The $x$ axis corresponds to the curvilinear coordinates along a magnetic field line and the $y$ axis is the radial direction. In the following equations, the curvature terms have been neglected. As the magnetic field lines above the limiter set are closed, periodic boundary conditions are imposed at $x = \pm 0.5$.

The dimensionless problem for the electrical potential reads:

$$\begin{cases} -\partial_t \partial_y^2 \phi_\eta - \frac{1}{\eta} \partial_x^2 \phi_\eta + \nu \partial_y^4 \phi_\eta = S & \text{in } ]0, T[ \times \Omega \\ \partial_y \phi_{\eta|t=0} = \partial_y \phi_{ini} & \text{in } \Omega \\ \partial_y \phi_{\eta|\Sigma_\parallel} = 0 \quad \text{and} \quad \partial_y^3 \phi_{\eta|\Sigma_\parallel} = 0 & \text{on } ]0, T[ \times \Sigma_\parallel \\ \partial_x \phi_{\eta|x=-L} = \eta \left(1 - e^{\Lambda - \phi_{\eta|x=-L}}\right) & \text{on } ]0, T[ \times ]0, l[ \times \{-L\} \\ \partial_x \phi_{\eta|x=L} = -\eta \left(1 - e^{\Lambda - \phi_{\eta|x=L}}\right) & \text{on } ]0, T[ \times ]0, l[ \times \{L\}, \end{cases} \tag{1}$$

where $\nu$ corresponds to the ionic viscosity in the perpendicular direction and $\Lambda$ stands for the reference potential inside the limiter. The initial condition is $\partial_y \phi_{\eta|t=0} = \partial_y \phi_{ini}$. Negulescu et al. [4] proved that, for a fixed value of $\eta > 0$, the problem (1) admits a unique weak solution, under suitable hypotheses on the data $\phi_{ini}$ and $S$.

The boundary conditions at the limiter interface $x = \pm L$, are nonlinear. Setting directly $\eta = 0$ in the system (1) (after multiplying the first equation by $\eta$) leads to an under-determined problem since there are only homogeneous Neumann boundary conditions at the limiter surface $x = \pm L$. Thus, when $\eta$ is small the numerical resolution of the problem (1) becomes stiff. This issue can be avoided by reformulating the problem (1) thanks to asymptotic-preserving methods.

# 3 The Micro-Macro Asymptotic-Preserving Method

We study the Asymptotic-Preserving (AP) method introduced by Degond et al. [3] for a linear anisotropic elliptic problem. It consists in a decomposition of the solution $\phi_\eta$ as $\phi_\eta = p_\eta + \eta q_\eta$ where $\partial_x p_\eta = 0$ and $q_{\eta|x=-L} = 0$. Then, it yields the problem below where the unknowns are $(\phi_\eta, q_\eta)$:

$$
\begin{cases}
-\partial_t \partial_y^2 \phi_\eta - \partial_x^2 q_\eta + \nu \partial_y^4 \phi_\eta = S & \text{in } ]0, T[ \times \Omega \\
\partial_x^2 \phi_\eta = \eta \partial_x^2 q_\eta & \text{in } ]0, T[ \times \Omega \\
\partial_x \phi_{\eta|x=-L} = \eta \partial_x q_{\eta|x=-L} & \text{on } ]0, T[ \times ]0, l[ \times \{-L\} \\
\partial_x \phi_{\eta|x=L} = \eta \partial_x q_{\eta|x=L} & \text{on } ]0, T[ \times ]0, l[ \times \{L\} \\
\partial_x \phi_{\eta|x=-0.5} = \eta \partial_x q_{\eta|x=-0.5} & \text{on } ]0, T[ \times ]l, 1[ \times \{-0.5\} \\
\partial_x \phi_{\eta|x=0.5} = \eta \partial_x q_{\eta|x=0.5} & \text{on } ]0, T[ \times ]l, 1[ \times \{0.5\} \\
\partial_y \phi_{\eta|t=0} = \partial_y \phi_{ini} & \text{in } \Omega \\
\partial_y \phi_{\eta|\Sigma_\|} = 0 \quad \text{and} \quad \partial_y^3 \phi_{|\Sigma_\|} = 0 & \text{on } ]0, T[ \times \Sigma_\| \\
\partial_x q_{\eta|x=-L} = \left(1 - e^{\Lambda - \phi_{\eta|x=-L}}\right) & \text{on } ]0, T[ \times ]0, l[ \times \{-L\} \\
\partial_x q_{\eta|x=L} = -\left(1 - e^{\Lambda - \phi_{\eta|x=L}}\right) & \text{on } ]0, T[ \times ]0, l[ \times \{L\},
\end{cases}
\tag{2}
$$

One important advantage of this AP method is that it can be easily implemented even if the mesh is not aligned with the directions $(Ox)$ and $(Oy)$. The main drawback is the need to compute two unknowns ($\phi_\eta$ and $q_\eta$) on the 2D domain though only $\phi_\eta$ is interesting for the physics.

Let us give the theoretical result which ensures that the modified problem is well-posed for $\eta = 0$, and that $\phi_\eta$ converges towards $\phi_0$. First, we provide the definitions of the spaces used for the variational formulation of the problem (2).

**Definition 1** Let us define the following Hilbert spaces:

- $V = \left\{ f \in H^1(\Omega), \partial_y^2 f \in L^2(\Omega), f \text{ periodic on } \{-0.5, 0.5\} \times ]l, 1[, \partial_y f = 0 \text{ on } \Sigma_\| \right\}$ with the scalar product:

$$
\langle f, u \rangle_V = \int_\Omega \partial_x f \, \partial_x u \, dy dx + \int_\Omega \partial_y^2 f \, \partial_y^2 u \, dy dx + 2 \int_0^l f_{|x=L} \, u_{|x=L} \, dy.
$$

- $Q = \left\{ f \in L^2(\Omega), \partial_x f \in L^2(\Omega), f_{|x=-L} = 0 \text{ on } ]0, 1[ \right\}$, with the scalar product:

$$
\langle f, u \rangle_Q = \int_\Omega \partial_x f \, \partial_x u \, dy dx.
$$

**Definition 2** The space $\mathscr{A}$ is the set of functions $\phi$ such that:

- $\phi \in L^2(0, T; V)$.
- $\partial_y \phi \in L^\infty(0, T; L^2(\Omega))$.
- $\partial_y \phi \in L^2\left(0, T; \{f \in H^1(\Omega), \partial_y^2 f \in L^2(\Omega), f_{|\Sigma_\|} = 0\}\right)$.

- $\partial_y^2\phi \in L^\infty(0, T; L^2(\Omega))$.
- $\partial_t\phi \in L^2(0, T; V)$.
- $\partial_y\partial_t\phi \in L^\infty(0, T; L^2(\Omega))$.

The weak solution $\phi_\eta$ of (2) is then searched in the space $\mathscr{A}$.

**Assumption 31** *Assume that $S$ and $\phi_{ini}$ verify:*

1. $S, \partial_y S, \partial_y^2 S, \partial_t S, \partial_t^2 S \in L^2(]0, T[\times\Omega)$, $\|S\|_{L^\infty(]0,T[\times\Omega)} \le C_s$ *and*
   $\|S_{|t=T}\|_{L^\infty(\Omega)} \le C_s$ *with $C_s$ sufficiently small.*
2. $\phi_{ini} \in H^4(\Omega)$.
3. $\phi_{ini}$ *does not depend on $x$.*
4. $\displaystyle\int_\Omega S_{|t=0}\,dydx = v\int_\Omega \partial_y^4\phi_{ini}\,dydx + 2\int_0^l \left(1 - e^{\Lambda-\phi_{ini|x=L}}\right) dy.$

The two last hypotheses are compatibility conditions for the initial and boundary conditions with the source term.

We can now write the theorem which asserts the convergence of $\phi_\eta$ to $\phi_0$ when $\eta$ goes to 0:

**Theorem 1** *With the assumption* 3.1, *the weak formulation of* (2):
*find $(\phi_\eta, q_\eta) \in \mathscr{A} \times L^2(0, T; Q)$ verifying*

$$
\begin{cases}
\forall\xi \in H^1(]0, T[), \forall u \in V \cap H^2(\Omega), \forall w \in Q, \\[4pt]
\displaystyle\int_\Omega \partial_y\phi_{\eta|t=T}\,\partial_y u\,dydx\,\xi(T) - \int_0^T \int_\Omega \partial_y\phi_{\eta|t=T}\,\partial_y u\,dydx\,\xi'\,dt \\[8pt]
\displaystyle+ \int_0^T \int_\Omega \partial_x q_\eta\,\partial_x u\,dydx\,\xi\,dt + v\int_0^T \int_\Omega \partial_y^2\phi_\eta\,\partial_y^2 u\,dydx\,\xi\,dt \\[8pt]
\displaystyle+ \int_0^T \int_0^l \left(1 - e^{\Lambda-\phi_{\eta|x=-L}}\right) u_{|x=-L}\,dy\,\xi\,dt + \int_0^T \int_0^l \left(1 - e^{\Lambda-\phi_{\eta|x=L}}\right) u_{|x=L}\,dy\,\xi\,dt \\[8pt]
\displaystyle= \int_\Omega \partial_y\phi_{ini}\,\partial_y u\,dydx\,\xi(0) + \int_0^T \int_\Omega S u\,dydx\,\xi\,dt \\[8pt]
\displaystyle\eta\int_0^T \int_\Omega \partial_x q_\eta\,\partial_x w\,dydx\,\xi\,dt = \int_0^T \int_\Omega \partial_x\phi_\eta\,\partial_x w\,dydx\,\xi\,dt,
\end{cases}
\tag{3}
$$

*admits a unique solution. Besides, $(\phi_\eta, q_\eta)$ converges weakly in $L^2(]0, T[\times\Omega)^2$, towards $(\phi_0, q_0) \in \mathscr{A} \times L^2(0, T; Q)$ the solution of* (3) *when $\eta$ equals 0.*
*Finally, the following error estimate holds:*

$$\|\phi_\eta - \phi_0\|_{L^1(0,T;L^2(\Omega))} \le c(T, \Omega, \phi_0, S, \Lambda)\,\sqrt{\eta},$$

*where $c(T, \Omega, \phi_0, S, \Lambda) > 0$ does not depend on $\eta$.*

Theorem 1 provides an error estimate for the norm in $L^1(0, T; L^2(\Omega))$, but not for the $L^2(]0, T[\times\Omega)$ norm. This point can be subject to further improvements.

This result is shown in [1, 2]. The proof of the existence and uniqueness of $\phi_0$ follows the same steps of [4], based on a fixed point method. The existence and

uniqueness of $q_0$ and the convergence of $(\phi_\eta, q_\eta)$ when $\eta$ goes to 0 are shown by extending to a nonlinear case the proof provided in [3] for a linear elliptic problem.

## 4 Numerical Experiments

In this section, some numerical tests are presented for the system (2). The space discretisation is done by the centred finite difference scheme. The time resolution uses Euler semi-implicit method.

At first glance, a directional splitting method seems to be interesting. But, the discrete problems obtained in the directions $x$ and $y$ are not invertible. The problem is thus discretised implicitly, except for the nonlinear term. At each time step, a linear system has to be solved to compute the approximations of $\phi_\eta$ and $q_\eta$.

Let us consider a rectangular mesh of the space domain $\Omega$ with a constant mesh step $\delta x$ (for the direction $(Ox)$) and $\delta y$ (for the direction $(Oy)$). The time step writes $\delta t$. The scalar quantities $\phi_{i,j}^n$, $q_{i,j}^n$ stands respectively for the approximations of $\phi_\eta(n\delta t, -0.5 + i\delta x, j\delta y)$ and $q_\eta(n\delta t, -0.5 + i\delta x, j\delta y)$. The boundary condition at $x = -L$ is discretised as:

$$\frac{q_{I_1+1,j}^{n+1} - q_{I_1-1,j}^{n+1}}{2\delta x} - \phi_{I_1,j}^{n+1} = \left(1 - e^{\Lambda - \phi_{I_1,j}^n} - \phi_{I_1,j}^n\right),$$

where $I_1$ is the index such that $-0.5 + I_1\delta x = -L$.

For the boundary condition at $x = L$, the same technique is used. This time linearisation enables us to have an invertible matrix which is the same at each time step.

The mesh convergence test is performed using a configuration where the limiter goes up to the top of the computational domain, i.e. $l = 1$. This does not change the results proven for $l < 1$. For $L = 0.4$, the chosen manufacturated solution is

$$\phi_\eta(t, x, y) = \eta \left(\frac{t}{\pi}\right)^2 \cos(\pi y) \cos(1.25\pi x) - \ln\left(1 - \frac{1.25t^2}{\pi \cos(\pi y)}\right) + \Lambda. \quad (4)$$

Let us note that the source term $S$ associated to the manufactured solution (4) depends on $\eta$ but is not singular when $\eta$ goes to 0. This differs from the hypotheses made for Theorem 1.

The plot of the approximated solution is shown in Fig. 2. Studying the $L^2$ error in Fig. 3, we observe that the numerical scheme is of second-order accuracy in space.

In Fig. 4, we observe that the condition number obtained with the AP method is high but it is bounded independently from $\eta$. This is not the case for the matrix obtained for the resolution of (1) without the asymptotic-preserving method. In order to avoid the issues due to the bad conditioning, we choose a LU method to solve the linear problem at each time step, which is faster than a GMRES solver with

phi approx vs x,y (dx=dy=0.003125 , dt=0.0001, eta=0.001, t=1)     q approx vs x,y (dx=dy=0.003125 , dt=0.0001, eta=0.001, t=1)

**Fig. 2** Approximate fields of $\phi_\eta$ and $q_\eta$ for $\delta x = \delta y = 0.003125$, $\delta t = 0.0001$ and $\eta = 0.001$. The reference solution is given by (4). Recall that the limiter area corresponds to $x \leq -0.4$ and $x \geq 0.4$: the values of $\phi_\eta$ do not have any physical sense in this zone

**Fig. 3** $\|\phi_\eta^{approx} - \phi_\eta\|_{L^2(\Omega)}$ at $t = 1$ as a function of the space step $\delta x = \delta y$ for different values of the time step and $\eta = 0.001$. The reference solution is given by (4)

L2 error (phi_eta approx - phi_eta exact) versus dx (= dy) ; t= 1

PETSc library. Finding an efficient preconditioner in order to use iterative methods is a future enhancement of this work.

For the convergence when $\eta$ tends to 0, the same domain is considered ($l = 1$, $L = 0.4$) but another source term is chosen:

$$S(t, x, y) = 40\, t \cos(2\pi\, y) \sin\left(\frac{\pi}{2L}\, x\right) \quad , \quad \phi_{ini}(x, y) = \Lambda = 0 \qquad (5)$$

This configuration (5) with $l = 1$ leads to $\phi_0(t, x, y) = 0$, which enables us to compute numerically $\|\phi_\eta - \phi_0\|_{L^1(0,T;L^2(\Omega))}$ and $\|\phi_\eta - \phi_0\|_{L^2(0,T;L^2(\Omega))}$. For these two norms, we observe a convergence in $\mathcal{O}(\eta)$, see Fig. 5. This suggests that the estimate of Theorem 1 might be improved.

**Fig. 4** Condition number
in the Euclidean norm as
a function of the parallel
resistivity $\eta$ for the linear
system approaching (the same
at each time step) the solution
(4) with $\delta x = \delta y = 0.025$ and
$\delta t = 0.001$



Condition number versus eta (dx= dy= 0.025, dt= 0.01)

**Fig. 5** $\|\phi_\eta - \phi_0\|_{L^1(0,T;L^2(\Omega))}$
($\Delta$) and $\|\phi_\eta - \phi_0\|_{L^2(]0,T[\times\Omega)}$
($+$) as a function of $\eta$. The
configuration is given by
Eq. (5) with $T = 1$, $\delta x =$
$\delta y = 0.003125$ and $\delta t =$
0.0001



Error (phi_eta approx- phi_0 exact) vs eta (dt= 0.0001, dx= dy= 0.003125)

## 5 Conclusion

The high anisotropy of the 2D model for the edge plasma electrical potential in a
tokamak leads to an ill-conditioned matrix for the numerical approximation using
classical methods. The micro-macro decomposition induced by Degond et al. [3]
for a linear anisotropic elliptic problem is studied and analysed for the nonlinear
evolution problem of the electrical potential. This method yields a weak formulation
which is not degenerated when the parallel resistivity $\eta$ tends to 0. Moreover, we
have the estimate

$$\|\phi_\eta - \phi_0\|_{L^1(0,T,L^2(\Omega))} = \mathcal{O}\left(\sqrt{\eta}\right),$$

which can probably be improved, as suggested by the numerical results.

# References

1. Angot, P., Auphan, T., Guès, O.: Analysis of asymptotic preserving methods for nonlinear anisotropic models of electrical potential in plasma (in preparation) (2014)
2. Auphan, T.: Analyse de modèles pour ITER; Traitement des conditions aux limites de systèmes modélisant le plasma de bord dans un tokamak. Ph.D. thesis in mathematics, Aix Marseille Université (2014)
3. Degond, P., Lozinski, A., Narski, J., Negulescu, C.: An asymptotic-preserving method for highly anisotropic elliptic equations based on a micromacro decomposition. J. Comput. Phys. **231**(7), 2724–2740 (2012)
4. Negulescu, C., Nouri, A., Ghendrih, P., Sarazin, Y.: Existence and uniqueness of the electric potential profile in the edge of tokamak plasmas when constrained by the plasma-wall boundary physics. Kinet. Relat. Models **1**(4), 619– 639 (2008)

# Semi-implicit Second Order Accurate Finite Volume Method for Advection-Diffusion Level Set Equation

**Martin Balažovjech, Peter Frolkovič, Richard Frolkovič and Karol Mikula**

**Abstract**    We present a second order accurate finite volume method for level set equation describing the motion in normal direction with the speed depending on external properties and curvature. A convenient combination of a Crank-Nicolson type of the time discretization for diffusion term [1] and an Inflow Implicit and Outflow Explicit scheme [6] for advection term is used. Numerical experiments for an example with the exact solution derived in this paper and for examples motivated by modeling of fire propagation in forests are presented.

## 1 Introduction

Although not in a divergence form, the level set equations are often solved with finite volume methods [3–5, 8]. The basic idea behind such approaches is to rewrite the level set equation in such a way that it can be approximated using integration by parts. In this paper we apply such approach with an aim to suggest a second order accurate finite volume method to solve level set equations that describe the motion in normal direction with the speed depending on external properties and on curvature.

In the level set equation one can recognize two terms that have a character of advection and diffusion, respectively. In [6, 7] a novel second order accurate

M. Balažovjech · P. Frolkovič (✉) · R. Frolkovič · K. Mikula
Department of Mathematics, Slovak University of Technology, Radlinského 11,
813 68 Bratislava, Slovak Republic
e-mail: peter.frolkovic@stuba.sk

M. Balažovjech
e-mail: balazovjech@stuba.sk

K. Mikula
e-mail: karol.mikula@stuba.sk

semi-implicit finite volume discretization is used for the advection where the inflow parts of finite volume boundaries are treated implicitly in time and the outflow parts are treated in an explicit way. Our idea is to combine such approach with a second order accurate approximation of the curvature term using a procedure similar to the Crank-Nicolson method. The latter method is successfully used in a Lagrangian type of method for curvature driven flow in [1].

In this paper we propose a particular finite volume scheme of this type. The scheme treats the advection and diffusion fluxes in a compatible way. The resulting system of semilinear algebraic equations has favorable properties that can be used conveniently to solve it. When fixing the nonlinear coefficients in algebraic equations, the resulting matrix is a M-matrix and iterative solvers like the Gauss-Seidel method can be used to solve the linearized algebraic system.

Second order accurate methods for purely advective type of equations need in general some stabilization ("limiter") techniques to suppress nonphysical oscillations in numerical solutions [4, 6, 7]. In the presence of curvature driven motion as in our case we need not to apply such techniques if the advection is not too strong.

The paper is organized as follows. In Sect. 2 we derive briefly the level set equation that we want to solve. In Sect. 3 the finite volume method is derived. The Sect. 4 introduces a method for the solution of nonlinear algebraic equations. In Sect. 5 we derive a representative exact solution of the level set equation and present experimental order of convergence for our numerical method. Moreover, examples motivated by the modeling of fire front propagation in forests are presented. Finally, in Sect. 6 we conclude briefly our results.

## 2 Mathematical Model

Let $u = u(x, t)$, $(x, t) \in D \times [0, T]$ be the so called level set function used e.g. to represent implicitly an evolving interface. We denote $\mathbf{n} := \nabla u / |\nabla u|$ when $|\nabla u| \neq 0$. Note that $\mathbf{n}(\bar{x})$ is the normal vector at $\bar{x}$ to the level set given by $u(x, t) = u(\bar{x}, t)$.

We search $u = u(x, t)$ for $(x, t) \in D \times (0, T]$ fulfilling the level set equation

$$\frac{\partial u}{\partial t} + (f + \delta k)\, \mathbf{n} \cdot \nabla u = 0\,, \quad u(x, 0) = u^0(x)\,. \tag{1}$$

In (1) the term $f + \delta k$ represents a speed in normal direction $\mathbf{n}$ with $f(x)$ and $\delta(x) > 0$ being given. The function $k$ denotes the curvature that is defined by

$$k = -\nabla \cdot \left( \frac{\nabla u}{|\nabla u|} \right)\,. \tag{2}$$

Substituting (2) to (1) one obtains the nonlinear advection-diffusion level set equation of the form

$$\frac{\partial u}{\partial t} + \left( f \frac{\nabla u}{|\nabla u|} \right) \cdot \nabla u = \delta |\nabla u| \nabla \cdot \left( \frac{\nabla u}{|\nabla u|} \right), \quad u(x, 0) = u^0(x). \tag{3}$$

## 3 Finite Volume Method

Before discretizing (3) we divide it by $|\nabla u|$ and rewrite the advection term as in [3] to obtain

$$\frac{1}{|\nabla u|} \frac{\partial u}{\partial t} + \nabla \cdot \left( u f \frac{\nabla u}{|\nabla u|^2} \right) - u \nabla \cdot \left( f \frac{\nabla u}{|\nabla u|^2} \right) = \delta \nabla \cdot \left( \frac{\nabla u}{|\nabla u|} \right). \tag{4}$$

For simplicity we consider the domain $D \subset R^2$ to be a square and the finite volume mesh to consist of squared elements $p_{ij}$, $i, j = 1, 2, \ldots, N$ having uniform length $h > 0$ for all edges. The edges of $p_{ij}$ are denoted by $l_k$, $k \in \Lambda_{ij}$ where $\Lambda_{ij} = \{(i + 1/2\, j), (i\, j + 1/2), (i - 1/2\, j), (i\, j - 1/2)\}$ is the set of indices for particular edges of $p_{ij}$.

Furthermore, we consider a uniform time step $\Delta t$ and $t^m = m \Delta t$. The numerical solution of (4) will be represented by the discrete unknown values $u_{ij}^m$ that approximates $u$ in $p_{ij} \times (t^{m-1}, t^m]$.

The idea of a finite volume discretization for (4) is to integrate it over $p_{ij}$ and to use appropriate quadrature rules that we explain for each term separately. Firstly,

$$\int_{p_{ij}} \frac{1}{|\nabla u|} \frac{\partial u}{\partial t} dx \approx \frac{h^2}{|\nabla u|_{ij}} \frac{du_{ij}}{dt}, \tag{5}$$

where the value $|\nabla u|_{ij}$ and the time discretization of $u_{ij} = u_{ij}(t)$ will be introduced later. Next,

$$\int_{p_{ij}} \nabla \cdot \left( f u \frac{\nabla u}{|\nabla u|^2} \right) dx = \sum_k \oint_{l_k} \left( \frac{f u}{|\nabla u|^2} \frac{\partial u}{\partial n} \right) ds \approx h \sum_k \left( \frac{f_k \bar{u}_k}{|\nabla u|_k^2} \frac{\partial u}{\partial n} \bigg|_{l_k} \right) \tag{6}$$

and

$$\int_{p_{ij}} u \nabla \cdot \left( \frac{f \nabla u}{|\nabla u|^2} \right) dx \approx \bar{u}_{ij} \int_{p_{ij}} \nabla \cdot \left( \frac{f \nabla u}{|\nabla u|^2} \right) dx \approx h \bar{u}_{ij} \sum_k \left( \frac{f_k}{|\nabla u|_k^2} \frac{\partial u}{\partial n} \bigg|_{l_k} \right). \tag{7}$$

The value $f_k$ denotes an averaged value of $f$ at $l_k$. Furthermore, $\bar{u}_k$ represents a reconstructed value of $u$ assigned to $l_k$ and $\bar{u}_{ij}$ is a reconstructed value of $u$ assigned to $p_{ij}$ [6, 7]. Particular choices for their computations, together with the approximations of $|\nabla u|_k$ and the normal derivatives $\partial u / \partial n$, will be introduced later.

Finally, analogous rules are applied for the last term in (4) to obtain

$$
\int_{p_{ij}} \delta \nabla \cdot \left( \frac{\nabla u}{|\nabla u|} \right) dx \approx \delta_{ij} \sum_k \oint_{l_k} \left( \frac{1}{|\nabla u|} \frac{\partial u}{\partial n} \right) ds \approx h \delta_{ij} \sum_k \left( \frac{1}{|\nabla u|_k} \left. \frac{\partial u}{\partial n} \right|_{l_k} \right),
$$
(8)

where $\delta_{ij}$ is an averaged value of $\delta$ with respect to $p_{ij}$.

Putting all approximations (5)–(8) together, we obtain a compact form of our finite volume discretization method

$$
\frac{h^2}{|\nabla u|_{ij}} \frac{du_{ij}}{dt} - h \sum_k \left( f_k \frac{\bar{u}_{ij} - \bar{u}_k}{|\nabla u|_k^2} \left. \frac{\partial u}{\partial n} \right|_{l_k} \right) = h \sum_k \left( \frac{\delta_{ij}}{|\nabla u|_k} \left. \frac{\partial u}{\partial n} \right|_{l_k} \right).
$$
(9)

We define now the missing approximations in (9). Firstly, we define $u_k$, $k \in \Lambda_{ij}$ by a linear interpolation,

$$
u_{i+1/2\, j} := \frac{u_{ij} + u_{i+1\, j}}{2}, \quad u_{i-1/2\, j} := \frac{u_{ij} + u_{i-1\, j}}{2}, \quad \text{and so on.}
$$

The normal derivatives are approximated in a standard way,

$$
\left. \frac{\partial u}{\partial n} \right|_{l_{i+1/2\, j}} \approx \frac{u_{i+1/2\, j} - u_{ij}}{h/2} = \frac{u_{i+1\, j} - u_{ij}}{h}, \quad \left. \frac{\partial u}{\partial n} \right|_{l_{i-1/2\, j}} \approx \frac{u_{i-1/2\, j} - u_{ij}}{h/2}, \quad \text{and so on.}
$$

To approximate $\nabla u$ at the edges $l_k$ of $p_{ij}$, we use the diamond cell formula. To do so we use the notation $u_{i\pm\frac{1}{2}\, j\pm\frac{1}{2}}$ for the four values of $u$ in the corners of $p_{ij}$ that are obtained as arithmetic averages

$$
u_{i\pm\frac{1}{2}\, j\pm\frac{1}{2}} := \frac{1}{4} \left( u_{ij} + u_{i\pm1\, j} + u_{i\, j\pm1} + u_{i\pm1\, j\pm1} \right).
$$

Using it, we can approximate $|\nabla u|$ at the edges $l_k$, $k \in \Lambda_{ij}$ of $p_{ij}$ by

$$
|\nabla u|_{i+1/2\, j} \approx \sqrt{ \left( \frac{u_{i+1\, j} - u_{ij}}{h} \right)^2 + \left( \frac{u_{i+\frac{1}{2}\, j+\frac{1}{2}} - u_{i+\frac{1}{2}\, j-\frac{1}{2}}}{h} \right)^2 + \varepsilon^2 }, \quad \text{and so on.}
$$

A regularization was introduced in above formula by choosing $0 < \varepsilon << 1$ to avoid a division by zero in (9). Furthermore,

$$
|\nabla u|_{ij} \approx \frac{1}{4} \sum_{k \in \Lambda_{ij}} |\nabla u|_k.
$$
(10)

Finally, we have to define in (9) the reconstructed values $\bar{u}_{ij}$ and $\bar{u}_k$. Following [6] we take simply $\bar{u}_{ij} = u_{ij}$ and $\bar{u}_k = u_k$. This choice works well when the advection

does not dominate the diffusion term in (3), in general more sophisticated choices have to be taken into account, see [6, 7].

Summarizing all approximations used in (9) we obtain

$$\frac{h^2}{|\nabla u|_{ij}} \frac{du_{ij}}{dt} = 2 \sum_{k \in \Lambda_{ij}} \frac{1}{|\nabla u|_k} \left( f_k \frac{u_{ij} - u_k}{|\nabla u|_k} + \delta_{ij} \right) (u_k - u_{ij}) . \qquad (11)$$

To introduce formally a second order accurate time discretization of (11) we treat the advection and diffusion term separately. We begin with the time discretization of the curvature term. Inspired by [1] we use a Crank-Nicolson type of time discretization that can be viewed as an arithmetic average of fully explicit and fully implicit time discretization scheme,

$$\frac{h^2}{2 \Delta t} \left( \frac{1}{|\nabla u|_{ij}^{m+1}} + \frac{1}{|\nabla u|_{ij}^{m}} \right) (u_{ij}^{m+1} - u_{ij}^{m}) =$$

$$\delta_{ij} \sum_{k \in \Lambda_{ij}} \frac{1}{|\nabla u|_k^{m+1}} \left( u_k^{m+1} - u_{ij}^{m+1} \right) + \delta_{ij} \sum_{k \in \Lambda_{ij}} \frac{1}{|\nabla u|_k^{m}} \left( u_k^{m} - u_{ij}^{m} \right), \qquad (12)$$

where $|\nabla u|_{ij}^{m}$ and $|\nabla u|_{ij}^{m+1}$ are computed from (10) at corresponding time levels.

To discretize the advection term in time we introduce the notation in which we distinguish between the edges $l_k$ of $p_{ij}$ with an inflow and outflow character, namely

$$a_k^{in} = max \left( f_k \frac{u_{ij}^{m+1} - u_k^{m+1}}{|\nabla u|_k^{m+1}}, 0 \right), \quad a_k^{out} = min \left( 0, f_k \frac{u_{ij}^{m} - u_k^{m}}{|\nabla u|_k^{m}} \right). \qquad (13)$$

The advection term can be approximated by the "Inflow Implicit/Outflow Explicit" time discretization [7] to obtain

$$\frac{h^2}{2 \Delta t} \left( \frac{1}{|\nabla u|_{ij}^{m+1}} + \frac{1}{|\nabla u|_{ij}^{m}} \right) (u_{ij}^{m+1} - u_{ij}^{m}) =$$

$$\sum_{k \in \Lambda_{ij}} \frac{2 a_k^{in}}{|\nabla u|_k^{m+1}} \left( u_k^{m+1} - u_{ij}^{m+1} \right) + \sum_{k \in \Lambda_{ij}} \frac{2 a_k^{out}}{|\nabla u|_k^{m}} \left( u_k^{m} - u_{ij}^{m} \right). \qquad (14)$$

Putting (12) and (14) together we obtain

$$\frac{h^2}{2 \Delta t} \left( \frac{1}{|\nabla u|_{ij}^{m+1}} + \frac{1}{|\nabla u|_{ij}^{m}} \right) \left( u_{ij}^{m+1} - u_{ij}^{m} \right) =$$

$$\sum_{k \in \Lambda_{ij}} \frac{2 a_k^{in} + \delta_{ij}}{|\nabla u|_k^{m+1}} \left( u_k^{m+1} - u_{ij}^{m+1} \right) + \sum_{k \in \Lambda_{ij}} \frac{2 a_k^{out} + \delta_{ij}}{|\nabla u|_k^{m}} \left( u_k^{m} - u_{ij}^{m} \right). \qquad (15)$$

## 4 Solution of Algebraic Equations

In this section we briefly comment how to solve the algebraic system of equations represented by the discretization scheme (15).

The values $u_{ij}^0$ are computed from the initial condition. In our numerical experiments we consider only the Dirichlet type of boundary conditions. Consequently, one has to solve (15) for the unknowns $\{u_{ij}^{m+1}, \ i, j = 1, 2, \ldots, N-1\}$ in a sequence for $m = 0, 1$ and so on.

We propose to solve (15) using a combination of fixed point iterations and Gauss-Seidel iterative method. To introduce it we define for $p = -1, 0, 1$ and $q = -1, 0, 1$ that fulfill $|p| + |q| = 1$ the following coefficients

$$\lambda_{ij} = \Delta t \frac{|\nabla u_{ij}|^{m+1} + |\nabla u_{ij}|^m}{h^2}, \quad M_{ij}^{pq} = \frac{2a_{i+p/2\,j+q/2}^{in} + \delta_{ij}}{|\nabla u|_{i+p/2\,j+q/2}^{m+1}} \tag{16}$$

$$M_{ij} = \sum_{|p|+|q|=1} M_{ij}^{pq}, \quad b_{ij} = \sum_{|p|+|q|=1} \frac{2a_{i+p/2\,j+q/2}^{out} + \delta_{ij}}{|\nabla u|_{i+p/2\,j+q/2}^m} \left( u_{i+p\,j+q}^m - u_{ij}^m \right). \tag{17}$$

Using (16)–(17) the scheme (15) can be written in the form

$$u_{ij}^{m+1} = \frac{1}{1 + \lambda_{ij} M_{ij}} \left( u_{ij}^m + \lambda_{ij} \left( b_{ij} + \sum_{|p|+|q|=1} M_{ij}^{pq} u_{i+p\,j+q}^{m+1} \right) \right). \tag{18}$$

We note that the coefficients defined in (16) are nonlinear and always positive.

The iterative method consists of the following steps. Firstly, an initial guess for the unknowns $u_{ij}^{m+1}$ is set to the available values $u_{ij}^m$ from the previous time step or from the initial conditions if $m = 0$. Moreover, the coefficients $b_{ij}$ in (17) are computed only once in each time step.

Each iteration of our iterative method is realized by computing the nonlinear coefficients defined in (16) using the values computed from the previous iteration. Fixing these coefficients one can update the values $u_{ij}^{m+1}$ according to (18) for $i, j = 1, 2, \ldots, N-1$ by evaluating the values of $u_{i+p\,j+q}^{m+1}$ on the right hand side of (18) in a manner of Gauss-Seidel iterative method.

## 5 Numerical Experiments

At first we derive an exact solution in a simplified situation when an evolving curve is a circle initially, and it evolves according to (1) with constant values of $f$ and $\delta$. In such case the evolving curve preserves its circular shape, so it can be described by its radius $r = R(t, r^0)$ where $r^0 = R(0, r^0)$.

**Table 1** The comparison of numerical solution obtained with (15) with the exact solution (20)

| N | Error | EOC | #$it$ |
|---|---|---|---|
| 16 | 6.27e-2 | – | 26 |
| 32 | 1.00e-2 | 2.64 | 38 |
| 64 | 1.87e-3 | 2.42 | 59 |
| 128 | 4.41e-4 | 2.09 | 102 |

Let $u^0(r)$ be a given increasing function and $u(x_1, x_2, 0) = u^0(r)$, $r = \left(x_1^2 + x_2^2\right)^{1/2}$. Clearly, any circle of radius $r^0$ consists of points $(x_1, x_2)$ such that $u(x_1, x_2, 0) = u^0(r^0)$. Our aim is to find $u(x_1, x_2, t)$ such that $u(x_1, x_2, t) = u^0(r^0)$ for all points $(x_1, x_2)$ that fulfill $\left(x^2 + y^2\right)^{1/2} = R(t, r^0)$. To do so the inverse function of $r = R(t, r^0)$ with respect to $r^0$ must exist, i.e. $r^0 = R^{-1}(t, r)$. Once available, one obtains $u(x_1, x_2, t) = u^0(R^{-1}(t, \sqrt{x^2 + y^2}))$.

If a circular curve expands or shrinks with a constant speed $f$ and $\delta$, the radius $r(t)$ shall fulfill the equation $\dot{r}(t) = f + \frac{\delta}{r}$, $r(0) = r_0$ which is solved by

$$R(t, r^0) = \frac{\delta}{f} + \frac{\delta}{f} W\left(\frac{1}{\delta}\left(fr^0 - \delta\right) e^{\frac{-\delta + fr^0 + f^2 t}{\delta}}\right) \tag{19}$$

where $W$ is the product log function, i.e. $W(z)$ is obtain such that $z = We^W$.

Let us choose as initial function $u^0(x_1, x_2) = \sqrt{x_1^2 + x_2^2}$. Using our approach one obtains the solution of (3) for constant $f$ and $\delta$ in the form

$$u(x_1, x_2, t) = \frac{\delta}{f} + \frac{\delta}{f} W\left(\frac{1}{\delta}\left(f\sqrt{x_1^2 + x_2^2} - \delta\right) e^{-\frac{\delta - f\sqrt{x_1^2 + x_2^2} + f^2 t}{\delta}}\right). \tag{20}$$

In Table 1 we present the comparison of numerical solution obtained by (15) with the exact solution (20) for $f = \delta = 1$ and $t \in [0, 1]$ using a standard $l_2$ discrete norm in time and space. The domain $D$ is a square with the side length $L = 8$. The Dirichlet boundary conditions defined by the available exact solution are used on $\partial D$. The discretization step is taken $h = 8/N$ for $N = 16, 32, 64, 128$, the time step is chosen $\Delta t = h/2$.

One can see from Table 1 that for this example the experimental order of convergence is approaching 2 from above. Moreover we present the number of iterations for each $N$ that were necessary to reduce the residuum below the value $10^{-10}$.

In the following illustrative examples we are motivated by numerical simulation of fire front propagation in forests [2]. The parameter $f(x)$ defines how fast the underlying forest can burn and $\delta(x) = \mu f(x)$ where $\mu > 0$, so the speed in normal direction **n** is given by $f(x)(1 + \mu k)$.

The first example shows a behavior for inhomogeneous forest, see Fig. 1. The second example illustrates a topological change when the evolving fire front, being a circle initially, has to surround later a small area that can not burn, see Fig. 2.

**Fig. 1** Picture of fire front position at different time levels. The *smallest circle* is the initial position of the front. The parameters are $f = 1$ left and $f = 0.2$ right, $\mu = 0.1$



**Fig. 2** Pictures of fire front position at 4 different time levels, the *top row* with 3D view, the *bottom row* 2D view. The small *black square* can not burn ($f = 0$), the north-east region is less burnable ($f = 0.2$) than the rest ($f = 1$). The *small circle* in 2D view is the initial position of the front

## 6 Conclusions

Our novel finite volume method combines conveniently explicit and implicit time discretization to obtain the second order accurate numerical solution of level set equation containing the terms of advection and diffusion character. For the chosen representative example for which the exact solution is derived, we can report experimental order of convergence approaching the value 2 from above.

# References

1. Balažovjech, M., Mikula, K.: A higher order scheme for a tangentially stabilized plane curve shortening flow with a driving force. SIAM J. Sci. Comp. **33**, 2277–2294 (2011)
2. Balažovjech, M., Mikula, K., Petrášová, M., Urbán, J.: Lagrangian methods with topological changes for numerical modelling of forest fire propagation. In: Handlovičova, A. et al. (eds.) Proceeding of Algoritmy 2012, pp. 42–52. Slovak University of Technology, Bratislava (2012)
3. Frolkovič, P., Mikula, K.: Flux-based level set method: a finite volume method for evolving interfaces. Appl. Numer. Math. **57**(4), 436–454 (2007)
4. Frolkovič, P., Mikula, K.: High-resolution flux-based level set method. SIAM J. Sci. Comp. **29**(2), 579–597 (2007)
5. Handlovičová, A., Mikula, K., Sgallari, F.: Semi-implicit complementary volume scheme for solving level set like equations in image processing and curve evolution. Numer. Math. **93**, 675–695 (2003)
6. Mikula, K., Ohlberger, M.: A new level set method for motion in normal direction based on a semi-implicit forward-backward diffusion approach. SIAM J. Sci. Comp. **32**(3), 1527–1544 (2010)
7. Mikula, K., Ohlberger, M.: Inflow-Implicit/Outflow-Explicit scheme for solving advection equations. In: Fort, J. et al. (ed.) Finite Volumes for Complex Applications VI, pp. 683–692. Springer, New York (2011)
8. Walkington, N.J.: Algorithms for computing motion by mean curvature. SIAM J. Numer. Anal. **33**, 2215–2238 (1996)

# Adaptive Time Discretization and Linearization Based on a Posteriori Estimates for the Richards Equation

**Vincent Baron, Yves Coudière and Pierre Sochala**

**Abstract** We derive some a posteriori error estimates for the Richards equation, based on the dual norm of the residual. This equation is nonlinear in space and in time, thus its resolution requires fixed-point iterations within each time step. We propose a strategy to decrease the computational cost relying on a splitting of the error terms in three parts: linearization, time discretization, and space discretization. In practice, we stop the fixed-point iterations after the linearization error becomes negligible, and choose the time step in order to balance the time and space errors.

## 1 Introduction

We focus on water infiltration modeled by the parabolic nonlinear Richards equation, written here on a polygonal domain $\Omega$ (in $\mathbb{R}^2$) with a finite time horizon $T > 0$ and mixed Dirichlet-Neumann boundary conditions:

$$\begin{cases} \partial_t \theta(\psi) - \nabla \cdot (\mathbb{K}(\psi)\nabla(\psi + z)) = f & \text{in } Q_T := \Omega \times (0, T), \\ \psi = \psi_{\mathrm{D}} & \text{on } \partial\Omega^{\mathrm{D}} \times (0, T), \\ -\mathbb{K}(\psi)\nabla(\psi + z) \cdot n = g & \text{on } \partial\Omega^{\mathrm{N}} \times (0, T), \\ \psi_{t=0} = \psi^0 & \text{in } \Omega \times \{0\}, \end{cases} \tag{1}$$

V. Baron (✉) · P. Sochala
BRGM, 3 avenue Claude Guillemin, 45060 Orléans, France
e-mail: vincent.baron@univ-nantes.fr

P. Sochala
e-mail: p.sochala@brgm.fr

Y. Coudière
IMB, 351 cours de la libération, 33405 Talence cedex, France
e-mail: yves.coudiere@inria.fr

where the unknown is the hydraulic head $\psi$. Implicit schemes are preferred to solve this equation because explicit ones are only valid in the vadose zone and have a restrictive CFL condition. Therefore, nonlinear systems are solved at each discrete time level, and we have to pay special attention to the computational cost. In order to optimize this cost, we can use local a posteriori estimates, which allow to control the error between the exact and approximate solutions. These estimates are local bounds that involve only the approximate solution and, ideally, some fully computable constants. Several methods are available to obtain such estimates. The current work deals with the *equilibrated fluxes method* [13]. For this method, the difficulty is to reconstruct some continuous equilibrated fluxes from the discrete ones.

This method has received particular attention in various studies over the last few years: finite elements for elasticity problems and the Poisson equation in [5, 11], DG methods for a reaction-diffusion-convection equation in [8], finite volumes for multiphase compositional flows in [6]. More recently, theoretical developments have unified various space discretizations, for the linear heat equation [9], and for a nonlinear parabolic problem [7]. In the latter robust fully computable lower and upper bounds were obtained using a space-time dual norm. This work includes the water content formulation of the Richards equation. Other related results concerning the Richards equation are available in [3], but based on the Kirchoff transform.

In this paper we consider a formulation based on the hydraulic head, which remains valid in saturated soils unlike the water content form. We also put aside the Kirchoff transform to benefit from conservative schemes especially designed for physical variables. The objective of this work is twofold: to derive a fully computable upper bound in the spirit of [7, 13], and to propose some space and time flux reconstructions for the Discrete Duality Finite Volume Scheme (DDFV) from [2]. Indeed, the nonlinearity in the time derivative term requires a special treatment, which we address by equilibrating the time flux as well as the space flux in accordance with the space-time norm used in the estimators. The upper bound can be split into spatial, temporal and linearization error components. An adaptive algorithm is proposed to choose a stopping criterion for the nonlinear algorithm and to adjust the time step during the simulation. Although we consider a DDFV scheme in this paper, our results remain general and not attached to a particular space discretization.

Section 2 introduces our upper bounds using a key space and time equilibrated flux assumption. Section 3 describes the appropriate flux reconstructions for the DDFV scheme. Section 4 presents an adaptive algorithm, and numerical results for two test cases: an infiltration problem with an analytical solution, and a stiff case. Section 5 draws some conclusions and perspectives.

## 2 A Posteriori Error Estimate

Consider a simplicial mesh $\mathcal{M}$ of $\Omega$. We denote by $K$ a cell of $\mathcal{M}$ of diameter $h_K$ and by $(t^n)_{1 \le n \le N}$ some discrete time levels. We set $I^n := (t^{n-1}, t^n)$ and $\delta t^n := t^n - t^{n-1}$ for $1 \le n \le N$. The weak formulation of (1) reads:

$$\int_0^T \{(f, \phi) + (\theta(\psi), \partial_t \phi) - (\mathbb{K}(\psi)\nabla(\psi + z), \nabla\phi) - (g, \phi)_{\partial\Omega}\}(t) \, dt$$

$$+ (\theta(\psi^0), \phi(., 0)) = 0 \quad (2)$$

where $(., .)$ denotes the usual $L^2$-inner product on $\Omega$, and the test functions $\phi$ belong to the space $Y := \{\phi \in L^2(0, T; H^1(\Omega)) \mid \partial_t \phi \in L^2(Q_T), \phi(., T) = 0, \phi(x, t) = 0 \, \forall(x, t) \in \partial\Omega^D \times ]0, T]\}$. Under reasonable assumptions expressed in [1], one can prove the existence of a solution $\psi$ to Eq. (2) in $L^2(0, T; H^1(\Omega)) \cap L^\infty(Q_T)$.

We assume that an approximate solution $\tilde{\psi}$ is available in the space $X := \{\phi \in L^2(0, T, H^1(\Omega)) \mid \partial_t \phi \in L^2(Q_T)\}$. The usual $L^2$-norm on $K \times I^n$ is denoted by $\|\cdot\|_{K \times I^n}$, we set $\|\phi\|_{Y, K \times I^n} := \left(\|\phi^2\|_{K \times I^n} + h_K^2 \|\nabla\phi^2\|_{K \times I^n} + (\delta t^n)^2 \|\partial_t \phi\|_{K \times I^n}^2\right)^{1/2}$ and $\|\phi\|_Y := \left(\sum_n \sum_{K \in \mathcal{M}} \|\phi^2\|_{Y, K \times I^n}\right)^{1/2}$. As we want local estimates, we define at each time level the subspace $Y^n := \{\phi \in Y \mid t \mapsto \phi(\cdot, t) \text{ vanishes outside } I_n\}$. The error we want to measure is defined as $\mathscr{E}^n(\tilde{\psi}) := \sup_{\phi \in Y^n, \|\phi\|_Y = 1} \langle R^n(\tilde{\psi}), \phi\rangle$, where $\langle R^n(\tilde{\psi}), \phi\rangle := \int_{I^n} \{(f, \phi) + (\theta(\tilde{\psi}), \partial_t \phi) - (\mathbb{K}(\tilde{\psi})\nabla(\tilde{\psi} + z), \nabla\phi) - (g, \phi)_{\partial\Omega}\}(t) \, dt$.

We now consider any space discretization of (1) that can be written as:

$$\frac{d}{dt} M\Theta(\Psi_h) + A(\Psi_h)\Psi_h = B(\Psi_h) \quad (3)$$

where $\Psi_h(t)$ is the vector of the degrees of freedom for the approximate solution, $M$ is a nonsingular mass matrix, the matrix $A(\Psi_h)$ approximates the semilinear diffusion term $-\nabla \cdot (\mathbb{K}(\cdot)\nabla\cdot)$, and the vector $B(\Psi_h)$ gathers the contributions of the gravity, source and boundary conditions. Then we use the Crank-Nicolson time-stepping algorithm, and obtain the following discrete problem for each $n \geq 1$:

$$\frac{M}{\delta t^n}\Theta(\Psi^n) + \frac{1}{2}A(\Psi^n)\Psi^n$$

$$= \frac{1}{2}\left(B(\Psi^n) + B(\Psi^{n-1})\right) + \frac{M}{\delta t^n}\Theta(\Psi^{n-1}) - \frac{1}{2}A(\Psi^{n-1})\Psi^{n-1} =: D(\Psi^n, \Psi^{n-1}),$$

where $\Psi^n \simeq \Psi_h(t^n)$. We solve this nonlinear system using linearizations like in [12]. For each time level $n$, we look for a sequence of vectors $(\Psi^{n,m})_{m\geq 1}$ solving

$$\left(\frac{M}{\delta t^n}\Theta'(\Psi^{n,m-1}) + \frac{1}{2}A(\Psi^{n,m-1})\right)\delta\Psi^{n,m}$$

$$= D(\Psi^{n,m-1}, \Psi^{n-1}) - \frac{1}{2}A(\Psi^{n,m-1})\Psi^{n,m-1} - \frac{M}{\delta t^n}\Theta(\Psi^{n,m-1}),$$

$$(4)$$

where $\delta\Psi^{n,m} = \Psi^{n,m} - \Psi^{n,m-1}$. In order to obtain this discrete system, we used the approximations $\Theta(\Psi^{n,m}) \simeq \Theta(\Psi^{n,m-1}) + \Theta'(\Psi^{n,m-1}) \cdot \delta\Psi^{n,m}$, $A(\Psi^{n,m})\Psi^{n,m} \simeq A(\Psi^{n,m-1})\Psi^{n,m}$, and $D(\Psi^{n,m}, \Psi^{n-1}) \simeq D(\Psi^{n,m-1}, \Psi^{n-1})$. This corresponds to a Newton linearization in time, and a Picard linearization in space.

We assume that we can associate with each discrete vector $\Psi^\bullet$ ($\bullet$ stands for $\{n, m\}$ or $n - 1$) of degrees of freedom a unique reconstructed function $\psi_h^\bullet \in H^1(\Omega)$. At time $t^n$ the value of $\Psi^{n-1}$ is known, and so is its reconstruction $\psi_h^{n-1} \in H^1(\Omega)$. After resolution of the linear system (4) for $1 \leq m$, the values of $\Psi^{n,m}$ are known, and we consider their reconstructions $\psi_h^{n,m} \in H^1(\Omega)$. Next, for each iteration level $m$, the space and time reconstruction is given by $\tilde{\psi}^{n,m}(t) = (t - t^{n-1})/\delta t^n \, \psi_h^{n,m} + (t^n - t)/\delta t^n \, \psi_h^{n-1} \in X$. The following theorem gives an a posteriori estimate of the error $\mathscr{E}^n(\tilde{\psi}^{n,m})$. In the resolution algorithm, the nonlinear iterations are stopped upon completion of a criterion based on this theorem, and we set $\Psi^n := \Psi^{n,m}$ before moving to the next time level.

**Theorem 1** *Let $n \geq 1$ be given and consider some reconstructions on $\Omega \times I^n$ of: the source term, $\tilde{f} \in L^2(\Omega \times I_n)$; the time flux at the current iteration $m$, $\tilde{\theta}^m(t) \in L^2(\Omega)$ and affine in time; the space flux at the current iteration $m$, $\tilde{\mathbf{t}}^m \in L^2(I^n, H(\mathrm{div}, \Omega))$; and the time and space fluxes obtained by linearization between the iterations $m - 1$ and $m$, $\tilde{\theta}_{lin}^m(t) \in L^2(\Omega)$ (affine in time) and $\tilde{\mathbf{t}}_{lin}^m \in L^2(I^n, H(\mathrm{div}, \Omega))$ and such that $\tilde{\theta}_{lin}^m(t^{n-1}) = \tilde{\theta}^m(t^{n-1})$, $\tilde{\mathbf{t}}_{lin}^m(t^{n-1}) = \tilde{\mathbf{t}}^m(t^{n-1})$, both independent of $m$ at time $t^{n-1}$.*
*Assume, that for any cell $K \in \mathscr{M}$, we have*

$$\int_{K \times I^n} \{\tilde{f} - \partial_t \tilde{\theta}_{lin}^m - \nabla \cdot \tilde{\mathbf{t}}_{lin}^m\} \, \mathrm{d}t = 0 \tag{5}$$

*(note that $\partial_t \tilde{\theta}_{lin}^m = 1/\delta t^n \left( \tilde{\theta}_{lin}^m(t^n) - \tilde{\theta}^m(t^{n-1}) \right)$). Then the following estimate holds:*

$$\mathscr{E}^n(\tilde{\psi}^{n,m}) \leq \eta_{residual}^{n,m} + \eta_{\tilde{\theta}}^{n,m} + \eta_{\tilde{\theta}_{lin}}^{n,m} + \eta_{\tilde{\mathbf{t}}}^{n,m} + \eta_{\tilde{\mathbf{t}}_{lin}}^{n,m} + \eta_{\tilde{f}}^{n} + \eta_{boundary}^{n,m},$$

*with $\eta_\bullet^{n,m} := \left( \sum_{K \in \mathscr{M}} (\eta_{\bullet, K}^{n,m})^2 \right)^{1/2}$ and $\eta_\bullet^{n,m} := \left( \sum_{\sigma \subset \partial \Omega} (\eta_{\bullet, \sigma}^{n,m})^2 \right)^{1/2}$, where*

$$\eta_{residual, K}^{n,m} := C^P \left\| \tilde{f} - \partial_t \tilde{\theta}_{lin}^m - \nabla \cdot \tilde{\mathbf{t}}_{lin}^m \right\|_{K \times I^n},$$

$$\eta_{\tilde{\theta}, K}^{n,m} := (\delta t^n)^{-1} \left\| \theta(\tilde{\psi}^{n,m}) - \tilde{\theta}^m \right\|_{K \times I^n}, \quad \eta_{\tilde{\theta}_{lin}, K}^{n,m} := (\delta t^n)^{-1} \left\| \tilde{\theta}_{lin}^m - \tilde{\theta}^m \right\|_{K \times I^n},$$

$$\eta_{\tilde{\mathbf{t}}, K}^{n,m} := h_K^{-1} \left\| \mathbb{K}(\tilde{\psi}^{n,m}) \nabla(\tilde{\psi}^{n,m} + z) + \tilde{\mathbf{t}}^m \right\|_{K \times I^n}, \quad \eta_{\tilde{\mathbf{t}}_{lin}, K}^{n,m} := h_K^{-1} \left\| \tilde{\mathbf{t}}^m - \tilde{\mathbf{t}}_{lin}^m \right\|_{K \times I^n},$$

$$\eta_{\tilde{f}, K}^{n} := \left\| f - \tilde{f} \right\|_{K \times I^n}, \quad \eta_{boundary, \sigma}^{n,m} := \left( \sqrt{2} C_{K, \sigma}^T \right)^{\frac{1}{2}} h_K^{-1/2} \left\| \tilde{\mathbf{t}}^m - g \right\|_{\sigma \times I^n}.$$

The proof is derived from [7]. It is based on Eq. (5), which is a re-expression of the scheme, along with some space and time Poincaré and trace inequalities from [7, 8].

**Fig. 1** **a** Notations for DDFV method—**b** Cells (in *gray*) of the secondary mesh, obtained by joining the center of each cell to the corresponding edges

**Fig. 2** Degrees of freedom for the reconstruction of $\psi_h^\bullet$ (*circle points*) and $\theta_h^\bullet$ (*square points*) on each triangle $D_{\sigma,K}$. Both functions are continuous across interfaces



## 3 Application to the DDFV Scheme

We discretize equation (1) with the DDFV scheme detailed in [2]. This method provides some approximate values $\Psi_h = (\psi_K, \psi_A)$ of $\psi$ at the centers $x_K$ and at the vertices $x_A$ of the cells $K \in \mathcal{M}$. The unknown vector $\Psi_h$ solves a system of finite volume equations (written here for $f = 0$ and a homogeneous Dirichlet problem):

$$|K|\frac{d\theta(\psi_K)}{dt} - \sum_{\sigma \in \partial K} \mathbb{K}(\psi_{\sigma,K})\left(\nabla_{\sigma,K}\Psi_h + e_z\right) \cdot N_{\sigma,K} = 0,$$

$$|A|\frac{d\theta(\psi_A)}{dt} - \sum_{\sigma \in \partial A} \left(\mathbb{K}(\psi_{\sigma,K})(\nabla_{\sigma,K}\Psi_h + e_z) \cdot N_K^A + \mathbb{K}(\psi_{\sigma,L})(\nabla_{\sigma,L}\Psi_h + e_z) \cdot N_L^A\right) = 0,$$

for all cells $K$ in $\mathcal{M}$ and all vertices $x_A$ of $\mathcal{M}$. The permeability matrix $\mathbb{K}$ is evaluated at the points $\psi_{\sigma,K} := (\psi_K + \psi_A + \psi_B)/3$ (Fig. 1). We construct two gradients $\nabla_{\sigma,K}$ and $\nabla_{\sigma,L}$ on each side of an interface $\sigma$ between some cells $K$ and $L$ (see Fig. 2):

$$\nabla_{\sigma,K}\Psi_h = \frac{1}{2|D_{\sigma,K}|}\left((\psi_\sigma - \psi_K)N_{\sigma,K} + (\psi_B - \psi_A)N_K^A\right),$$

$$\nabla_{\sigma,L}\Psi_h = \frac{1}{2|D_{\sigma,L}|}\left((\psi_L - \psi_\sigma)N_{\sigma,L} + (\psi_B - \psi_A)N_L^A\right).$$

We get rid of the auxiliary unknown $\psi_\sigma$ by solving the linear local conservativity condition $\mathbb{K}(\psi_{\sigma,K})\left(\nabla_{\sigma,K}\Psi_h + e_z\right) \cdot N_{\sigma,K} + \mathbb{K}(\psi_{\sigma,L})\left(\nabla_{\sigma,L}\Psi_h + e_z\right) \cdot N_{\sigma,L} = 0$.

Given $1 \leq n$ and the vector $\Psi^{n-1}$ at time $t^{n-1}$, the fully discretized scheme provides the sequence $(\Psi^{n,m})_{m \geq 1}$ at time $t^n$. The function $\psi_h^{n-1} \in H^1(\Omega)$ at time $t^{n-1}$ is the piecewise affine function on the triangles $D_{\sigma,K}$ that interpolates $\psi_K^{n-1}$, $\psi_A^{n-1}$ and $\psi_B^{n-1}$ at the vertices $x_K$, $x_A$ and $x_B$ of $D_{\sigma,K}$. The functions $\psi_h^{n,m} \in H^1(\Omega)$ are constructed similarly from the vectors $\Psi_h^{n,m}$. It remains to explain how the reconstructions $\tilde{\theta}^m$, $\tilde{\theta}_{\mathrm{lin}}^m$, $\tilde{\mathbf{t}}^m$, $\tilde{\mathbf{t}}_{\mathrm{lin}}^m$ and $\tilde{f}$ used to obtain (5) are built. All these reconstructions are affine in time, hence completely determined by their values at time $t^{n-1}$ and at time $t^n$, nonlinear iterate number $m$, denoted respectively by $\theta_h^{n-1}$, $\theta_{h,\mathrm{lin}}^{n-1}$, $\mathbf{t}_h^{n-1}$, $\mathbf{t}_{h,\mathrm{lin}}^{n-1}$, $f_h^{n-1}$, and $\theta_h^{n,m}$, $\theta_{h,\mathrm{lin}}^{n,m}$, $\mathbf{t}_h^{n,m}$, $\mathbf{t}_{h,\mathrm{lin}}^{n,m}$, $f_h^n$.

The functions $f_h^{n-1}$ and $f_h^n$ in $L^2(\Omega)$ are piecewise constant on the cells $K$ with $(f_h^{n-1})_{|K} = f_K^{n-1}$ and $(f_h^n)_{|K} = f_K^n$, where $f_K^\bullet$ approximates $1/|K| \int_K f(x, t^\bullet) \, dx$.

The functions $\theta_h^\bullet \in L^2(\Omega)$ are piecewise polynomials of degree 3 on the triangles $D_{\sigma,K}$, uniquely defined by their values $\theta(\psi(x_i))$ at the 9 degrees of freedom depicted on Fig. 2 and the equality $1/|D_{\sigma,K}| \int_{D_{\sigma,K}} \theta_h^\bullet(x) \, dx = \theta(\psi_K^\bullet)$. The function $\theta_{h,\mathrm{lin}}^{n,m}$ is defined similarly, bar the condition $1/|D_{\sigma,K}| \int_{D_{\sigma,K}} \theta_{h,\mathrm{lin}}^{n,m}(x) \, dx = \theta(\psi_K^{n,m-1}) + \theta'(\psi_K^{n,m-1})(\psi_K^{n,m} - \psi_K^{n,m-1})$. In addition, we take $\theta_{h,\mathrm{lin}}^{n-1} = \theta_h^{n-1}$.

The functions $\mathbf{t}_h^\bullet \in H(\mathrm{div}, \Omega)$ are piecewise in the usual Raviart-Thomas-Nédélec space $\mathbf{RTN}_1(K)$ on each triangle $K$ and uniquely defined by the conditions (see [4]):

$$\frac{1}{|\sigma|} \int_\sigma \mathbf{t}_h^\bullet \cdot N_{\sigma,K} \, d\sigma = -\mathbb{K}(\psi_{\sigma,K}^\bullet) \left( \nabla_{\sigma,K} \Psi^\bullet + e_z \right) \cdot N_{\sigma,K},$$

$$\frac{1}{|\sigma|} \int_\sigma x \, \mathbf{t}_h^\bullet \cdot N_{\sigma,K} \, d\sigma = -\frac{1}{2}\mathbb{K}(\psi_{\sigma,K}^\bullet) \left( \nabla_{\sigma,K} \Psi^\bullet + e_z \right) \cdot N_{\sigma,K},$$

$$\int_K \mathbf{t}_h^\bullet \, dx = - \sum_{\sigma \subset \partial K} |D_{\sigma,K}| \mathbb{K}(\psi_{\sigma,K}^\bullet) \left( \nabla_{\sigma,K} \Psi^\bullet + e_z \right).$$

Note that the last condition is straightforward to derive with the DDFV scheme, as the discrete gradient is piecewise constant on each triangle $D_{\sigma,K}$. $\mathbf{t}_{h,\mathrm{lin}}^{n,m}$ is set similarly, using $\mathbb{K}(\psi_{\sigma,K}^{n,m-1})$ instead of $\mathbb{K}(\psi_{\sigma,K}^{n,m})$.

# 4 Results

In order to define the adaptive algorithm below, the total error estimate from Theorem 1 is split into the contribution of the time-stepping method $\eta_{\mathrm{time}}^{n,m} := \eta_{\mathrm{residual}}^{n,m}$, the contribution of the space discretization $\eta_{\mathrm{space}}^{n,m} := \eta_{\tilde{\theta}}^{n,m} + \eta_{\tilde{\mathbf{t}}}^{n,m} + \eta_{\tilde{f}}^n + \eta_{\mathrm{boundary}}^{n,m}$, and the contribution of the linearization $\eta_{\mathrm{lin}}^{n,m} := \eta_{\tilde{\theta}_{\mathrm{lin}}}^{n,m} + \eta_{\tilde{\mathbf{t}}_{\mathrm{lin}}}^{n,m}$. At time $t^n$, we start with the current time step, then stop the nonlinear iterations for $m = m^\star$ such that the estimate $\eta_{\mathrm{lin}}^{n,m^\star}$ becomes small with respect to $\eta_{\mathrm{space}}^{n,m^\star} + \eta_{\mathrm{time}}^{n,m^\star}$, and finally adjust the time step so as to balance $\eta_{\mathrm{time}}^{n,m^\star}$ and $\eta_{\mathrm{space}}^{n,m^\star}$.

The first test case features a downward infiltration problem with an analytical solution given by $\psi(z, t) = 20.4 \tanh (0.5 (z + t/12 - 15)) - 41.1$, which determines adequate source term and Dirichlet boundary condition enforced on $\partial\Omega$. The

**Table 1** L$^2$ norm of the error $\|\psi - \psi_h\|_{L^2(\Omega \times (0,T))}$

| $N_{triangles}$ | $\delta t^1$ | Error | Order |
|---|---|---|---|
| 118 | 4 | 3.02e-03 | |
| 430 | 2 | 9.16e-04 | 1.88 |
| 1688 | 1 | 2.27e-04 | 2.24 |
| 6474 | 0.5 | 5.75e-05 | 1.99 |

Both space and time diameters are halved at each refinement. The parameters are $\gamma_{\text{lin}} = 1/100$ and $\gamma_{\text{time}} = 2$



**Fig. 3** **a** Adapted time step over time (in seconds)—**b** Cumulated number of Picard iterations, adapted case (*dashed line*) versus non adapted case (*solid line*)

---

**Algorithm 1** Adaptive algorithm (Crank-Nicolson) at time level $t^n$

---

**Require:** $t^{n-1}, \delta t, \Psi^{n-1}$
1: **repeat**
2:     $m \leftarrow 0$, initialize $\Psi^n$
3:     **repeat**
4:         $\delta\Psi^{n,m} \leftarrow$ solution to the linearized Richards equation (4) with the time step $\delta t$
5:         $m \leftarrow m+1$, $\Psi^{n,m} \leftarrow \Psi^{n,m} + \delta\Psi^{n,m}$
6:         $\theta_h^{n,m}, \theta_{h,\text{lin}}^{n,m}, \mathbf{t}_h^{n,m}, \mathbf{t}_{h,\text{lin}}^{n,m} \leftarrow$ reconstructions computed as explained in section 3
7:         $\eta_{\text{lin}}^{n,m}, \eta_{\text{time}}^{n,m}, \eta_{\text{space}}^{n,m} \leftarrow$ the estimators from theorem 1
8:     **until** $\eta_{\text{lin}}^{n,m} < \gamma_{\text{lin}}(\eta_{\text{time}}^{n,m} + \eta_{\text{space}}^{n,m})$
9:     $\delta t^n \leftarrow \delta t$
10:     $\delta t \leftarrow \eta_{\text{space}}^{n,m}/\eta_{\text{time}}^{n,m} * \delta t$
11: **until** $\eta_{\text{time}}^{n,m} < \gamma_{\text{time}}\eta_{\text{space}}^{n,m}$
12: $t^n \leftarrow t^{n-1} + \delta t^n$

---

water content and the hydraulic conductivity are defined by Haverkamp's constitutive relationships from [10]. Table 1 shows that the expected second-order convergence remains valid when the adaptive algorithm is used.

We then propose a stiffer case characterized by a strong overpressure imposed on the top of the column through a Dirichlet condition. This stiffness requires a fairly low time step at the beginning of the simulation for the iterative procedure to converge. But the time step can be increased afterwards, and a time adaptation is especially relevant in this case. Figure 3 displays the gain measured in terms of the cumulated number of Picard iterations over time.

## 5  Conclusions

We presented a fully computable global in space and local in time upper bound for the Richards equation, using the Crank-Nicolson time scheme and any space discretization. The estimate was decoupled into three error components, which we equilibrated in an adaptive time-stepping algorithm. Our results showed that the second-order convergence still holds true, and that the benefit in terms of total number of iterations can reach a full order of magnitude on stiff cases. In a future work we plan to derive a lower bound to confirm the robustness of this estimate.

## References

1. Alt Wilhelm, H., Luckhaus, S.: Quasilinear elliptic-parabolic differential equations. Math. Z. **183**(3), 311–341 (1983)
2. Baron, V., Coudière, Y., Sochala, P.: Comparison of DDFV and DG methods for flow in anisotropic heterogeneous porous media. Oil Gas Sci. Technol. Rev. IFP En. nouvelles (2013)
3. Bernardi, C., El Alaoui, L., Mghazli, Z.: A posteriori analysis of a space and time discretization of a nonlinear model for the flow in variably saturated porous media. IMA J. Numer. Anal. (2013)
4. Brezzi, F., Fortin, M.: Mixed and Hybrid Finite Element Methods. Springer, New York (1991)
5. Destuynder, P., Métivet, B.: Explicit error bounds in a conforming finite element method. Math. Comput. Am. Math. Soc. **68**(228), 1379–1396 (1999)
6. Di Pietro, D.A., Vohralík, M., Yousef, S., et al.: A posteriori error estimates with application of adaptive mesh refinement for thermal multiphase compositional flows in porous media. Comput. Math. Appl. (2013)
7. Dolejší, V., Ern, A., Vohralík, M.: A framework for robust a posteriori error control in unsteady nonlinear advection-diffusion problems. SIAM J. Numer. Anal. **51**(2), 773–793 (2013)
8. Ern, A., Stephansen, A.F., Vohralík, M.: Guaranteed and robust discontinuous Galerkin a posteriori error estimates for convection-diffusion-reaction problems. J. Comput. Appl. Math. **234**(1), 114–130 (2010)
9. Ern, A., Vohralík, M.: A posteriori error estimation based on potential and flux reconstruction for the heat equation. SIAM J. Numer. Anal. **48**(1), 198–223 (2010)
10. Haverkamp, R., Vauclin, M., Touma, J., Wierenga, P., Vachaud, G.: A comparison of numerical simulation models for one-dimensional infiltration. Soil Sci. Soc. America J. **41**(2), 285–294 (1977)
11. Ladevèze, P.: Comparaison de modèles de milieux continus. Ph.D. thesis (1975)
12. Manzini, G., Ferraris, S.: Mass-conservative finite volume methods on 2-D unstructured grids for the Richards equation. Adv. Water Res. **27**(12), 1199–1215 (2004)
13. Prager, W., Synge, J.L.: Approximations in elasticity based on the concept of function space. Q. Appl. Math. **5**(3), 1–21 (1947)

# Monotone Combined Finite Volume-Finite Element Scheme for a Bone Healing Model

**Marianne Bessemoulin-Chatard and Mazen Saad**

**Abstract**    We define a combined edge FV-FE scheme for a bone healing model. This choice of discretization allows to take into account anisotropic diffusions and does not impose any restrictions on the mesh. Moreover, following [3], we propose a nonlinear correction to obtain a monotone scheme. We present some numerical experiments which show its good behavior.

## 1 Introduction

We consider a bone growth model based on [1]. It describes the evolution of the concentrations of the following quantities: the mesenchymal stem cells $s$, the osteoblasts $b$, the bone matrix $m$ and the osteogenic growth factor $g$. Bone healing begins by the migration of the stem cells to the site of the injury. Then along the bone, these cells differentiate into osteoblasts which start to synthetize the bone matrix. This cell differentiation is only possible in presence of the growth factor.

The proposed model takes into account several phenomena: the diffusion of the stem cells and the growth factor, the migration of the stem cells towards the bone matrix, the proliferation and the differentiation of stem cells. The osteoblasts are considered without movement since they are fixed at the bone matrix. Moreover, the model includes the case of heterogeneous domains, with possibly anisotropic diffusions.

M. Bessemoulin-Chatard (✉)
Université de Nantes, LMJL-UMR6629, 2 rue de la Houssinière, BP 92208, 44322
Nantes Cedex 3, France
e-mail: Marianne.Bessemoulin@univ-nantes.fr

M. Saad
Ecole Centrale de Nantes, LMJL-UMR6629,
1 rue de la Noé, BP 92101, 44321 Nantes Cedex 3, France
e-mail: Mazen.Saad@ec-nantes.fr

It is given by the following nonlinear coupled system: for $t > 0$ and $x \in \Omega$, where $\Omega$ is an open bounded polyhedral subset of $\mathbb{R}^d$, $d = 2, 3$,

$$\partial_t s - \operatorname{div}\left(\mathbf{S}(x)\left(\Lambda(m)\nabla s - V(m)\chi(s)\nabla m\right)\right) = K_1(m)\chi(s) - H(g)s, \quad (1)$$

$$\partial_t b = K_2(m)\chi(b) + \rho H(g)s - \delta_1 b, \quad (2)$$

$$\partial_t m = \lambda(1 - m)b, \quad (3)$$

$$\partial_t g - \operatorname{div}\left(\mathbf{S}(x)\Lambda_g \nabla g\right) = P(g)b - \delta_2 g. \quad (4)$$

The functions $K_1(m)$, $K_2(m)$, $H(g)$, $P(m)$ and the positive parameters $\rho, \delta_1, \lambda$ and $\delta_2$ are given (see [1]). The diffusion coefficient $\Lambda(m)$ and the haptotaxis velocity $V(m)$ are given by

$$\Lambda(m) = \frac{\chi_h}{\zeta_h^2 + m^2}(m + \Lambda_0)(1 - m), \quad V(m) = \frac{\chi_k}{(\zeta_k + m)^2},$$

with $\chi_h, \zeta_h, \Lambda_0, \chi_k, \zeta_k > 0$. The diffusion coefficient $\Lambda_g$ for the growth factor is a positive constant. Moreover, the accumulation of stem cells is limited by the factor $\chi(s) = s(1 - s)$. The permeability $\mathbf{S}(x)$ is a symmetric $d \times d$ matrix, with $\mathbf{S} \in L^\infty(\Omega)$, and we assume that $\exists C_S > 0$ such that $\forall x \in \Omega, \forall \xi \in \mathbb{R}^d, \mathbf{S}(x)\xi \cdot \xi \geq C_S|\xi|^2$.

This nonlinear system (1–4) is supplemented with initial conditions $s_0$, $b_0$, $m_0$, $g_0$ and with homogeneous Neumann boundary conditions on $s$ and $g$:

$$\mathbf{S}(x)\left(\Lambda_1(m)\nabla s - V(m)\chi(s)\nabla m\right) \cdot \mathbf{n} = 0, \quad \mathbf{S}(x)\Lambda_g \nabla g \cdot \mathbf{n} = 0, \quad (5)$$

for $t \in (0, T)$ and $x \in \partial\Omega$, where $\mathbf{n}$ is the outward unit normal of $\partial\Omega$. Following [5], a solution $u = (s, b, m, g)$ is said to be physically admissible if $u \in \mathscr{A} = [0, 1] \times [0, \bar{b}] \times [0, 1] \times [0, \bar{g}]$, where $\bar{b}$ and $\bar{g}$ depend on the physical parameters.

In this paper, we propose a numerical scheme for this bone growth model. A finite volume (FV) scheme was proposed in [5] for this model in homogeneous domains where the diffusion tensor $\mathbf{S} = Id$. The cell-centered FV method with an upwind discretization of the convective terms provides the stability and is extremely robust. However in this case, the mesh is assumed to be admissible [7, Definition 9.1]. In particular, this implies that the orthogonality condition has to be satisfied. As mentioned in [5], a difficulty in the implementation is to construct such admissible meshes. Structured rectangular meshes are admissible, but they cannot be used for complex geometries arising in physical contexts. Furthermore, the finite element (FE) method allows for an easy discretization of diffusive terms with full tensors without imposing any restrictions on the meshes. However, some numerical instabilities may arise in the convection-dominated case.

The idea is hence to combine a FE discretization of diffusive terms with a FV discretization of the other terms. Such schemes were proposed and studied in [9]

**Fig. 1** Triangles $K$, $L$ and $M \in \mathscr{T}$ and diamonds $D$, $E \in \mathscr{D}$ associated with edges $\sigma_D, \sigma_E \in \mathscr{E}$



for fluid mechanics equations in the case of diffusion terms with $\mathbf{S} = Id$ and in [4] for anisotropic Keller–Segel model. This idea was extended in [8] to inhomogeneous and anisotropic diffusion-dispersion tensors and to very general meshes only satisfying the shape regularity condition (6). However, the maximum principle is no more guaranteed if there exist negative transmissibilities.

We first introduce in Sect. 2 the combined FV-FE scheme for the bone healing model (1–4). Then in Sect. 3 we apply the method described in [3] to construct a nonlinear correction providing a discrete maximum principle. Finally in Sect. 4 we present some numerical experiments showing the efficiency of the scheme.

## 2 The Combined FV-FE Scheme

A mesh of $\Omega$ is a family $\mathscr{T}$ of closed simplices $K$ such that $\overline{\Omega} = \cup_{K \in \mathscr{T}} K$. We denote by $\mathscr{E}$ the set of all edges, by $\mathscr{E}^{int}$ the set of interior edges, by $\mathscr{E}^{ext}$ the set of boundary edges and by $\mathscr{E}_K$ the set of all edges of $K \in \mathscr{T}$. The size of the mesh is defined by $h := \max \operatorname{diam}(K)$. We assume that there exists a positive constant $k_{\mathscr{T}}$ such that:

$$reg(\mathscr{T}) := \min_{K \in \mathscr{T}} \frac{|K|}{(\operatorname{diam}(K))^d} \geq k_{\mathscr{T}}. \tag{6}$$

We also use a dual partition $\mathscr{D}$ of control volumes $D$ of $\Omega$ called diamonds such that $\overline{\Omega} = \cup_{D \in \mathscr{D}} D$. Each diamond $D$ is associated with one edge $\sigma_D \in \mathscr{E}$. We construct it by connecting the barycenters of every $K \in \mathscr{T}$ that contains $\sigma_D$ through the vertices of $\sigma_D$ (Fig. 1). For $\sigma_D \in \mathscr{E}^{ext}$, the contour of $D$ is completed by the edge $\sigma_D$ itself. We define $\mathscr{D}^{int}$ and $\mathscr{D}^{ext}$ the set of all interior and boundary dual volumes respectively. For $K \in \mathscr{T}$, we set $\mathscr{D}_K := \{D \in \mathscr{D}; \sigma_D \in \mathscr{E}_K\}$. We denote by $|D|$ the $d$-dimensional Lebesgue measure of $D$ and $|\sigma|$ the $(d-1)$-dimensional measure of $\sigma$. For all $D \in \mathscr{D}$, $P_D$ is the barycenter of $\sigma_D$ and $\mathscr{N}(D)$ is the set of neighbours of $D$. For all $D \in \mathscr{D}$ and all $E \in \mathscr{N}(D)$, $\sigma_{D,E}$ is the interface between $D$ and $E$ and $\mathbf{n}_{D,E}$ is the unit normal vector to $\sigma_{D,E}$ outward to $D$.

Next we define the following finite-dimensional space of piecewise linear non-conforming FE [6]:
$X := \{\varphi \in L^2(\Omega); \varphi|_K \text{ linear } \forall K \in \mathscr{T}, \varphi \text{ continuous at the points } P_D, D \in \mathscr{D}^{int}\}$, equipped with the seminorm $\|u\|_X^2 := \sum_{K \in \mathscr{T}} \int_K |\nabla u|^2 dx.$

The basis of $X$ is spanned by the shape functions $\varphi_D$, $D \in \mathscr{D}$, such that $\varphi_D(P_E) = \delta_{DE}$, $E \in \mathscr{D}$. The approximations in this space are nonconforming since $X \not\subseteq H^1(\Omega)$.

Finally, we define the time step $\Delta t$ and the increasing sequence $(t^n)_{0 \leq n \leq N+1}$, where $t^n = n\Delta t$ and $N$ is the smallest integer such that $(N+1)\Delta t \geq T$.
The discrete unknowns are denoted by $\{w_D^n,\ D \in \mathscr{D},\ n \in \{0 \cdots N+1\}\}$, where the value $w_D^n$ is an approximation of $w(P_D, t^n)$, $w = s,\ b,\ m,\ g$.

We now define the semi-implicit in time and combined FV-FE in space discretization for (1)–(4). The initial conditions are approximated by $(s_D^0, b_D^0, m_D^0, g_D^0)_{D \in \mathscr{D}}$ by taking the mean values of $s_0, b_0, m_0$ and $g_0$ on each dual cell $D$. Then the scheme is given by the following set of equations: for all $n \in \{0, ..., N\}$ and all $D \in \mathscr{D}^{int}$,

$$|D|\left(s_D^{n+1} - s_D^n\right) - \Delta t \sum_{E \in \mathscr{D}^{int}} \Lambda_{D,E}^n s_E^{n+1} + \Delta t \sum_{E \in \mathscr{N}(D)} G\left(s_D^{n+1}, s_E^{n+1}, V_{D,E}^n\left(m_E^{n+1} - m_D^{n+1}\right)\right)$$

$$= \Delta t |D| \left(K_1(m_D^n)s_D^{n+1}(1 - s_D^n) - H(g_D^n)s_D^{n+1}\right), \tag{7}$$

$$|D|\left(b_D^{n+1} - b_D^n\right) = \Delta t |D| \left(K_2(m_D^n)\chi(b_D^{n+1}) + H(g_D^n)s_D^{n+1} - \delta_1 b_D^{n+1}\right), \tag{8}$$

$$|D|\left(m_D^{n+1} - m_D^n\right) = \Delta t |D| \lambda (1 - m_D^{n+1})b_D^{n+1}, \tag{9}$$

$$|D|\left(g_D^{n+1} - g_D^n\right) - \Delta t \sum_{E \in \mathscr{D}^{int}} \mathscr{S}_{D,E} \Lambda_g g_E^{n+1} = \Delta t |D| \left(P(g_D^n)b_D^n - \delta_2 g_D^{n+1}\right), \tag{10}$$

where for $U = \Lambda,\ V$,

$$U_{D,E} = -\sum_{K \in \mathscr{T}} U_K \left(\mathbf{S}(x)\nabla\varphi_E, \nabla\varphi_D\right)_{0,K}, \quad \mathscr{S}_{D,E} = -\sum_{K \in \mathscr{T}} \left(\mathbf{S}(x)\nabla\varphi_E, \nabla\varphi_D\right)_{0,K},$$

with $U_K = \dfrac{\sum_{D \in \mathscr{D}_K} U(m_D)}{\text{card}(\mathscr{E}_K)}$.

The flux function $G$ is supposed to be monotone, consistent, conservative and locally Lipschitz continuous. For example, we consider in the following

$$G(a, b, c) = c^+ \left(\chi_\uparrow(a) + \chi_\downarrow(b)\right) - c^- \left(\chi_\uparrow(b) + \chi_\downarrow(a)\right),$$

where $c^+ = \max(c, 0)$, $c^- = \max(-c, 0)$, $\chi_\uparrow$ and $\chi_\downarrow$ are respectively the nondecreasing and nonincreasing parts of $\chi$.

**Definition 1** (Approximate solution) Using the values $(u_D)_{D \in \mathscr{D}}$, $u = s,\ b,\ m,\ g$, we define a nonconforming FE solution $u_h$ as a function piecewise linear and continuous in the barycenters $P_D$ of interior edges such that

$$u_h(x) = \sum_{D \in \mathscr{D}} u_D \varphi_D(x), \quad x \in \Omega.$$

**Properties of the discrete diffusive operators.** We define

$$\mathscr{A}^{\mathscr{D}} : \mathbb{R}^{Card(\mathscr{D})} \quad \to \mathbb{R}^{Card(\mathscr{D})} \qquad \mathscr{L}^{\mathscr{D}} : \mathbb{R}^{Card(\mathscr{D})} \quad \to \mathbb{R}^{Card(\mathscr{D})}$$
$$s_h = (s_D)_{D \in \mathscr{D}} \mapsto (A_D(s_h))_{D \in \mathscr{D}}, \qquad g_h = (g_D)_{D \in \mathscr{D}} \mapsto (L_D(g_h))_{D \in \mathscr{D}},$$

the discrete diffusive operators appearing in (7), (10), with for all $D \in \mathscr{D}$,

$$A_D(s_h) = \sum_{E \in \mathscr{D}^{int}} \Lambda_{D,E} s_E, \quad L_D(g_h) = \sum_{E \in \mathscr{D}^{int}} \Lambda_g \mathscr{S}_{D,E} g_E.$$

We give in the following proposition some properties of $\mathscr{A}^{\mathscr{D}}$ which are crucial to get the convergence of the scheme. The same results hold for $\mathscr{L}^{\mathscr{D}}$ too.

**Proposition 1** *The discrete diffusive operator $\mathscr{A}^{\mathscr{D}}$ is*

- **Conservative:** $\forall D \in \mathscr{D}, \quad A_D(s_h) = \sum_{E \in \mathscr{N}(D)} \Lambda_{D,E}(s_E - s_D),$
- **Coercive:** $\exists C_A > 0$ *such that* $- \sum_{D \in \mathscr{D}} A_D(s_h) s_D \geq C_A \|s_h\|_X^2 \quad \forall s_h \in X.$

## 3 Monotone Correction

At this stage, the constructed scheme is valid both for full anisotropic diffusion tensors and for general meshes satisfying only assumption (6). However, it possesses a discrete maximum principle only if all transmissibilities $\Lambda_{D,E}, \mathscr{S}_{D,E}$ are nonnegative, which is not guaranteed in the general case. Following [3], we now define a nonlinear correction which gives monotone scheme while preserving the properties described in Proposition 1.

We replace the operator $\mathscr{A}^{\mathscr{D}}$ in (7) by the corrected operator $\mathscr{B}^{\mathscr{D}}$ defined by

$$B_D(s) = A_D(s) + \sum_{E \in \mathscr{N}(D)} \beta_{D,E}^{\varepsilon}(s)(s_D - s_E) \quad \forall D \in \mathscr{D},$$

where $\beta_{D,E}(s)$ is the regularized correction proposed in [3]:

$$\beta_{D,E}^{\varepsilon}(s) = \max \left( \frac{|A_D(s)|}{Card_{\varepsilon} V(D, s)^*}, \frac{|A_E(s)|}{Card_{\varepsilon} V(E, s)^*} \right) \frac{1}{|s_D - s_E| + \varepsilon},$$

with $Card_{\varepsilon} V(D, s)^* = \sum_{E \in \mathscr{N}(D)} \frac{|s_D - s_E|}{|s_D - s_E| + \varepsilon}.$

This corrected diffusive operator is monotone since $\beta_{D,E}(s) > |A_D(s)|$ for all $D \in \mathscr{D}$, all $E \in \mathscr{N}(D)$. Moreover, the corrected diffusive operator $\mathscr{B}^{\mathscr{D}}$ still satisfies the properties described in Proposition 1:

- It is conservative, since $\beta_{D,E} = \beta_{E,D}$ for all $D \in \mathscr{D}$, $E \in \mathscr{N}(D)$,
- It is coercive, since $\beta_{D,E} \geq 0$ for all $D \in \mathscr{D}$, $E \in \mathscr{N}(D)$.

The diffusive operator $\mathscr{L}^{\mathscr{D}}$ can also be corrected in the same way.

**Theorem 1** *If $(s_D^0, b_D^0, m_D^0, g_D^0) \in \mathscr{A}$ for all $D \in \mathscr{D}$, then the discrete problem* (7)–(10) *with monotone correction has a physically admissible solution* $(s_D^n, b_D^n, m_D^n, g_D^n) \in \mathscr{A}$, *for all $n \geq 0$ and all $D \in \mathscr{D}$.*

The proof of this result can be done by introducing a truncated version of the scheme (see [5, Theorem 5]), and using the properties of the corrected diffusive operators (monotony, conservativity, continuity). Following the same lines as [5, Theorem 7], we can also prove some energy estimates:

**Theorem 2** *Let $(s_D^n, b_D^n, m_D^n, g_D^n)_{D \in \mathscr{D}, n \geq 0}$ be a solution of the corrected scheme. Then $\exists C > 0$ not depending on the discretization parameters such that*

$$\sum_{n=0}^{N-1} \Delta t \left( \|s_h^n\|_X^2 + \|b_h^n\|_X^2 + \|m_h^n\|_X^2 + \|g_h^n\|_X^2 \right) \leq C.$$

Starting from this result, one can obtain some compactness estimates on discrete solutions. The complete study of convergence of the corrected scheme, which requires some additional numerical assumptions [3], is done in [2].

## 4 Numerical Experiments

We simulate the healing of a long bone fracture in rats [10]. The simulation corresponds to a 0.07 cm fracture. To implement the semi-implicit scheme (7)–(10), we use the Newton's method coupled with a biconjugate gradient method to solve the nonlinear system. While the discrete maximum principle is not satisfied, the monotone correction is computed using the iterative algorithm described in [3]. The geometry of the fracture and the initial condition are described on Fig. 2.

We assume that $\mathbf{S}(x) = I_2$. We first consider an admissible mesh made of 14336 triangles and 21632 edges. Especially, all the angles are acute, which ensures in this case that the combined FV-FE scheme without correction satisfies the maximum principle. In particular, we observe that the discrete unknowns remain nonnegative (Table 1).

**Fig. 2** Geometry and initial condition: the *black* area corresponds to the bone matrix ($m_0 = 1$) and the *grey* area to the cellular cluster ($s_0 = 1$, $g_0 = 20$). Elsewhere there is nothing initially

**Table 1** Results obtained with the non corrected scheme on an admissible mesh

|          | Min. Val. $s$        | Max. Val. $s$ | Min. Val. $g$         | Max. Val. $g$ |
|----------|----------------------|---------------|-----------------------|---------------|
| Iter. 1  | $9.47 \times 10^{-21}$ | 0.999         | $9.9 \times 10^{-21}$  | 19.8          |
| Iter. 10 | $5.83 \times 10^{-21}$ | 0.991         | $9.05 \times 10^{-21}$ | 17.99         |

**Table 2** Numerical results with the original and the corrected schemes after 10 iterations

|                  |                | Mesh 1                  | Mesh 2                  | Mesh 3                  |
|------------------|----------------|-------------------------|-------------------------|-------------------------|
| Without          | Undershoots $s$ | 16                      | 16                      | 92                      |
| correction       | Min. Val. $s$   | $-2.67 \times 10^{-4}$  | $-8.95 \times 10^{-7}$  | $-3.01 \times 10^{-4}$  |
| after 10 it.     | Max. Val. $s$   | 0.990                   | 0.991                   | 0.992                   |
|                  | Undershoots $g$ | 70                      | 71                      | 144                     |
|                  | Min. Val. $g$   | $-0.27$                 | $-1.01 \times 10^{-2}$  | $-8.56 \times 10^{-3}$  |
|                  | Max. Val. $g$   | 17.96                   | 18.41                   | 20.27                   |
| With             | Min. Val. $s$   | $9.58 \times 10^{-6}$   | $1.18 \times 10^{-6}$   | $3.93 \times 10^{-6}$   |
| correction       | Max. Val. $s$   | 0.989                   | 0.990                   | 0.99                    |
| after 10 it.     | Min. Val. $g$   | $1.51 \times 10^{-4}$   | $8.02 \times 10^{-5}$   | $7.68 \times 10^{-5}$   |
|                  | Max. Val. $g$   | 17.81                   | 18.33                   | 19.37                   |

Then we consider three general unstructured meshes that contain obtuse angles. Mesh 1 is made of 1539 triangles and 2346 edges, mesh 2 is made of 3132 triangles and 4756 edges, and mesh 3 is made of 15568 triangles and 23479 edges. In Table 2, we present the minimum and maximum values obtained with the scheme before and after correction, after 1 and 10 iterations. We clearly observe that the discrete maximum principle is well respected after correction, with disappearance of the undershoots.

We now consider the corrected scheme on the finest mesh 3. After 2 days, we observe the formation of osteoblasts where the stem cells were initially concentrated

**Fig. 3** Bone matrix density, concentrations of stem cells, osteoblasts and growth factor at $T = 2$ days **a** Concentration of stem cells $s$, **b** Concentration of osteoblasts $b$, **c** Bone matrix density $m$, **d** Concentration of the growth factor $g$



**Fig. 4** Evolution of the bone matrix density **a** $T = 12\,h$, **b** $T = 48\,h$

(see Fig. 3). These osteoblasts synthetized the new bone matrix, which evolution is shown on Fig. 4. The stem cells moved towards the center of the fracture. These results are in agreement with previous results [5, 10].

# References

1. Bailón-Plaza, A., Van Der Meulen, M.: A mathematical framework to study the effects of growth factor influences on fracture healing. J. Theor. Biol. **212**(2) (2001)
2. Bessemoulin-Chatard, M., Saad, M.: Analysis of a monotone combined finite volume-finite element scheme for a bone healing model. In preparation
3. Cancès, C., Cathala, M., Le Potier, C.: Monotone coercive cell-centered finite volume schemes for anisotropic diffusion equations. Numer. Math. **125**(3) (2013)

4. Chamoun, G., Saad, M., Talhouk, R.: Monotone combined edge finite volume-finite element scheme for anisotropic Keller-Segel model. Numer. Meth. Part. Diff. Equ. **30**(3) (2014)
5. Coudière, Y., Saad, M., Uzureau, A.: Analysis of a finite volume method for a bone growth system in vivo. Comput. Math. Appl. **66** (2013)
6. Crouzeix, M., Raviart, P.A.: Conforming and nonconforming finite element methods for solving the stationary Stokes equations. I. Rev. Française Automat. Informat. Recherche Opérationnelle Sér. Rouge **7** (1973)
7. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Handbook of numerical analysis, vol. 7, North-Holland, Amsterdam (2000)
8. Eymard, R., Hilhorst, D., Vohralik, M.: A combined finite volume-nonconforming/mixed hybrid finite element scheme for degenerate parabolic problems. Numer. Math. **105** (2006)
9. Feistauer, M., Felcman, J., Lukáčová-Medvidová, M.: Combined finite element-finite volume solution of compressible flow. J. Comput. Appl. Math. **63** (1995)
10. Uzureau, A.: Modélisations et calculs pour la cicatrisation osseuse. Application à la modélisation d'un bioréacteur. Ph.D. thesis, Université de Nantes (2012)

# Vertex Approximate Gradient Scheme for Hybrid Dimensional Two-Phase Darcy Flows in Fractured Porous Media

**Konstantin Brenner, Mayya Groza, Cindy Guichard and Roland Masson**

**Abstract**    This paper presents the Vertex Approximate Gradient (VAG) discretization of a two-phase Darcy flow in discrete fracture networks (DFN) taking into account the mass exchange between the matrix and the fracture. We consider the asymptotic model for which the fractures are represented as interfaces of codimension one immersed in the matrix domain with continuous pressures at the matrix fracture interface. Compared with Control Volume Finite Element (CVFE) approaches, the VAG scheme has the advantage to avoid the mixing of the fracture and matrix rocktypes at the interfaces between the matrix and the fractures, while keeping the low cost of a nodal discretization on unstructured meshes. The convergence of the scheme is proved under the assumption that the relative permeabilities are bounded from below by a strictly positive constant but cover the case of discontinuous capillary pressures. The efficiency of our approach compared with CVFE discretizations is shown on a 3D fracture network with very low matrix permeability.

K. Brenner (✉) · M. Groza · R. Masson
Laboratoire de Mathématiques J.A. Dieudonné UMR CNRS 7251 and team Coffee,
Université Nice Sophia Antipolis, CNRS and INRIA Sophia Antipolis Méditerranée, Parc
Valrose, 06108 Nice, France
e-mail: roland.masson@unice.fr

K. Brenner
e-mail: konstantin.brenner@unice.fr

M. Groza
e-mail: mayya.groza@unice.fr

C. Guichard
Laboratoire Jacques-Louis Lions, CNRS, UMR 7598, Sorbonne Universités,
UMPC Univ Paris 06, F-75005 Paris, France
e-mail: guichard@ljll.math.upmc.fr

# 1 Hybrid Dimensional Two-Phase Darcy Flow Model in Fractured Porous Media

Let $\Omega$ denotes a bounded polyhedral domain of $\mathbb{R}^d$, $d = 2, 3$. We consider the asymptotic model introduced in [1] where fractures are represented as interfaces of codimension 1. Let $\overline{\Gamma} = \bigcup_{i \in I} \overline{\Gamma}_i$ denotes the network of fractures $\Gamma_i \subset \Omega, i \in I$, such that each $\Gamma_i$ is a planar polygonal simply connected open domain. It is assumed that the angles of $\Gamma_i$ are strictly lower than $2\pi$ and that $\Gamma_i \cap \Gamma_j = \emptyset$ for all $i \neq j$. It is also assumed that $\Sigma_{i,0} = \overline{\Gamma}_i \cap \partial\Omega$ has a vanishing $d - 1$ measure. We will denote by $d\tau(\mathbf{x})$ the $d - 1$ dimensional Lebesgue measure on $\Gamma$. Let $H^1(\Gamma)$ denote the set of functions $v = (v_i)_{i \in I}$ such that $v_i \in H^1(\Gamma_i)$, $i \in I$ with continuous traces at the fracture intersections, and endowed with the norm $\|v\|_{H^1(\Gamma)}^2 = \sum_{i \in I} \|v_i\|_{H^1(\Gamma_i)}^2$. Its subspace with vanishing traces on $\Sigma_0 = \bigcup_{i \in I} \Sigma_{i,0}$ is denoted by $H_{\Sigma_0}^1(\Gamma)$. The gradient operator from $H^1(\Omega)$ to $L^2(\Omega)^d$ is denoted by $\nabla$, and the tangential gradient from $H^1(\Gamma)$ to $L^2(\Gamma)^{d-1}$ by $\nabla_\tau$. Let us also consider the trace operator $\gamma$ from $H^1(\Omega)$ to $L^2(\Gamma)$. We can now define the hybrid dimensional function spaces that will be used in the variational formulation of the two-phase Darcy flow model:

$$V = \{v \in H^1(\Omega),\ \gamma v \in H^1(\Gamma)\}, \text{ and its subspace}$$
$$V_0 = \{v \in H_0^1(\Omega),\ \gamma v \in H_{\Sigma_0}^1(\Gamma)\}.$$

The space $V_0$ is endowed with the norm $\|v\|_V^2 = \|\nabla v\|_{L^2(\Omega)^d}^2 + \|\nabla_\tau \gamma v\|_{L^2(\Gamma)^{d-1}}^2$.

Let $u^2$ (resp. $u^1$) denote the wetting (resp. non wetting) phase pressure, $p = u^1 - u^2$ the capillary pressure, and $p_{\text{ini}} \in V$ the initial capillary pressure distribution. For the sake of simplicity in the convergence analysis, homogeneous Dirichlet boundary conditions are assumed for $u^1$ and $u^2$ at the boundary $\partial\Omega$, as well as at $\Sigma_0$ for $\gamma u^1$ and $\gamma u^2$. Let us denote by $S_m^1(\mathbf{x}, p)$ (resp. $S_f^1(\mathbf{x}, p)$) the inverses of the monotone graph extension of the capillary pressure curves in the matrix domain $\Omega$ (resp. in the fracture network $\Gamma$), and let us set $S_m^2 = 1 - S_m^1$ (resp. $S_f^2 = 1 - S_f^1$). In the matrix domain $\Omega$ (resp. in the fracture network $\Gamma$), let us denote by $k_m^\alpha(\mathbf{x}, S_m^\alpha)$ (resp. $k_f^\alpha(\mathbf{x}, S_f^\alpha)$), $\alpha \in \{1, 2\}$, the phase mobilities, by $\phi_m(\mathbf{x})$ (resp. $\phi_f(\mathbf{x})$) the porosity, and by $\Lambda_m(\mathbf{x})$ (resp. $\Lambda_f(\mathbf{x})$) the permeability tensor. We also denote by $d_f(\mathbf{x})$, $\mathbf{x} \in \Gamma$ the width of the fractures, and by $d\tau_f(\mathbf{x})$ the weighted Lebesgue $d - 1$ dimensional measure on $\Gamma$ defined by $d\tau_f(\mathbf{x}) = d_f(\mathbf{x})d\tau(\mathbf{x})$. The hybrid dimensional phase pressures weak formulation amounts to find $u^1, u^2 \in L^2(0, T; V_0)$ satisfying the following variational equalities for $\alpha \in \{1, 2\}$, and for all $\varphi \in C_c^\infty([0, T[\times\Omega):$

$$
\begin{cases}
\displaystyle\int_0^T \int_\Omega \Big( -\phi_m(\mathbf{x})S_m^\alpha(\mathbf{x},p)\partial_t\varphi(\mathbf{x},t) + k_m^\alpha(\mathbf{x},S_m^\alpha(\mathbf{x},p))\Lambda_m(\mathbf{x})\nabla u^\alpha(\mathbf{x},t) \\
\qquad\cdot\nabla\varphi(\mathbf{x},t)\Big)d\mathbf{x}dt \\
\displaystyle + \int_0^T \int_\Gamma -\phi_f(\mathbf{x})S_f^\alpha(\mathbf{x},\gamma p)\partial_t\gamma\varphi(\mathbf{x},t)d\tau_f(\mathbf{x})dt \\
\displaystyle + \int_0^T \int_\Gamma k_f^\alpha(\mathbf{x},S_f^\alpha(\mathbf{x},\gamma p))\Lambda_f(\mathbf{x})\nabla_\tau\gamma u^\alpha(\mathbf{x},t)\cdot\nabla_\tau\gamma\varphi(\mathbf{x},t)d\tau_f(\mathbf{x})dt \\
\displaystyle + \int_\Omega \phi_m(\mathbf{x})S_m^\alpha(\mathbf{x},p_{\mathrm{ini}})\varphi(\mathbf{x},0)d\mathbf{x}dt + \int_\Gamma \phi_f(\mathbf{x})S_f^\alpha(\mathbf{x},\gamma p_{\mathrm{ini}})\varphi(\mathbf{x},0)d\tau_f(\mathbf{x})dt \\
\displaystyle - \int_0^T \int_\Omega h_m^\alpha(\mathbf{x},t)\varphi(\mathbf{x},t)d\mathbf{x}dt - \int_0^T \int_\Gamma h_f^\alpha(\mathbf{x},t)\gamma\varphi(\mathbf{x},t)d\tau_f(\mathbf{x})dt = 0,
\end{cases}
\tag{1}
$$

where the function $h_m^\alpha$ (resp. $h_f^\alpha$), $\alpha \in \{1,2\}$ stands for the source term in the matrix domain $\Omega$ (resp. in the fracture network $\Gamma$).

## 2 Vertex Approximate Gradient Discretization

In the spirit of [3], we consider generalised polyhedral meshes of $\Omega$. Let $\mathscr{M}$ be the set of cells that are disjoint open polyhedral subsets of $\Omega$ such that $\bigcup_{K\in\mathscr{M}} \overline{K} = \overline{\Omega}$. For all $K \in \mathscr{M}$, $\mathbf{x}_K$ denotes the so-called "centre" of the cell $K$ under the assumption that $K$ is star-shaped with respect to $\mathbf{x}_K$. We then denote by $\mathscr{F}_K$ the set of interfaces of non zero $d - 1$ dimensional measure among the interior faces $\overline{K} \cap \overline{L}$, $L \in \mathscr{M}$, and the boundary interface $\overline{K} \cap \partial\Omega$, which possibly splits in several boundary faces. Let us denote by $\mathscr{F} = \bigcup_{K\in\mathscr{M}} \mathscr{F}_K$ the set of all faces of the mesh. Remark that the faces are not assumed to be planar, hence the term "generalised polyhedral mesh". For $\sigma \in \mathscr{F}$, let $\mathscr{E}_\sigma$ be the set of interfaces of non zero $d - 2$ dimensional measure among the interfaces $\sigma \cap \sigma'$, $\sigma' \in \mathscr{F}$. Then, we denote by $\mathscr{E} = \bigcup_{\sigma\in\mathscr{F}} \mathscr{E}_\sigma$ the set of all edges of the mesh. Let $\mathscr{V}_\sigma = \bigcup_{e,e'\in\mathscr{E}_\sigma,e\neq e'} (e \cap e')$ be the set of vertices of $\sigma$, for each $K \in \mathscr{M}$ we define $\mathscr{V}_K = \bigcup_{\sigma\in\mathscr{F}_K} \mathscr{V}_\sigma$, and we also denote by $\mathscr{V} = \bigcup_{K\in\mathscr{M}} \mathscr{V}_K$ the set of all vertices of the mesh. It is then assumed that for each face $\sigma \in \mathscr{F}$, there exists a so-called "centre" of the face $\mathbf{x}_\sigma \in \sigma \setminus \bigcup_{e\in\mathscr{E}_\sigma} e$ such that $\mathbf{x}_\sigma = \sum_{\mathbf{s}\in\mathscr{V}_\sigma} \beta_{\sigma,\mathbf{s}}\,\mathbf{x_s}$, with $\sum_{\mathbf{s}\in\mathscr{V}_\sigma} \beta_{\sigma,\mathbf{s}} = 1$, and $\beta_{\sigma,\mathbf{s}} \geq 0$ for all $\mathbf{s} \in \mathscr{V}_\sigma$; moreover the face $\sigma$ is assumed to match with the union of the triangles $T_{\sigma,e}$ defined by the face centre $\mathbf{x}_\sigma$ and each edge $e \in \mathscr{E}_\sigma$. The mesh is also supposed to be conforming w.r.t. the fracture network $\Gamma$ in the sense that for all $i \in I$ there exist the subsets $\mathscr{F}_{\Gamma_i}$ of $\mathscr{F}$ such that $\overline{\Gamma}_i = \bigcup_{\sigma\in\mathscr{F}_{\Gamma_i}} \overline{\sigma}$. We will denote by $\mathscr{F}_\Gamma$ the set of fracture faces $\bigcup_{i\in I} \mathscr{F}_{\Gamma_i}$. This geometrical discretization of $\Omega$ and $\Gamma$ is denoted in the following by $\mathscr{D}$.

The VAG discretization has been introduced in [3] for diffusive problems on heterogeneous anisotropic media. Its extension to the hybrid dimensional Darcy model is based on the following vector space of degrees of freedom:

$$X_{\mathcal{D}} = \{v_K, v_{\mathbf{s}}, v_\sigma \in \mathbb{R}, K \in \mathcal{M}, \mathbf{s} \in \mathcal{V}, \sigma \in \mathcal{F}_\Gamma\},$$

and its subspace with homogeneous Dirichlet boundary conditions on $\partial\Omega$:

$$X_{\mathcal{D}}^0 = \{v \in X_{\mathcal{D}} \mid v_{\mathbf{s}} = 0 \text{ for } \mathbf{s} \in \mathcal{V}_{ext}\}.$$

where $\mathcal{V}_{ext} = \mathcal{V} \cap \partial\Omega$ denotes the set of boundary vertices, and $\mathcal{V}_{int} = \mathcal{V} \setminus \partial\Omega$ denotes the set of interior vertices.

A finite element discretization of $V$ is built using a tetrahedral sub-mesh of $\mathcal{M}$ and a second order interpolation at the face centres $\mathbf{x}_\sigma$, $\sigma \in \mathcal{F} \setminus \mathcal{F}_\Gamma$ defined by the operator $I_\sigma : X_{\mathcal{D}} \to \mathbb{R}$ such that $I_\sigma(v) = \sum_{\mathbf{s}\in\mathcal{V}_\sigma} \beta_{\sigma,\mathbf{s}} v_{\mathbf{s}}$. The tetrahedral sub-mesh is defined by $\mathcal{T} = \{T_{K,\sigma,e}, e \in \mathcal{E}_\sigma, \sigma \in \mathcal{F}_K, K \in \mathcal{M}\}$ where $T_{K,\sigma,e}$ is the tetrahedron joining the cell centre $\mathbf{x}_K$ to the triangle $T_{\sigma,e}$.

For a given $v \in X_{\mathcal{D}}$, we define the function $\pi_{\mathcal{T}} v \in V$ as the continuous piecewise affine function on each tetrahedron of $\mathcal{T}$ such that $\pi_{\mathcal{T}} v(\mathbf{x}_K) = v_K$, $\pi_{\mathcal{T}} v(\mathbf{s}) = v_{\mathbf{s}}$, $\pi_{\mathcal{T}} v(\mathbf{x}_\sigma) = v_\sigma$, and $\pi_{\mathcal{T}} v(\mathbf{x}_{\sigma'}) = I_{\sigma'}(v)$ for all $K \in \mathcal{M}$, $\mathbf{s} \in \mathcal{V}$, $\sigma \in \mathcal{F}_\Gamma$, and $\sigma' \in \mathcal{F} \setminus \mathcal{F}_\Gamma$. Discrete gradient operators are defined this from finite element discretization of $V$ by

$$\nabla_{\mathcal{D}_m} : X_{\mathcal{D}} \to L^2(\Omega)^d \text{ such that } \nabla_{\mathcal{D}_m} v = \nabla \pi_{\mathcal{T}} v,$$

in the matrix, and by

$$\nabla_{\mathcal{D}_f} : X_{\mathcal{D}} \to L^2(\Gamma)^{d-1} \text{ such that } \nabla_{\mathcal{D}_f} v = \nabla_\tau \gamma \pi_{\mathcal{T}} v,$$

in the fracture network. In addition, the VAG discretization uses two non conforming piecewise constant reconstructions of functions from $X_{\mathcal{D}}$ into respectively $L^2(\Omega)$ and $L^2(\Gamma)$ based on partitions of each cell and of each fracture face denoted by $K = \omega_K \bigcup \left(\bigcup_{\mathbf{s}\in\mathcal{V}_K\cap\mathcal{V}_{int}} \omega_{K,\mathbf{s}}\right) \bigcup \left(\bigcup_{\sigma\in\mathcal{F}_K\cap\mathcal{F}_\Gamma} \omega_{K,\sigma}\right)$, for all $K \in \mathcal{M}$, and by $\sigma = \Sigma_\sigma \bigcup \left(\bigcup_{\mathbf{s}\in\mathcal{V}_\sigma\cap\mathcal{V}_{int}} \Sigma_{\sigma,\mathbf{s}}\right)$ for all $\sigma \in \mathcal{F}_\Gamma$. Then, the function reconstruction operators are defined by $\pi_{\mathcal{D}_m} v(\mathbf{x}) = \begin{cases} v_K \text{ for all } \mathbf{x} \in \omega_K, \ K \in \mathcal{M}, \\ v_{\mathbf{s}} \text{ for all } \mathbf{x} \in \omega_{K,\mathbf{s}}, \ \mathbf{s} \in \mathcal{V}_K \cap \mathcal{V}_{int}, \ K \in \mathcal{M}, \\ v_\sigma \text{ for all } \mathbf{x} \in \omega_{K,\sigma}, \ \sigma \in \mathcal{F}_K \cap \mathcal{F}_\Gamma, \ K \in \mathcal{M}, \end{cases}$

and $\pi_{\mathcal{D}_f} v(\mathbf{x}) = \begin{cases} v_\sigma \text{ for all } \mathbf{x} \in \Sigma_\sigma, \ \sigma \in \mathcal{F}_\Gamma, \\ v_{\mathbf{s}} \text{ for all } \mathbf{x} \in \Sigma_{\sigma,\mathbf{s}}, \ \mathbf{s} \in \mathcal{V}_\sigma \cap \mathcal{V}_{int}, \ \sigma \in \mathcal{F}_\Gamma. \end{cases}$

It is important to notice that, in the practical case when the space discretization is conforming with respect to the heterogeneities and when the source term $h_m^\alpha$ (resp. $h_f^\alpha$) is a cellwise (resp. facewise) constant function, the implementation of the VAG scheme does not require to build these partitions. In that case, it is sufficient to define the matrix volume fractions $\alpha_{K,\mathbf{s}} = \frac{\int_{\omega_{K,\mathbf{s}}} d\mathbf{x}}{\int_K d\mathbf{x}}, \mathbf{s} \in \mathcal{V}_K \cap \mathcal{V}_{int}, K \in \mathcal{M}, \alpha_{K,\sigma} = \frac{\int_{\omega_{K,\sigma}} d\mathbf{x}}{\int_K d\mathbf{x}}, \sigma \in \mathcal{F}_K \cap \mathcal{F}_\Gamma, K \in \mathcal{M}$, constrained to satisfy $\alpha_{K,\mathbf{s}} \geq 0, \alpha_{K,\sigma} \geq 0$, and $\sum_{\mathbf{s}\in\mathcal{V}_K\cap\mathcal{V}_{int}} \alpha_{K,\mathbf{s}} + \sum_{\sigma\in\mathcal{F}_K\cap\mathcal{F}_\Gamma} \alpha_{K,\sigma} \leq 1$, as well as the fracture volume fractions

$\alpha_{\sigma,\mathbf{s}} = \frac{\int_{\Sigma_{\sigma,\mathbf{s}}} d\tau_f(\mathbf{x})}{\int_\sigma d\tau_f(\mathbf{x})}$, $\mathbf{s} \in \mathscr{V}_\sigma \cap \mathscr{V}_{int}$, $\sigma \in \mathscr{F}_\Gamma$, constrained to satisfy $\alpha_{\sigma,\mathbf{s}} \geq 0$, and $\sum_{\mathbf{s}\in\mathscr{V}_\sigma\cap\mathscr{V}_{int}} \alpha_{\sigma,\mathbf{s}} \leq 1$. The convergence of the scheme will be shown to hold whatever the choice of these volume fractions. As will be detailed in the numerical section, this flexibility is a crucial asset, compared with usual CVFE approaches [5, 6], in order to improve the accuracy of the scheme for highly heterogeneous test cases.

For $N \in \mathbb{N}^*$, let us consider generally nonuniform discretization $t^0 = 0 < t^1 < \cdots < t^{n-1} < t^n \cdots < t^N = T$ of the time interval $[0, T]$. We denote the time steps by $\Delta t^n = t^n - t^{n-1}$ for all $n \in \{1, \cdots, N\}$ while $\Delta t$ stands for the whole sequence $(\Delta t^n)_{n\in\{1,\dots,N\}}$. Let us denote by $u^{\alpha,n} \in X_{\mathscr{D}}^0$, $\alpha \in \{1, 2\}$ the discrete phase pressures, and by $p^n = u^{1,n} - u^{2,n}$ the discrete capillary pressure at time $t^n$ for all $n \in \{1, \cdots, N\}$. Given an approximation $p^0 \in X_{\mathscr{D}}$ of the initial capillary pressure $p_{\mathrm{ini}}$, the VAG discretization of the two-phase Darcy flow model in phase pressures formulation (1) looks for $u^\alpha = \left(u^{\alpha,n} \in X_{\mathscr{D}}^0\right)_{n\in\{1,\cdots,N\}}$, $\alpha \in \{1, 2\}$, such that for $\alpha \in \{1, 2\}$, and for all $v \in X_{\mathscr{D}}^0$ one has

$$
\begin{cases}
\displaystyle \int_\Omega \phi_m \frac{S_{\mathscr{D}_m}^{\alpha,n} - S_{\mathscr{D}_m}^{\alpha,n-1}}{\Delta t^n} \pi_{\mathscr{D}_m} v \, d\mathbf{x} + \int_\Omega k_{\mathscr{D}_m}^{\alpha,n} \Lambda_m \nabla_{\mathscr{D}_m} u^{\alpha,n} \cdot \nabla_{\mathscr{D}_m} v \, d\mathbf{x} \\
\displaystyle + \int_\Gamma \phi_f \frac{S_{\mathscr{D}_f}^{\alpha,n} - S_{\mathscr{D}_f}^{\alpha,n-1}}{\Delta t^n} \pi_{\mathscr{D}_f} v \, d\tau_f(\mathbf{x}) + \int_\Gamma k_{\mathscr{D}_f}^{\alpha,n} \Lambda_f \nabla_{\mathscr{D}_f} u^{\alpha,n} \cdot \nabla_{\mathscr{D}_f} v \, d\tau_f(\mathbf{x}) \\
\displaystyle = \frac{1}{\Delta t^n} \int_{t^{n-1}}^{t^n} \left( \int_\Omega h_m^\alpha \pi_{\mathscr{D}_m} v \, d\mathbf{x} + \int_\Gamma h_f^\alpha \pi_{\mathscr{D}_f} v \, d\tau_f(\mathbf{x}) \right) dt,
\end{cases} \tag{2}
$$

where $S_{\mathscr{D}_m}^{\alpha,n}(\mathbf{x}) = S_m^\alpha(\mathbf{x}, \pi_{\mathscr{D}_m} p^n(\mathbf{x}))$, $S_{\mathscr{D}_f}^{\alpha,n}(\mathbf{x}) = S_f^\alpha(\mathbf{x}, \pi_{\mathscr{D}_f} p^n(\mathbf{x}))$, and $k_{\mathscr{D}_m}^{\alpha,n}(\mathbf{x}) = k_m^\alpha(\mathbf{x}, S_{\mathscr{D}_m}^{\alpha,n}(\mathbf{x}))$, $k_{\mathscr{D}_f}^{\alpha,n}(\mathbf{x}) = k_f^\alpha(\mathbf{x}, S_{\mathscr{D}_f}^{\alpha,n}(\mathbf{x}))$.

**Convergence analysis**: We present in Theorem 1 below the main theoretical result obtained in [2]. Let $\rho_T$ denote the insphere diameter of a given tetrahedron $T$, $h_T$ its diameter, and $h_{\mathscr{T}} = \max_{T\in\mathscr{T}} h_T$. We will assume in the convergence analysis that the family of tetrahedral submeshes $\mathscr{T}$ is shape regular in the sense that $\theta_{\mathscr{T}} = \max_{T\in\mathscr{T}} \frac{h_T}{\rho_T}$ and $\gamma_{\mathscr{M}} = \max_{K\in\mathscr{M}} \mathrm{Card}(\mathscr{V}_K)$ are uniformly bounded. The assumptions on the data are natural extensions to our hybrid dimensional model (1) of the assumptions stated in [4]. They are quite general, except for the assumption $k_m^\alpha(\mathbf{x}, s)$ (resp. $k_f^\alpha(\mathbf{x}, s)$) $\in [k_{\min}, k_{\max}]$ for $(\mathbf{x}, s) \in \Omega \times [0, 1]$ (resp. $(\mathbf{x}, s) \in \Gamma \times [0, 1]$) which is needed in the following convergence analysis but not in the practical implementation of the scheme. Using the discrete phase pressures as test functions in the discrete variational formulation (2), we deduce the following a priori estimate.

**Lemma 1** *Let $u^\alpha$, $\alpha \in \{1, 2\}$, be a solution to (2), then, there exists $C > 0$ depending only on the data and on $\gamma_{\mathscr{M}}$ and $\theta_{\mathscr{T}}$ such that $\sum_{\alpha\in\{1,2\}} \sum_{n=1}^N \Delta t^n \|\pi_{\mathscr{T}} u^{\alpha,n}\|_V^2 \leq C$.*

Using a topological degree argument, this estimate allows to obtain the existence of a discrete solution to (2). For all $v = \left(v^n \in X_{\mathscr{D},\Delta t}\right)_{n\in\{1,\cdots,N\}}$ let us define

$\pi_{\mathcal{D}_m, \Delta t} v(\mathbf{x}, t) = \pi_{\mathcal{D}_m} v^n(\mathbf{x}), \pi_{\mathcal{D}_f, \Delta t} v(\mathbf{x}, t) = \pi_{\mathcal{D}_f} v^n(\mathbf{x}), \pi_{\mathcal{T}, \Delta t} v(\mathbf{x}, t) = \pi_{\mathcal{T}} v^n(\mathbf{x})$
for all $(\mathbf{x}, t) \in \Omega \times (t^{n-1}, t^n], n \in \{1, \cdots, N\}$. We also define for $\alpha \in \{1, 2\}$ the
functions $S^\alpha_{\mathcal{D}_m, \Delta t}(\mathbf{x}, t) = S^\alpha(\mathbf{x}, \pi_{\mathcal{D}_m, \Delta t} p(\mathbf{x}, t))$ and $S^\alpha_{\mathcal{D}_f, \Delta t}(\mathbf{x}, t) = S^\alpha(\mathbf{x}, \pi_{\mathcal{D}_f, \Delta t} p(\mathbf{x}, t))$.

**Theorem 1** *Let $(\mathcal{D}^{(k)}, \Delta t^{(k)})_{k \in \mathbb{N}}$ be a sequence of space-time discretizations such
that there exist two positive constants $\theta$ and $\gamma$ satisfying $\theta_{\mathcal{T}^{(k)}} \leq \theta, \gamma_{\mathcal{M}^{(k)}} \leq \gamma$ for all
$k \in \mathbb{N}$ and such that $h_{\mathcal{T}^{(k)}}, \max_n \Delta t^{(k),n} \to 0$ as $m \to \infty$. Let $u^{\alpha,(k)} \in X_{\mathcal{D}^{(k)}, \Delta t^{(k)}}$,
$S^\alpha_{\mathcal{D}_m^{(k)}, \Delta t^{(k)}}$ and $S^\alpha_{\mathcal{D}_f^{(k)}, \Delta t^{(k)}}, \alpha \in \{1, 2\}$, be s.t (2) holds for all $m \in \mathbb{N}$. It is also
assumed that $\pi_{\mathcal{D}_m^{(k)}} p^{0,(k)}$ converges strongly to $p^{ini}$ in $L^2(\Omega)$, and that $\pi_{\mathcal{D}_f^{(k)}} p^{0,(k)}$
converges strongly to $\gamma p^{ini}$ in $L^2(\Gamma)$. Then, there exists a weak solution $(\overline{u}^1, \overline{u}^2)$
with $\overline{p} = \overline{u}^1 - \overline{u}^2$ to the problem (1) such that for each phase $\alpha \in \{1, 2\}$
and up to a subsequence, one has $\pi_{\mathcal{T}^{(k)}, \Delta t^{(k)}} u^{\alpha,(k)} \rightharpoonup \overline{u}^\alpha$ in $L^2(\Omega \times (0, T))$ and
$\gamma \pi_{\mathcal{T}^{(k)}, \Delta t^{(k)}} u^{\alpha,(k)} \rightharpoonup \gamma \overline{u}^\alpha$ in $L^2(\Gamma \times (0, T)), S^\alpha_{\mathcal{D}_m^{(k)}, \Delta t^{(k)}} \to S^\alpha_m(., \overline{p})$ in $L^2(\Omega \times (0, T))$ and $S^\alpha_{\mathcal{D}_f^{(k)}, \Delta t^{(k)}} \to S^\alpha_f(., \gamma \overline{p})$ in $L^2(\Gamma \times (0, T))$.*

## 3 Numerical Experiments

This test case considers the migration of oil in a 3D basin $\Omega = (0, L) \times (0, L) \times (0, H)$
with $H = L = 100$ m. The family of 4 tetrahedral meshes is generated using TetGen
[7] in order to be refined at the neighbourhood of the fracture network with a number
of cells ranging from 47670 to 3076262 and a factor of refinement at the matrix
fracture interface of roughly 12 (see Fig. 1 for the coarsest mesh $i_{mesh} = 1$).

The permeability of the matrix $\Lambda_m = \lambda_m \text{Id}$ and the permeability of the fractures
$\Lambda_f = \lambda_f \text{Id}$ are highly contrasted with $\Lambda_m = 10^{-17} \text{m}^2, \Lambda_f = 10^{-11} \text{m}^2$. The
width of the fractures is fixed to $d_f = 0.01$ m and their porosity to $\phi_f = 0.3$. The
porosity of the matrix is set to $\phi_m = 0.1$. The inverses of the capillary pressure
monotone graph in the matrix $(j = m)$ and in the fractures $(j = f)$ are defined
by the Corey law $S^1_j(p) = 0$ if $p < 0$ and $S^1_j(p) = (1 - s^2_{r,j})(1 - e^{\frac{-p}{b_j}})$ if $p \geq 0$
with the rocktype $b_m = 5 \cdot 10^3$ Pa, $s^2_{r,m} = 0.2, s^1_{r,m} = 0$ in the matrix and the
rocktype $b_f = 10^2$ Pa, $s^2_{r,f} = s^1_{r,f} = 0$ in the fractures. The mobilities are defined
for $j = m$ and $j = f$ by the following Corey law $k^\alpha_j(\mathbf{x}, s^\alpha) = 0$ if $\overline{s}^\alpha < 0$,
$k^\alpha_j(\mathbf{x}, s^\alpha) = \frac{1}{\mu^\alpha}$ if $\overline{s}^\alpha > 1$ and $k^\alpha_j(\mathbf{x}, s^\alpha) = \frac{(\overline{s}^\alpha)^2}{\mu^\alpha}$ if $\overline{s}^\alpha \in (0, 1)$, for phase $\alpha = 1$
(oil), and phase $\alpha = 2$ (water) where $\overline{s}^1 = \frac{s^1 - s^1_{r,j}}{1 - s^1_{r,j} - s^2_{r,j}}$, and $\overline{s}^2 = \frac{s^2 - s^2_{r,j}}{1 - s^2_{r,j} - s^1_{r,j}}$ are the
reduced saturations, and $\mu^1 = 0.005$ Pa.s and $\mu^2 = 0.001$ Pa.s are the viscosities of
the phases. The densities of phases are fixed to $\rho^1 = 700$ Kg/m$^3$ for the oil phase and
$\rho^2 = 1000$ Kg/m$^3$ for the water phase. Phase 1 is injected at the bottom boundary with
imposed pressures $u^2(\mathbf{x}) = 8.1 \cdot 10^6 + \rho^2 g H$ Pa, $u^1(\mathbf{x}) = u^2(\mathbf{x}) + (S^1_f)^{-1}(0.999999)$

**Fig. 1** Coarsest mesh
$i_{mesh} = 1$ and discrete solu-
tion obtained by the VAG-1
scheme with this mesh at final
time: oil saturation in the frac-
ture network and in the matrix
using the lower threshold in
the matrix equal to 0.001



corresponding to an input phase 1 saturation $s^1 = 0.999999$ in the fractures. At the
top boundary, the phase pressures are fixed to $u^2(\mathbf{x}) = 8 \cdot 10^6$ and $u^1(\mathbf{x}) = u^2(\mathbf{x})$. The
remaining boundaries of the basin are assumed to be impervious. The boundaries
of the fracture network not located at the top or bottom boundaries of the basin are
also assumed impervious. At initial time the porous media is saturated with phase 2
with a hydrostatic pressure $u_{ini}^2(\mathbf{x}) = 8 \cdot 10^6 + \rho_2 g(H - y)$, and a phase 1 pressure
defined by $u_{ini}^1(\mathbf{x}) = u_{ini}^2(\mathbf{x})$.

The implementation of the VAG scheme is based on a flux formulation with
upwinding of the mobilities rather than the discrete variational formulation (2) in
order to improve the stability of the solution on coarse meshes for convective dom-
inant regimes. The nonlinear systems obtained at each time step are solved by a
Newton Raphson algorithm. The time stepping is defined by an initial time step, a
maximum time step and the following rule: if the Newton solver does not converge
after 20 iterations, the time step is chopped by a factor 2 and recomputed. The time
step is increased by a factor 1.2 after each successful time step until it reaches again
the maximum time step. The stopping criteria are fixed to $10^{-7}$ for the GMRes solver
and to $10^{-6}$ for the Newton solver. A CPR-AMG right preconditioner is used in the
GMRes iterative solver. Let us also stress that, using the two equations in each cell,
the cell unknowns are eliminated from the discrete linearized system at each New-
ton iteration without any fill-in, reducing the Jacobian system to nodal and fracture
face unknowns only. The simulation is run over a period of 10 years with an initial
time step of 0.2 days, and a maximum time step fixed to 5 days, except on mesh 4
for which a smaller maximum time step of 2.5 days is used. All the numerical tests
have been performed on the Cicada Cluster located at the University Nice Sophia-
Antipolis and which includes 1152 nodes equipped with two eight-core Intel(R)
E5-2670 processors. Figure 1 exhibits the oil saturation obtained on the coarsest
mesh $i_{mesh} = 1$ at final simulation time. We observe that the oil phase injected at
the bottom side in the domain initially saturated with water, quickly rises by gravity
along the faults and slowly penetrate in the matrix.

Figure 2 compares the convergence of the oil saturation on the family of refined
meshes for two choices of the volume fractions VAG-1 and VAG-2. The choice VAG-2
simply set $\alpha_{K,\mathbf{s}} = \alpha_{K,\sigma} = 0.1$ and $\alpha_{\sigma,e} = 0.075$ on the whole mesh while the choice
VAG-1 does not mix the fracture and matrix porous volume taking $\alpha_{K,\mathbf{s}} = \alpha_{K,\sigma} = 0$
for all $\mathbf{s} \in \mathcal{V}_\Gamma$ and $\sigma \in \mathcal{F}_\Gamma$. One can see that, for such high ratio of the fracture

**Fig. 2** Volumes of oil in the fracture and in the matrix function of time for the family of meshes $i_{mesh} = 1, \ldots, 4$, and for both choices VAG-1 and VAG-2 of the volume distribution

**Table 1** For each choice VAG-1 and VAG-2 of the volume distribution and for each mesh $i_{mesh} = 1, \ldots, 4$: number $N_{\Delta t}$ of successful time steps, number $N_{Chop}$ of time step chops, number $N_{Newton}$ of Newton iterations per successful time step, number $N_{GMRes}$ of GMRes iterations by Newton iteration, CPU time in seconds

| $i_{mesh}$ | | $N_{\Delta t}$ | $N_{Chop}$ | $N_{Newton}$ | $N_{GMRes}$ | CPU (s) | | $N_{\Delta t}$ | $N_{Chop}$ | $N_{Newton}$ | $N_{GMRes}$ | CPU (s) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | VAG-1 | 384 | 6 | 2.20 | 10.05 | 588 | VAG-2 | 373 | 0 | 1.87 | 6.94 | 482 |
| 2 | VAG-1 | 390 | 10 | 3.08 | 15.11 | 5 898 | VAG-2 | 373 | 0 | 2.42 | 13.05 | 4 452 |
| 3 | VAG-1 | 415 | 21 | 4.02 | 15.93 | 31 806 | VAG-2 | 375 | 1 | 3.02 | 14.56 | 21 645 |
| 4 | VAG-1 | 784 | 30 | 3.37 | 16.75 | 209 485 | VAG-2 | 747 | 13 | 2.92 | 16.55 | 172 946 |

and matrix permeabilities, VAG-1 seems to provide a much better convergence than VAG-2 since it does not mix porous volumes from the matrix and the fracture network. This is confirmed by the 2D test case presented in [2] for which a reference solution is computed on a fine grid. It illustrates again the advantage of the VAG scheme compared with CVFE discretizations which cannot avoid such mixing of porous volumes [5, 6]. Table 1 presents the numerical behaviour of the simulations for both choices of the distribution of the volumes and for the family of meshes. The results obtained demonstrate the good robustness and scalability of the proposed numerical scheme both in terms of Newton convergence, linear solver convergence and CPU time.

# References

1. Alboin, C., Jaffré, J., Roberts, J., Serres, C.: Modeling fractures as interfaces for flow and transport in porous media. Fluid flow transp. porous media **295**, 13–24 (2002)
2. Brenner, K., Groza, M., Guichard, C., Masson, R.: Vertex approximate gradient scheme for hybrid dimensional two-phase darcy flows in fractured porous media. Preprinted,

http://hal.archives-ouvertes.fr/hal-00910939 (2013)

3. Eymard, R., Guichard, C., Herbin, R.: Small-stencil 3d schemes for diffusive flows in porous media. ESAIM: M2AN **46**, 265–290 (2010)
4. Eymard, R., Guichard, C., Herbin, R., Masson, R.: Gradient schemes for two-phase flow in heterogeneous porous media and Richards equation. ZAMM—J. Appl. Math. Mech. (2013). doi:10.1002/zamm.201200206
5. Firoozabadi, A., Monteagudo, J.E.: Control-volume model for simulation of water injection in fractured media: incorporating matrix heterogeneity and reservoir wettability effects. SPE J. **12**, 355–366 (2007)
6. Reichenberger, V., Jakobs, H., Bastian, P., Helmig, R.: A mixed-dimensional finite volume method for multiphase flow in fractured porous media. Adv. Water Resour. **29**, 1020–1036 (2006)
7. Si, H.: TetGen. A quality tetrahedral mesh generator and three-dimensional delaunay triangulator. (2007). http://tetgen.berlios.de

# Coupling of a Two Phase Gas Liquid Compositional 3D Darcy Flow with a 1D Compositional Free Gas Flow

**Konstantin Brenner, Roland Masson, Laurent Trenty and Yumeng Zhang**

**Abstract** A model coupling a three dimensional gas liquid compositional Darcy flow and a one dimensional compositional free gas flow is presented. The coupling conditions at the interface between the gallery and the porous media account for the molar normal fluxes continuity for each component, the gas liquid thermodynamical equilibrium, the gas pressure continuity and the gas and liquid molar fractions continuity. This model is applied to the simulation of the mass exchanges at the interface between the repository and the ventilation excavated gallery in a nuclear waste geological repository. The convergence of the Vertex Approximate Gradient discretization is analysed for a simplified model coupling the Richards approximation in the porous media and the gas pressure equation in the gallery.

## 1 Model

Let $\omega$ and $S \subset \omega$ be two simply connected domains of $\mathbb{R}^2$ and $\Omega = (0, L) \times (\omega \backslash \overline{S})$ be the cylindrical domain defining the porous media. The excavated gallery corresponds

K. Brenner · R. Masson (✉) · Y. Zhang
Laboratoire de Mathématiques J.A. Dieudonné UMR CNRS 7251 and Team Coffee,
Université Nice Sophia Antipolis, CNRS and INRIA Sophia Antipolis Méditerranée, Parc
Valrose, 06108 Nice, France
e-mail: roland.masson@unice.fr

K. Brenner
e-mail: konstantin.brenner@unice.fr

Y. Zhang
e-mail: y.zhang56@yahoo.fr

L. Trenty
Andra, 1 rue Jean Monnet, 92290 Chatenay-Malabry, France
e-mail: laurent.trenty@andra.fr

to the domain $(0, L) \times S$ and it is assumed that the free flow in the gallery depends only on the $x$ coordinate along the gallery and on the time $t$. Let us denote by $\Gamma = (0, L) \times \partial S$ the interface between the gallery and the porous media and by $\Gamma_D = ((0, L) \times \partial \omega) \cup (\{0\} \times (\omega \setminus \overline{S})) \cup (\{L\} \times (\omega \setminus \overline{S}))$ the remaining boundaries of $\Omega$.

Let $\alpha = g, l$ denote the gas and liquid phases assumed to be both defined by a mixture of components $i \in \mathscr{C}$ among which the water component denoted by $e$ which can vaporize in the gas phase, and a set of gaseous components $j \in \mathscr{C} \setminus \{e\}$ which can dissolve in the liquid phase. For the sake of simplicity, the model will be assumed to be isothermal with a fixed temperature $T$. We will denote by $c^\alpha = \left(c_i^\alpha, i \in \mathscr{C}\right)$ the vector of molar fractions of the components in the phase $\alpha = g, l$ with $\sum_{i \in \mathscr{C}} c_i^\alpha = 1$, and by $P^g$ and $P^l$ the two phase pressures. The mass densities of the phases are denoted by $\rho^\alpha(P^\alpha, c^\alpha)$ and the molar densities by $\zeta^\alpha(P^\alpha, c^\alpha)$, $\alpha \in \mathscr{P}$. They are related by $\rho^\alpha(P^\alpha, c^\alpha) = \left(\sum_{i \in \mathscr{C}} c_i^\alpha M_i\right) \zeta^\alpha(P^\alpha, c^\alpha)$, where $M_i, i \in \mathscr{C}$ are the molar masses of the components. For the sake of simplicity, it is assumed that the liquid molar density $\zeta^l$ is constant as well as the viscosities $\mu^\alpha$, $\alpha = g, l$.

The two phase Darcy's laws are characterized by the relative permeability functions $k_r^\alpha(\mathbf{x}, S^\alpha)$, for both phases $\alpha = g, l$, and by the capillary pressure function $P_c(\mathbf{x}, S^l)$, where $S^\alpha, \alpha = l, g$ denote the saturations of the phases with $S^g + S^l = 1$.

Each component $i \in \mathscr{C}$ is assumed to be at thermodynamical equilibrium between both phases which is characterized by the equality of its fugacities $f_i^\alpha$, $\alpha = g, l$ if both phases are present. The fugacities of the components in the gas phase are given by Dalton's law for an ideal mixture of perfect gas $f_i^g = c_i^g P^g$, $i \in \mathscr{C}$. The fugacities of the components in the liquid phase are given by Henry's law for the dissolution of the gaseous components in the liquid phase $f_j^l = c_j^l H_j(T)$, $j \in \mathscr{C} \setminus \{e\}$, and by Raoult-Kelvin's law for the water component in the liquid phase $f_e^l = c_e^l P_{sat}(T) \exp\left(\frac{-(P^g - P^l)}{\zeta^l RT}\right)$, where $P_{sat}(T)$ is the vapor pressure of the pure water.

Following [3], the gas liquid Darcy flow formulation uses both phase pressures $P^g$ and $P^l$ and the component fugacities $f = (f_i, i \in \mathscr{C})$ as set of primary unknowns. For this set of unknowns, the component molar fractions of an absent phase are extended by those at equilibrium with the present phase leading to define $c_i^\alpha(f, P^g, P^l)$, $\alpha = g, l, i \in \mathscr{C}$ by

$$\begin{cases} c_e^l = \frac{f_e}{P_{sat}(T)} \exp\left(\frac{(P^g - P^l)}{\zeta^l RT}\right), & c_j^l = \frac{f_j}{H_j(T)}, \; j \in \mathscr{C} \setminus \{e\} \\ c_e^g = \frac{f_e}{P^g}, & c_j^g = \frac{f_j}{P^g}, \; j \in \mathscr{C} \setminus \{e\}. \end{cases} \quad (1)$$

The pressure of an absent phase is also extended by the buble (for gas) and by the dew (for liquid) pressure leading to the equations $\sum_{i \in \mathscr{C}} c_i^\alpha(f, P^g, P^l) = 1$, $\alpha = g, l$. Finally, we define $\mathscr{S}^l(\mathbf{x}, .)$ as the inverse of the monotone graph extension of

the opposite of the capillary pressure $-P_c(\mathbf{x}, .)$. This leads to the following set of equations in the porous media:

$$
\begin{cases}
\phi \partial_t \sum_{\alpha \in \mathscr{P}} \zeta^\alpha S^\alpha c_i^\alpha + \mathrm{div}\left( \sum_{\alpha \in \mathscr{P}} \zeta^\alpha c_i^\alpha \mathbf{V}^\alpha \right) = 0, \ i \in \mathscr{C}, \\
\mathbf{V}^\alpha = -\frac{k_r^\alpha(\mathbf{x}, S^\alpha)}{\mu^\alpha} \mathbf{K}\left( \nabla P^\alpha - \rho^\alpha \mathbf{g} \right), \ \alpha = g, l, \\
S^g + S^l = 1, \quad S^l = \mathcal{S}^l(\mathbf{x}, P^l - P^g), \quad \sum_{i \in \mathscr{C}} c_i^\alpha(f, P^g, P^l) = 1, \ \alpha = g, l.
\end{cases}
\tag{2}
$$

In the gallery, the primary unknowns, depending only on the $x$ coordinate along the gallery and on the time $t$, are chosen to be the gas pressure $p$ and the gas molar fractions $c = (c_i, i \in \mathscr{C})$. The set of equations is defined by the following no pressure wave isothermal pipe flow model where $\alpha > 0, \beta > 0$ are parameters for the pressure drop along the gallery, $\mathbf{n}$ is the unit normal vector at $\Gamma$ outward to $\Omega$, and $|S|$ is the surface of the section $S$.

$$
\begin{cases}
\partial_t \left( |S|\zeta^g(p)c_i \right) + \partial_x \left( |S|\zeta^g(p)c_i w \right) = \int_{\partial S} \sum_{\alpha = g, l} \zeta^\alpha c_i^\alpha \mathbf{V}^\alpha \cdot \mathbf{n} \, ds, \\
\sum_{i \in \mathscr{C}} c_i = 1, \quad (\alpha w + \beta |w|w) = -\partial_x p.
\end{cases}
\tag{3}
$$

At the interface $\Gamma$ between the gallery and the porous media the coupling conditions are an adaptation to a 1D model in the gallery of [4]. Compared with [4], the gas pressure jump at the interface is neglected since a small flow rate between the porous media and the gallery is assumed due to the low permeability of the storage. Hence the coupling conditions account first for the continuity of the gas phase pressure $P^g = p$. Second, we impose the continuity of the gas molar fractions $c^g = c$, and third the thermodynamical equilibrium $f_i = f_i^l = f_i^g = p c_i$ for all $i \in \mathscr{C}$ together with $\sum_{i \in \mathscr{C}} c_i^l = 1$ which provides the additional equation (using (1)):

$$
P^g - P^l = -\zeta^l RT \ln \left( \frac{f_e}{P_{sat}(T)(1 - \sum_{j \in \mathscr{C} \setminus \{e\}} \frac{f_j}{H_j(T)})} \right).
\tag{4}
$$

## 2 Numerical Test

Let $\omega$ and $S$ be the disks of center $0$ and radius respectively $r_\omega = 10$ m and $r_S = 2$ m. We consider a radial mesh of the domain $(0, L) \times (\omega \setminus \bar{S})$, $L = 100$ m, exponentially refined at the interface of the gallery $\Gamma$ to account for the steep gradient of the capillary pressure. In addition to the water component $e$, we consider the air gaseous component denoted by $a$ with the Henry constant $H_a = 6 \cdot 10^9$ Pa. The gas molar density is given by $\zeta^g = \frac{p}{RT}$. The porous medium is initially saturated by the liquid phase with

**Fig. 1** $(x, r)$ cut of the storage and initial and boundary conditions of the test case



imposed pressure $P_{init}^l = 40 \cdot 10^5$ Pa and composition $c_a^l = 0$, $c_e^l = 1$. At the external boundary $r = r_\omega$ the water pressure is fixed to $P_0^l = P_{init}^l$, with an input composition $c_a^l = 0$, $c_e^l = 1$. On both sides $x = 0$ and $x = L$ of the porous media, zero flux boundary conditions are imposed. At the left side of the gallery $x = 0$, we consider a given velocity $w = w_0$ and an input relative humidity $H_r = H_{r,0} = 0.5$ with $H_r = \frac{c_e p}{P_{sat}(T)}$. The initial condition in the gallery is given by $p_{init} = 10^5$ Pa and $H_r = H_{r,0}$, and the pressure $p = p_{init}$ is fixed at the right side of the gallery (see Fig. 1). The relative permeabilities and capillary pressure are given by the Van-Genuchten laws with the parameters $n = 1.54$, $S_r^l = 0.01$, $S_r^g = 0$, $P_r = 2 \cdot 10^6$ Pa accounting for concrete rocktype with homogeneous isotropic permeability $K = 10^{-18}$ m$^2$ and porosity $\phi = 0.3$. For the pressure load we have taken $\alpha = 0$ and $\beta = 10^{-3}$ kg m$^{-4}$. The simulation is run over a period of 1,500 days with an initial time step of 1 s and a maximum time step of 10 days. The input velocity $w_0$ is fixed to $w_0 = 1$ m/s during the first 400 days, $w_0 = 0.01$ m/s during the next 600 days, and $w_0 = 0$ m/s during the last 500 days. In order to validate the simulation, an approximate stationary solution is computed for each $w_0$ assuming that we can neglect the dissolution of air, the gravity, the pressure drop in the gallery, and that the longitudinal derivatives are small compared with the radial derivatives in the porous media. Then, the stationary solution $c_e(x)$, $x \in (0, L)$ can be approximated by the solution of the following ODE for $w_0 > 0$: $\zeta^g w_0 (1 - c_{e,0}) \partial_x \left( \frac{c_e(x)}{1 - c_e(x)} \right) = \frac{2}{r_S^2} V_T(p_c(c_e(x)))$, $x \in (0, L)$ with $V_T(p_c) = \frac{\zeta^l K}{\mu^l \log(\frac{r_\omega}{r_S})} \left( P_0^l - p_{init} + \int_0^{p_c} k_r^l (\mathcal{S}^l(-u)) du \right)$, and $p_c(c_e) = -\zeta^l RT \ln \left( \frac{p_{init} c_e}{P_{sat}(T)} \right)$, using the boundary condition $c_e(0) = c_{e,0} = \frac{H_{r,0} P_{sat}(T)}{p_{init}}$, and by $H_r(x) = \exp \left( \frac{P_0^l - p_{init}}{\zeta^l RT} \right)$, $x \in (0, L)$ for $w_0 = 0$. In Fig. 2, we plot the average relative humidity in the gallery as well as the volume of gas in the porous medium function of time. Figure 3 plots the stationary numerical solution obtained for the gas saturation at the interface and in the porous media for each $w_0$. At the opening of the gallery at $t = 0$, we observe in Fig. 2 an increase of $H_r$ up to almost 1 in average in a few seconds due to a large liquid flow rate at the interface. Then, the flow rate decreases and we observe a drying of the gallery due to the ventilation at $w_0 = 1$ m/s down to an average relative humidity slightly above $H_{r,0}$ in a few days. Meanwhile the gas penetrate slowly into the porous medium reaching a stationnary state with around 13.5 m$^3$ of gas in say 400 days. When the input velocity is reduced to 0.01 m/s, we observe first a rapid increase of $H_r$ in say 1 day due to the reduced ventilation followed by a convergence to a second stationnary state with $H_r = 0.77$

**Fig. 2** Average of the relative humidity in the gallery and volume of gas in the porous medium function of time (*left*); stationary relative humidity for each $w_0$ compared with its approximate "analytical" solution (*right*)



**Fig. 3** $(x, r)$ cut of the stationary gas saturation at the interface (depending only on $x$ for $r \in (0, r_S)$) and in the porous medium (for $r > r_S$) for $w_0 = 1$ m/s (*left*), and $w_0 = 0.01$ m/s (*right*). Only the values above the threshold $10^{-3}$ are plotted

in average in the gallery and $12\,\mathrm{m}^3$ of gas in the porous medium. When $w_0$ is set to 0, $H_r$ reaches a value above 1 corresponding to $S^l = 1$ at the interface and the gas disappears from the porous medium in around $100\,\mathrm{days}$.

Figure 2 also compares the stationary numerical relative humidity obtained for each $w_0$ with its approximate "analytical" solution. A very good match is obtained.

# 3 Convergence Analysis of a Simplified Model

We consider the following simplified model using the Richards approximation in the porous medium and a single component equation in the gallery with linear pressure drop

$$
\begin{cases}
\phi \partial_t (\zeta^l S^l(., u)) + \mathrm{div}(\zeta^l \mathbf{V}^l) = 0, \\
\partial_t (|S| \zeta^g(p)) + \partial_x (-\frac{1}{\alpha} |S| \tilde{\zeta}^g(p) \partial_x p) = \int_{\partial S} \zeta^l \mathbf{V}^l \cdot \mathbf{n} \, ds, \\
\mathbf{V}^l = -\frac{k_r^\alpha(., S^l(., u))}{\mu^l} \mathbf{K}(\nabla u - M^l \zeta^l \mathbf{g}), \quad p = g(\gamma(u)),
\end{cases}
\tag{5}
$$

where $\gamma$ denotes the trace operator from $H^1(\Omega)$ to $H^{1/2}(\Gamma)$. The only primary unknown in the porous media is the liquid pressure denoted by $u$. The liquid mass

density is assumed to be fixed to $M^l \zeta^l$ where $M^l$ is the molar mass of the liquid phase. The thermodynamical equilibrium at the interface $\Gamma$ is accounted for by the relation $p = g(\gamma(u))$ with $g \in C^1(\mathbb{R}, \mathbb{R})$, $0 < c_1 \leq g'(q) \leq c_2$ for all $q \in \mathbb{R}$ and for given constants $c_1, c_2$. The function $g$ is a regularization for large positive and negative $u$ of $p = \frac{P_{sat}}{c_e} e^{\frac{u}{\zeta^l RT}}$ for given constants $1 \geq c_e > 0$ and $T > 0$. The molar gas density is set to $\zeta^g(p) = \frac{p}{RT}$ and is truncated in the flux term such that $\widetilde{\zeta}^g(p)$ is assumed to be a non decreasing function in $C^1(\mathbb{R}, \mathbb{R})$ bounded from below and above by two strictly positive constants and with a bounded derivative.

Let $\gamma_e$ be the trace operator from $H^1(\Omega)$ to $H^{1/2}(\Gamma_D)$. We define the function space $V = \{u \in H^1(\Omega) \mid \gamma u \in H^1(\Gamma), \partial_s \gamma u = 0\}$, where $s$ denotes the curvilinear coordinate along $\partial S$. Taking into account homogeneous Dirichlet boundary conditions, its subspace is denoted by $V^0 = \{u \in V \mid \gamma_e u = 0, (\gamma u)(0) = (\gamma u)(L) = 0\}$, endowed with the norm $\|u\|_{V^0}^2 = \int_\Omega |\nabla u(\mathbf{x})|^2 d\mathbf{x} + \int_0^L |\frac{d}{dx} \gamma u(x)|^2 dx$.

Let $\mathscr{C}(\Omega \times [0, T_f))$ be the subspace of functions $\varphi$ of $C^\infty\left(\overline{\Omega} \times [0, T_f]\right)$ vanishing at $t = T_f$ and at $\Gamma_D$ and such that $\partial_s \varphi = 0$ on $(0, L) \times \partial S$. Given $\bar{u} \in V$ and $u_{ini} \in V$, the variational formulation of the simplified coupled model amounts to find $u$ with $u - \bar{u} \in L^2\left(0, T_f; V^0\right)$ such that for all $\varphi \in \mathscr{C}(\Omega \times [0, T_f))$ one has

$$
\begin{cases}
-\int_0^{T_f} \int_\Omega \phi(\mathbf{x}) \zeta^l \mathcal{S}^l(\mathbf{x}, u(\mathbf{x}, t)) \partial_t \varphi(\mathbf{x}, t) d\mathbf{x} dt - \int_\Omega \phi \zeta^l \mathcal{S}^l(\mathbf{x}, u_{ini}(\mathbf{x})) \varphi(\mathbf{x}, 0) d\mathbf{x} \\
-\int_0^{T_f} \int_0^L |S| \zeta^g(g(\gamma u)(x, t)) \partial_t \gamma \varphi(x, t) dx dt - \int_0^L |S| \zeta^g(g(\gamma(u_{ini}))(x)) \gamma \varphi(x, 0) dx \\
+\int_0^{T_f} \int_\Omega \zeta^l \frac{k_r^l(\mathbf{x}, \mathcal{S}^l(u(\mathbf{x}, t)))}{\mu^l} \mathbf{K}(\nabla u(\mathbf{x}, t) - M^l \zeta^l \mathbf{g}) \cdot \nabla \varphi(\mathbf{x}, t) d\mathbf{x} dt \\
+\int_0^{T_f} \int_0^L \frac{1}{\alpha(x)} |S| \widetilde{\zeta}^g(g(\gamma u)(x, t)) \partial_x g(\gamma u)(x, t) \partial_x \gamma \varphi(x, t) dx dt = 0.
\end{cases}
\tag{6}
$$

We make the following additional *assumptions on the data*:

- It is assumed that $k_r^l(\mathbf{x}, s)$ is a measurable function w.r.t. $\mathbf{x}$ and continuous w.r.t. $s$, and such that $0 < k_{min} \leq k_r^l(\mathbf{x}, s) \leq k_{max}$ for all $(\mathbf{x}, s) \in \Omega \times [0, 1]$.
- $\mathcal{S}^l(\mathbf{x}, u) \in [0, 1]$ for all $(\mathbf{x}, u) \in \Omega \times \mathbb{R}$ with $\mathcal{S}^l(\mathbf{x}, u) = \mathcal{S}_j^l(u)$ for a.e. $\mathbf{x} \in \Omega_j$ and all $u \in \mathbb{R}$, where $\mathcal{S}_j^l$ is a non decreasing Lipschitz continuous function with constant $L_S$ and $(\Omega_j)_{j \in J}$ is a finite family of disjoint connected polyhedral open sets such that $\bigcup_{j \in J} \overline{\Omega}_j = \overline{\Omega}$.
- The permeability tensor $\mathbf{K}$ is a measurable function on the space of symmetric 3 dimensional matrices such that there exist $0 < \lambda_{min} \leq \lambda_{max}$ with $\lambda_{min}|\xi|^2 \leq (\mathbf{K}(\mathbf{x})\xi, \xi) \leq \lambda_{max}|\xi|^2$ for all $\mathbf{x} \in \overline{\Omega}$.
- $\alpha \in L^\infty(0, L)$ is such that $0 < \alpha_{min} \leq \alpha(x) \leq \alpha_{max}$ for all $x \in (0, L)$.
- The porosity $\phi$ belongs to $L^\infty(\Omega)$ with $0 < \phi_{min} \leq \phi(\mathbf{x}) \leq \phi_{max}$ for all $\mathbf{x} \in \Omega$.

**Vertex Approximate Gradient (VAG) discretization**: We assume that $\omega$ and $S$ are polygonal domains of $\mathbb{R}^2$ and we consider a conforming polyhedral mesh of the domain $\Omega$. It is assumed that the intersection of the mesh with the boundary $\Gamma$ of

the gallery is the tensor product of the 1D mesh of $(0, L)$ defined by $0 = x_0 < x_1 < \cdots < x_{n_x+1} = L$ by the 1D mesh of $\partial S$ defined by the set of distinct points $\mathbf{s}_1, \mathbf{s}_2 \cdots, \mathbf{s}_{n_S}, \mathbf{s}_{n_S+1} = \mathbf{s}_1$ of $\partial S$ in cyclic order.

Let $\mathcal{M}$ denote the set of cells $K$, $\mathcal{V}$ the set of nodes $\mathbf{s}$, $\mathcal{E}$ the set of edges $e$, and $\mathcal{F}$ the set of faces $\sigma$, of the mesh. We denote by $\mathcal{V}_K$ the set of nodes of the cell $K \in \mathcal{M}$, by $\mathcal{V}_\sigma$ the set of nodes and by $\mathcal{E}_\sigma$ the set of edges of the face $\sigma \in \mathcal{F}$. The set of nodes of the mesh belonging to $\{x_i\} \times \partial S$ is denoted by $\mathcal{V}_i$ for all $i = 0, \cdots, n_x + 1$.

It is assumed for each face $\sigma \in \mathcal{F}$, that there exists a so-called "centre" of the face $\mathbf{x}_\sigma$ such that $\mathbf{x}_\sigma = \sum_{\mathbf{s} \in \mathcal{V}_\sigma} \beta_{\sigma,\mathbf{s}} \mathbf{x}_\mathbf{s}$, with $\sum_{\mathbf{s} \in \mathcal{V}_\sigma} \beta_{\sigma,\mathbf{s}} = 1$, where $\beta_{\sigma,\mathbf{s}} \geq 0$ for all $\mathbf{s} \in \mathcal{V}_\sigma$. The face $\sigma$ is assumed to be star-shaped w.r.t. its centre $\mathbf{x}_\sigma$ which means that the face $\sigma$ matches with the union of the triangles $\tau_{\sigma,e}$ defined by the face centre $\mathbf{x}_\sigma$ and each of its edge $e \in \mathcal{E}_\sigma$.

The previous discretization is denoted by $\mathcal{D}$, and we define the discrete space

$$X_\mathcal{D} = \{v_K \in \mathbb{R}, v_\mathbf{s} \in \mathbb{R}, v_i \in \mathbb{R}, K \in \mathcal{M}, \mathbf{s} \in \mathcal{V}, i = 0, \cdots, n_x + 1 \\ \mid v_\mathbf{s} = v_i \text{ for all } \mathbf{s} \in \mathcal{V}_i, i = 0, \cdots, n_x + 1\}, \tag{7}$$

and its subspace with homogeneous Dirichlet boundary conditions on $\Gamma_D$ and at $x = 0$, $x = L$

$$X_\mathcal{D}^0 = \{v \in X_\mathcal{D} \mid v_\mathbf{s} = 0 \text{ for all } \mathbf{s} \in \mathcal{V}_D\},$$

where $\mathcal{V}_D = \mathcal{V} \cap \overline{\Gamma}_D$ are the Dirichlet boundary nodes.

Following [1], the extension of the VAG discretization to the coupled model (6) is based on conforming Finite Element reconstructions of the gradient operators on $\Omega$ and on $(0, L)$, and on non conforming piecewise constant function reconstructions on $\Omega$ and on $(0, L)$.

For all $\sigma \in \mathcal{F}$, let us first define the operator $I_\sigma : X_\mathcal{D} \to \mathbb{R}$ such that $I_\sigma(v) = \sum_{\mathbf{s} \in \mathcal{V}_\sigma} \beta_{\sigma,\mathbf{s}} v_\mathbf{s}$, which is by definition of $\mathbf{x}_\sigma$ a second order interpolation operator at point $\mathbf{x}_\sigma$.

Let us introduce the tetrahedral sub-mesh $\mathcal{T} = \{\tau_{K,\sigma,e}, e \in \mathcal{E}_\sigma, \sigma \in \mathcal{F}_K, K \in \mathcal{M}\}$ of the mesh $\mathcal{M}$, where $\tau_{K,\sigma,e}$ is the tetrahedron defined by the cell center $\mathbf{x}_K$ and the triangle $\tau_{\sigma,e}$. For a given $v \in X_\mathcal{D}$, we define the function $\Pi_\mathcal{T} v \in V$ as the continuous piecewise affine function on each tetrahedron $\tau$ of $\mathcal{T}$ such that $\Pi_\mathcal{T} v(\mathbf{x}_K) = v_K$, $\Pi_\mathcal{T} v(\mathbf{s}) = v_\mathbf{s}$, and $\Pi_\mathcal{T} v(\mathbf{x}_\sigma) = I_\sigma(v)$ for all $K \in \mathcal{M}$, $\mathbf{s} \in \mathcal{V}$, $\sigma \in \mathcal{F}$.

It is easily checked that $\partial_s \gamma \Pi_\mathcal{T} v = 0$ which shows that $\Pi_\mathcal{T} v \in V$ for all $v \in X_\mathcal{D}$. Then, the gradient operators are defined for all $v \in X_\mathcal{D}$ by

$$\nabla_\mathcal{D} v = \nabla \Pi_\mathcal{T} v \text{ and } \nabla_{x,\mathcal{D}} v = \partial_x \gamma \Pi_\mathcal{T} v.$$

One can easily check that $\nabla_{x,D} v = \frac{v_{i+1} - v_i}{x_{i+1} - x_i}$ on $(x_i, x_{i+1})$ for all $i = 0, \cdots, n_x$. For the reconstructions of functions operators, we first set

$$\Pi_\mathcal{D} v(\mathbf{x}) = v_K \text{ for all } \mathbf{x} \in K, \ K \in \mathcal{M}.$$

Next, let us define the points $x_{i+\frac{1}{2}} = \frac{x_i + x_{i+1}}{2}, i = 1, \cdots, n_x - 1, x_{\frac{1}{2}} = 0, x_{n_x + \frac{1}{2}} = L$, we set

$$\Pi_{x,\mathscr{D}} v(x) = v_i \text{ for all } x \in (x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}), \ i = 1, \cdots, n_x.$$

Let $\rho_\tau$ denote the insphere diameter of a given tetrahedron $\tau \in \mathscr{T}$, $h_\tau$ its diameter, $h_{\mathscr{T}} = \max_{\tau \in \mathscr{T}} h_\tau$, and $\theta_{\mathscr{T}} = \max_{\tau \in \mathscr{T}} \frac{h_\tau}{\rho_\tau}$ and $\gamma_{\mathscr{M}} = \max_{K \in \mathscr{M}} \text{Card}(\mathscr{V}_K)$. For $N \in \mathbb{N}^*$, let us consider the time discretization $t^0 = 0 < t^1 < \cdots < t^{n-1} < t^n \cdots < t^N = T_f$ of the time interval $[0, T_f]$. We denote the time steps by $\Delta t^n = t^n - t^{n-1}$ for all $n = 1, \cdots, N$. For $v \in X_{\mathscr{D}}$, and a function $k \in C^0(\mathbb{R}, \mathbb{R})$, we define $k(v) \in X_{\mathscr{D}}$ as follows: $k(v)_\mathbf{s} = k(v_\mathbf{s})$ for all $\mathbf{s} \in \mathscr{V}$, $k(v)_K = k(v_K)$ for all $K \in \mathscr{M}$, and $k(v)_i = k(v_i)$ for all $i = 0, \cdots, n_x + 1$.

Given $u_{\mathscr{D}}^0 \in X_{\mathscr{D}}$ and $\bar{u}_{\mathscr{D}} \in X_{\mathscr{D}}$, the discretization of the coupled model (6) looks for $u_{\mathscr{D}}^n \in X_{\mathscr{D}}$ with $u_{\mathscr{D}}^n - \bar{u}_{\mathscr{D}} \in X_{\mathscr{D}}^0$ for all $n = 1, \cdots, N$ such that for all $v_{\mathscr{D}} \in X_{\mathscr{D}}^0$

$$
\begin{cases}
\displaystyle\int_\Omega \phi(\mathbf{x}) \zeta^l \frac{\mathcal{S}^l(\mathbf{x}, \Pi_{\mathscr{D}} u_{\mathscr{D}}^n(\mathbf{x})) - \mathcal{S}^l(\mathbf{x}, \Pi_{\mathscr{D}} u_{\mathscr{D}}^{n-1}(\mathbf{x}))}{\Delta t^n} \Pi_{\mathscr{D}} v_{\mathscr{D}}(\mathbf{x}) d\mathbf{x} \\
\displaystyle + \int_0^L |S| \frac{\zeta^g(\Pi_{x,\mathscr{D}} g(u_{\mathscr{D}}^n)(x)) - \zeta^g(\Pi_{x,\mathscr{D}} g(u_{\mathscr{D}}^{n-1})(x))}{\Delta t^n} \Pi_{x,\mathscr{D}} v_{\mathscr{D}}(x) dx \\
\displaystyle + \int_\Omega \zeta^l \frac{k_r^l(\mathbf{x}, \mathcal{S}^l(\mathbf{x}, \Pi_{\mathscr{D}} u_{\mathscr{D}}^n(\mathbf{x})))}{\mu^l} \mathbf{K}(\nabla_{\mathscr{D}} u_{\mathscr{D}}^n(\mathbf{x}) - M^l \zeta^l \mathbf{g}) \cdot \nabla_{\mathscr{D}} v_{\mathscr{D}}(\mathbf{x}) d\mathbf{x} \\
\displaystyle + \int_0^L \frac{1}{\alpha(x)} |S| \widetilde{\zeta}^g(\Pi_{x,\mathscr{D}} g(u_{\mathscr{D}}^n)(x)) \nabla_{x,\mathscr{D}} g(u_{\mathscr{D}}^n)(x) \nabla_{x,\mathscr{D}} v_{\mathscr{D}}(x) dx = 0.
\end{cases}
\tag{8}
$$

**Convergence analysis**: Let us set $X_{\mathscr{D},\Delta t} = (X_{\mathscr{D}})^N$, and for all $v_{\mathscr{D}} = (v_{\mathscr{D}}^n)_{n=1,\cdots,N} \in X_{\mathscr{D},\Delta t}$ let us define for all $n = 1, \cdots, N$, and for all $(\mathbf{x}, t) \in \Omega \times (t^{n-1}, t^n]$ the functions $\Pi_{\mathscr{D},\Delta t} v_{\mathscr{D}}(\mathbf{x}, t) = \Pi_{\mathscr{D}} v_{\mathscr{D}}^n(\mathbf{x})$, $\Pi_{x,\mathscr{D},\Delta t} v_{\mathscr{D}}(x, t) = \Pi_{x,\mathscr{D}} v_{\mathscr{D}}^n(x)$, $\Pi_{\mathscr{T},\Delta t} v_{\mathscr{D}}(\mathbf{x}, t) = \Pi_{\mathscr{T}} v_{\mathscr{D}}^n(\mathbf{x})$. Let $u_{\mathscr{D}} = (u_{\mathscr{D}}^n)_{n=1,\cdots,N}$, the given solution to (8), we also define the functions $S_{\mathscr{D},\Delta t}^l(\mathbf{x}, t) = \mathcal{S}^l(\mathbf{x}, \Pi_{\mathscr{D},\Delta t} u_{\mathscr{D}}(\mathbf{x}, t))$, $p_{x,\mathscr{D},\Delta t}(x, t) = g(\Pi_{x,\mathscr{D},\Delta t} u_{\mathscr{D}}(x, t))$. Using similar techniques as in [2], we can prove the following convergence theorem.

**Theorem 1** *Let $\mathscr{D}^{(m)}, \Delta t^{n,(m)}, n = 1, \cdots, N^{(m)}, m \in \mathbb{N}$ be a sequence of space time discretizations such that there exist $\theta > 0$, $\gamma > 0$ with $\theta_{\mathscr{T}^{(m)}} \leq \theta$, $\gamma_{\mathscr{T}^{(m)}} \leq \gamma$. It is assumed that $\lim_{m \to +\infty} h_{\mathscr{T}^{(m)}} = 0$, and that $\Delta t^{(m)} = \max_{n=1,\cdots,N^{(m)}} \Delta t^{n,(m)}$ tends to zero when $m \to +\infty$, and that $\|\Pi_{\mathscr{D}^{(m)}} u_{\mathscr{D}^{(m)}}^0 - u_{ini}\|_{L^2(\Omega)}$, $\|\Pi_{x,\mathscr{D}^{(m)}} u_{\mathscr{D}^{(m)}}^0 - \gamma u_{ini}\|_{L^2(0,L)}$, $\|\Pi_{\mathscr{T}^{(m)}} \bar{u}_{\mathscr{D}^{(m)}} - \bar{u}\|_{V^0}$, $\|\Pi_{\mathscr{D}^{(m)}} \bar{u}_{\mathscr{D}^{(m)}} - \bar{u}\|_{L^2(\Omega)}$, $\|\Pi_{x,\mathscr{D}^{(m)}} \bar{u}_{\mathscr{D}^{(m)}} - \gamma \bar{u}\|_{L^2(0,L)}$ tends to zero when $m \to +\infty$. Then, there exist a subsequence of $m \in \mathbb{N}$ and a function $u \in L^2(0, T_f; V)$ solution of (6) such that up to this subsequence $S_{\mathscr{D}^{(m)},\Delta t^{(m)}}^l \to \mathcal{S}^l(., u)$ strongly in $L^2(\Omega \times (0, T_f))$, $\Pi_{\mathscr{D}^{(m)},\Delta t^{(m)}} u_{\mathscr{D}^{(m)}} \rightharpoonup u$ weakly in $L^2(\Omega \times (0, T_f))$, and $p_{x,\mathscr{D}^{(m)},\Delta t^{(m)}} \to g(\gamma(u))$ strongly in $L^2((0, L) \times (0, T_f))$.*

# References

1. Eymard, R., Guichard, C., Herbin, R.: Small-stencil 3d schemes for diffusive flows in porous media. ESAIM Math. Model. Numer. Anal. **46**, 265–290 (2010)
2. Eymard, R., Guichard, C., Herbin, R., Masson, R.: Gradient schemes for two-phase flow in heterogeneous porous media and richards equation. ZAMM J. Appl. Math. Mech. (2013). doi:10.1002/zamm.201200206
3. Masson, R., Trenty, L., Zhang, Y.: Formulations of two-phase liquid gas compositional darcy flows with phase transitions (2013). http://hal.archives-ouvertes.fr/hal-00910366
4. Mosthaf, K., Baber, K., Flemisch, B., Helmig, R., Leijnse, A., Rybak, I., Wohlmuth, B.: A coupling concept for two-phase compositional porous-medium and single-phase compositional free flow. Water Resour. Res. **47**(10), 16 (2011)

# Gradient Discretization of Hybrid Dimensional Darcy Flows in Fractured Porous Media

**Konstantin Brenner, Mayya Groza, Cindy Guichard, Gilles Lebeau and Roland Masson**

**Abstract** This article deals with the discretization of hybrid dimensional model of Darcy flow in fractured porous media. These models couple the flow in the fractures represented as the surfaces of codimension one with the flow in the surrounding matrix. The convergence analysis is carried out in the framework of Gradient schemes which accounts for a large family of conforming and nonconforming discretizations. The Vertex Approximate Gradient (VAG) scheme and the Hybrid Finite Volume (HFV) scheme are applied to such models and are shown to verify the Gradient scheme framework. Our theoretical results are confirmed by a few numerical experiments performed both on tetrahedral and hexahedral meshes in heterogeneous isotropic and anisotropic media.

K. Brenner · M. Groza · R. Masson
Laboratoire de Mathématiques J.A. Dieudonné UMR CNRS 7251 and Team COFFEE,
University Nice Sophia Antipolis, CNRS and INRIA Sophia Antipolis Méditerranée, Nice, France
e-mail: konstantin.brenner@unice.fr

M. Groza
e-mail: mayya.groza@unice.fr

R. Masson
e-mail: roland.masson@unice.fr

C. Guichard (✉)
Laboratoire Jacques-Louis Lions, CNRS, UMR 7598, Sorbonne Universités,
UPMC Univ Paris 06, F-75005 Paris, France
e-mail: guichard@ljll.math.upmc.fr

G. Lebeau
Laboratoire de Mathématiques J.A. Dieudonné UMR CNRS 7251, University Nice Sophia
Antipolis and CNRS, Nice, France
e-mail: gilles.lebeau@unice.fr

**Fig. 1** Example of a 2D
domain with 3 intersecting
fractures and 2 connected
components



# 1 Hybrid Dimensional Darcy Flow in Fractured Porous Media

Let $\Omega$ denote a bounded polyhedral domain of $\mathbb{R}^d$, $d = 2, 3$. We consider the asymptotic model introduced in [1] where fractures are represented as interfaces of codimension 1. Let $\overline{\Gamma} = \bigcup_{i \in I} \overline{\Gamma}_i$ denotes the network of fractures $\Gamma_i \subset \Omega$, $i \in I$, such that each $\Gamma_i$ is a planar polygonal simply connected open domain. It is assumed that the angles of $\Gamma_i$ are strictly lower than $2\pi$ and that $\Gamma_i \cap \Gamma_j = \emptyset$ for all $i \neq j$. For all $i \in I$, let us set $\Sigma_i = \partial \Gamma_i$, $\Sigma_{i,j} = \Sigma_i \cap \Sigma_j$, $j \in I$, $\Sigma_{i,0} = \Sigma_i \cap \partial \Omega$, $\Sigma_{i,N} = \Sigma_i \setminus (\bigcup_{j \in I} \Sigma_{i,j} \cup \Sigma_{i,0})$, and $\Sigma = \bigcup_{(i,j) \in I \times I, i \neq j} \Sigma_{i,j}$ (Fig. 1). It is assumed that $\Sigma_{i,0} = \overline{\Gamma}_i \cap \partial \Omega$, and that $\bigcup_{i \in I} \Gamma_i = \Gamma \setminus \Sigma$. We will denote by $d\tau(\mathbf{x})$ the $d-1$ dimensional Lebesgue measure on $\Gamma$. Let $H^1(\Gamma)$ denote the set of functions $v = (v_i)_{i \in I}$ such that $v_i \in H^1(\Gamma_i), i \in I$ with continuous traces at the fracture intersections, and endowed with the norm $\|v\|_{H^1(\Gamma)}^2 = \sum_{i \in I} \|v_i\|_{H^1(\Gamma_i)}^2$. Its subspace with vanishing traces on $\Sigma_0 = \bigcup_{i \in I} \Sigma_{i,0}$ is denoted by $H_{\Sigma_0}^1(\Gamma)$. The gradient operator from $H^1(\Omega)$ to $L^2(\Omega)^d$ is denoted by $\nabla$, and the tangential gradient from $H^1(\Gamma)$ to $L^2(\Gamma)^{d-1}$ by $\nabla_\tau$. Let us also consider the trace operator $\gamma$ from $H^1(\Omega)$ to $L^2(\Gamma)$. The function spaces used in the variational formulation of the hybrid dimensional Darcy flow model are defined by

$$V = \{v \in H^1(\Omega), \ \gamma v \in H^1(\Gamma)\}, \text{ and its subspace}$$
$$V_0 = \{v \in H_0^1(\Omega), \ \gamma v \in H_{\Sigma_0}^1(\Gamma)\}.$$

The space $V_0$ is endowed with the norm $\|v\|_{V_0}^2 = \|\nabla v\|_{L^2(\Omega)^d}^2 + \|\nabla_\tau \gamma v\|_{L^2(\Gamma)^{d-1}}^2$ and the space $V$ with the norm $\|v\|_V^2 = \|v\|_{V_0}^2 + \|v\|_{L^2(\Omega)}^2$. Let $\Omega_\alpha, \alpha \in \Xi$ denote the connected components of $\Omega \setminus \overline{\Gamma}$, and let us define the space $H_{\mathrm{div}}(\Omega \setminus \overline{\Gamma}) = \{\mathbf{q}_m = (\mathbf{q}_{m,\alpha})_{\alpha \in \Xi} \mid \mathbf{q}_{m,\alpha} \in H_{\mathrm{div}}(\Omega_\alpha)\}$. For all $i \in I$, we can define the two sides $\pm$ of the fracture $\Gamma_i$ and the corresponding unit normal vector $n_i^\pm$ at $\Gamma_i$ outward to the sides $\pm$. For all $\mathbf{q}_m \in H_{\mathrm{div}}(\Omega \setminus \overline{\Gamma})$, let $\mathbf{q}_m^\pm \cdot n_i^\pm|_{\Gamma_i}$ denote the two normal traces at the fracture $\Gamma_i$ and let us define the jump operator $H_{\mathrm{div}}(\Omega \setminus \overline{\Gamma}) \to (H_{00}^{1/2}(\Gamma_i))'$ by $[\![\mathbf{q}_m \cdot \mathbf{n}_i]\!] = \mathbf{q}_m^+ \cdot n_i^+|_{\Gamma_i} + \mathbf{q}_m^- \cdot n_i^-|_{\Gamma_i}$. For all fractures $\Gamma_i, i \in I$, we denote by $\mathbf{n}_{\Sigma_i}$ the unit vector normal to $\Sigma_i$ outward to $\Gamma_i$.

**Hybrid Dimensional Darcy Flow Model**: In the matrix domain $\Omega \setminus \overline{\Gamma}$ (resp. in the fracture network $\Gamma$), let us denote by $\Lambda_m(\mathbf{x})$ (resp. $\Lambda_f(\mathbf{x})$) the permeability

tensor. We also denote by $d_f(\mathbf{x}), \mathbf{x} \in \Gamma$ the width of the fractures, and by $d\tau_f(\mathbf{x})$ the weighted Lebesgue $d-1$ dimensional measure on $\Gamma$ defined by $d\tau_f(\mathbf{x}) = d_f(\mathbf{x})d\tau(\mathbf{x})$. We consider the source terms $h_m \in L^2(\Omega)$ (resp. $h_f \in L^2(\Gamma)$) in the matrix domain $\Omega \setminus \overline{\Gamma}$ (resp. in the fracture network $\Gamma$). The strong formulation of the model amounts to find $u \in V_0$, $(\mathbf{q}_m, \mathbf{q}_f) \in W(\Omega, \Gamma)$ such that

$$
\begin{cases}
\operatorname{div}(\mathbf{q}_{m,\alpha}) = h_m & \text{on } \Omega_\alpha, \alpha \in \Xi, \\
\mathbf{q}_{m,\alpha} = -\Lambda_m \nabla u & \text{on } \Omega_\alpha, \alpha \in \Xi, \\
\operatorname{div}_\tau(\mathbf{q}_{f,i}) - [\![\mathbf{q}_m \cdot \mathbf{n}_i]\!] = d_f h_f & \text{on } \Gamma_i, i \in I, \\
\mathbf{q}_{f,i} = -d_f \Lambda_f \nabla_\tau \gamma u & \text{on } \Gamma_i, i \in I,
\end{cases}
\tag{1}
$$

where the function space $W(\Omega, \Gamma)$ is defined by

$$
\begin{aligned}
W(\Omega, \Gamma) = \{ \ & \mathbf{q}_m = (\mathbf{q}_{m,\alpha})_{\alpha \in \Xi}, \mathbf{q}_f = (\mathbf{q}_{f,i})_{i \in I} \mid \mathbf{q}_m \in H_{\operatorname{div}}(\Omega \setminus \overline{\Gamma}), \\
& \mathbf{q}_{f,i} \in L^2(\Gamma_i)^{d-1}, r_{f,i} = \operatorname{div}_\tau(\mathbf{q}_{f,i}) - [\![\mathbf{q}_m \cdot \mathbf{n}_i]\!] \in L^2(\Gamma_i), i \in I, \\
& \sum_{\alpha \in \Xi} \int_{\Omega_\alpha} (\mathbf{q}_{m,\alpha} \cdot \nabla v + \operatorname{div}(\mathbf{q}_{m,\alpha})v)d\mathbf{x} \\
& + \sum_{i \in I} \int_{\Gamma_i} (\mathbf{q}_{f,i} \cdot \nabla_\tau \gamma v + r_{f,i}\gamma v)d\tau = 0 \text{ for all } v \in V_0\}.
\end{aligned}
$$

The last condition corresponds to impose in a weak sense that $\sum_{i \in I} \mathbf{q}_{f,i} \cdot \mathbf{n}_{\Sigma_i} = 0$ on $\Sigma$ and $\mathbf{q}_{f,i} \cdot \mathbf{n}_{\Sigma_i} = 0$ on $\Sigma_{i,N}, i \in I$.

In variational form, (1) amounts to find $u \in V_0$ such that for all $v \in V_0$:

$$
\begin{cases}
\displaystyle \int_\Omega \Lambda_m(\mathbf{x})\nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x})d\mathbf{x} + \int_\Gamma \Lambda_f(\mathbf{x})\nabla_\tau \gamma u(\mathbf{x}) \cdot \nabla_\tau \gamma v(\mathbf{x})d\tau_f(\mathbf{x}) \\
\displaystyle - \int_\Omega h_m(\mathbf{x})v(\mathbf{x})d\mathbf{x} - \int_\Gamma h_f(\mathbf{x})\gamma v(\mathbf{x})d\tau_f(\mathbf{x}) = 0.
\end{cases}
\tag{2}
$$

**Proposition 1** *From the Lax-Milgram theorem, the variational problem* (2) *has a unique solution $u \in V_0$ which satisfies the a priori estimate*
$\|u\|_V \leq C\left(\|h_m\|_{L^2(\Omega)} + \|h_f\|_{L^2(\Gamma)}\right)$, *with C depending only on $\Omega$, $\Gamma$, $\Lambda_m$, $\Lambda_f$, $d_f$. In addition $(\mathbf{q}_m = -\Lambda_m \nabla u, \mathbf{q}_f = -d_f \Lambda_f \nabla_\tau \gamma u)$ belongs to $W(\Omega, \Gamma)$.*

## 2 Gradient Discretization

A gradient discretization $\mathcal{D}$ of (2) is defined by a vector space of degrees of freedom $X_\mathcal{D}$, its subspace associated with homogeneous Dirichlet boundary conditions $X_\mathcal{D}^0$, and the following set of linear operators:

- Gradient operator on the matrix domain: $\nabla_{\mathcal{D}_m} : X_\mathcal{D} \to L^2(\Omega)^d$
- Gradient operator on the fracture network: $\nabla_{\mathcal{D}_f} : X_\mathcal{D} \to L^2(\Gamma)^{d-1}$

- A function reconstruction operator on the matrix domain: $\Pi_{\mathscr{D}_m} : X_{\mathscr{D}} \to L^2(\Omega)$
- A function reconstruction operator on the fracture network: $\Pi_{\mathscr{D}_f} : X_{\mathscr{D}} \to L^2(\Gamma)$.

$X_{\mathscr{D}}$ is endowed with the semi-norm $\|v_{\mathscr{D}}\|_{\mathscr{D}}^2 = \|\nabla_{\mathscr{D}_m} v_{\mathscr{D}}\|_{L^2(\Omega)^d}^2 + \|\nabla_{\mathscr{D}_f} v_{\mathscr{D}}\|_{L^2(\Gamma)^{d-1}}^2$ which is assumed to define a norm on $X_{\mathscr{D}}^0$. Next, we define the coercivity, consistency, limit conformity and compactness properties of the gradient discretization.

**Coercivity**: There exists $C_{\mathscr{D}} \geq 0$ such that for all $v \in X_{\mathscr{D}}^0$ one has

$$\|\Pi_{\mathscr{D}_m} v_{\mathscr{D}}\|_{L^2(\Omega)} + \|\Pi_{\mathscr{D}_f} v_{\mathscr{D}}\|_{L^2(\Gamma)} \leq C_{\mathscr{D}} \|v_{\mathscr{D}}\|_{\mathscr{D}}.$$

**Consistency**: Let $u \in V_0$, and let us define

$$\mathscr{S}_{\mathscr{D}}(u) = \inf_{v_{\mathscr{D}} \in X_{\mathscr{D}}^0} \Big( \|\nabla_{\mathscr{D}_m} v_{\mathscr{D}} - \nabla u\|_{L^2(\Omega)^d} + \|\nabla_{\mathscr{D}_f} v_{\mathscr{D}} - \nabla_\tau \gamma u\|_{L^2(\Gamma)^{d-1}}$$
$$+ \|\Pi_{\mathscr{D}_m} v_{\mathscr{D}} - u\|_{L^2(\Omega)} + \|\Pi_{\mathscr{D}_f} v_{\mathscr{D}} - \gamma u\|_{L^2(\Gamma)} \Big)$$

Then, a sequence of gradient discretizations $(\mathscr{D}^l)_{l \in \mathbb{N}}$ is said to be consistent if for all $u \in V_0$ one has $\lim_{l \to +\infty} \mathscr{S}_{\mathscr{D}^l}(u) = 0$.

**Limit Conformity**: For all $(\mathbf{q}_m, \mathbf{q}_f) \in W(\Omega, \Gamma)$, we define

$$\mathscr{W}_{\mathscr{D}}(\mathbf{q}_m, \mathbf{q}_f) = \sup_{0 \neq v_{\mathscr{D}} \in X_{\mathscr{D}}^0} \frac{1}{\|v_{\mathscr{D}}\|_{\mathscr{D}}} \Big( \sum_{\alpha \in \Xi} \int_{\Omega_\alpha} (\nabla_{\mathscr{D}_m} v_{\mathscr{D}} \cdot \mathbf{q}_{m,\alpha} + (\Pi_{\mathscr{D}_m} v_{\mathscr{D}}) \mathrm{div}(\mathbf{q}_{m,\alpha}))(\mathbf{x}) d\mathbf{x}$$
$$+ \sum_{i \in I} \int_{\Gamma_i} (\nabla_{\mathscr{D}_f} v_{\mathscr{D}} \cdot \mathbf{q}_f + \Pi_{\mathscr{D}_f} v_{\mathscr{D}} (\mathrm{div}_{\tau_i}(\mathbf{q}_{f,i}) - [\![\mathbf{q}_m \cdot \mathbf{n}_i]\!]))(\mathbf{x}) d\tau(\mathbf{x}) \Big). \tag{3}$$

Then, a sequence of gradient discretizations $(\mathscr{D}^l)_{l \in \mathbb{N}}$ is said to be limit conforming if for all $(\mathbf{q}_m, \mathbf{q}_f) \in W(\Omega, \Gamma)$ one has $\lim_{l \to +\infty} \mathscr{W}_{\mathscr{D}^l}(\mathbf{q}_m, \mathbf{q}_f) = 0$.

**Compactness**: A sequence of gradient discretizations $(\mathscr{D}^l)_{l \in \mathbb{N}}$ is said to be compact if for all sequences $v_{\mathscr{D}^l} \in X_{\mathscr{D}^l}^0, l \in \mathbb{N}$ such that there exists $C > 0$ with $\|v_{\mathscr{D}^l}\|_{\mathscr{D}^l} \leq C$ for all $l \in \mathbb{N}$, then there exist $u_m \in L^2(\Omega)$ and $u_f \in L^2(\Gamma)$ with

$$\lim_{l \to +\infty} \|\Pi_{\mathscr{D}_m^l} v_{\mathscr{D}^l} - u_m\|_{L^2(\Omega)} = 0 \text{ and } \lim_{l \to +\infty} \|\Pi_{\mathscr{D}_f^l} v_{\mathscr{D}^l} - u_f\|_{L^2(\Gamma)} = 0.$$

The discretization of (2) using the Gradient Scheme framework is defined by: find $u \in X_{\mathscr{D}}^0$ such that for all $v_{\mathscr{D}} \in X_{\mathscr{D}}^0$:

$$\begin{cases} \displaystyle\int_{\Omega} \Lambda_m(\mathbf{x}) \nabla_{\mathscr{D}_m} u_{\mathscr{D}}(\mathbf{x}) \cdot \nabla_{\mathscr{D}_m} v_{\mathscr{D}}(\mathbf{x}) d\mathbf{x} + \int_{\Gamma} \Lambda_f(\mathbf{x}) \nabla_{\mathscr{D}_f} u_{\mathscr{D}}(\mathbf{x}) \cdot \nabla_{\mathscr{D}_f} v_{\mathscr{D}}(\mathbf{x}) d\tau_f(\mathbf{x}) \\ \displaystyle- \int_{\Omega} h_m(\mathbf{x}) \Pi_{\mathscr{D}_m} v_{\mathscr{D}}(\mathbf{x}) d\mathbf{x} - \int_{\Gamma} h_f(\mathbf{x}) \Pi_{\mathscr{D}_f} v_{\mathscr{D}}(\mathbf{x}) d\tau_f(\mathbf{x}) = 0. \end{cases} \tag{4}$$

**Proposition 2** *Let $\mathscr{D}$ be a gradient discretization of (2) assumed to be coercive. Then (4) has a unique solution $u_{\mathscr{D}} \in X_{\mathscr{D}}^0$ satisfying the a priori estimate $\|u_{\mathscr{D}}\|_{\mathscr{D}} \leq C\left(\|h_m\|_{L^2(\Omega)} + \|h_f\|_{L^2(\Gamma)}\right)$ with $C$ depending only on $C_{\mathscr{D}}$, $\Lambda_m$, $\Lambda_f$, $d_f$.*

**Proposition 3 Error Estimates**. *Let $u \in V_0$, $(\mathbf{q}_m, \mathbf{q}_f) \in W(\Omega, \Gamma)$ be the solution of (2). Let $\mathscr{D}$ be a gradient discretization of (2) assumed to be coercive, and let $u_{\mathscr{D}} \in X_{\mathscr{D}}^0$ be the solution of (4). Then, there exist $C_1$, $C_2$, $C_3$, $C_4$ depending only on $C_{\mathscr{D}}$, $\Lambda_m$, $\Lambda_f$, $d_f$ such that one has the following error estimates:*

$$\begin{cases} \|\nabla u - \nabla_{\mathscr{D}_m} u_{\mathscr{D}}\|_{L^2(\Omega)^d} + \|\nabla_\tau \gamma u - \nabla_{\mathscr{D}_f} u_{\mathscr{D}}\|_{L^2(\Gamma)^{d-1}} \leq C_1 \mathscr{S}_{\mathscr{D}}(u) + C_2 \mathscr{W}(\mathbf{q}_m, \mathbf{q}_f), \\ \|\Pi_{\mathscr{D}_m} u_{\mathscr{D}} - u\|_{L^2(\Omega)} + \|\Pi_{\mathscr{D}_f} u_{\mathscr{D}} - \gamma u\|_{L^2(\Gamma)} \leq C_3 \mathscr{S}_{\mathscr{D}}(u) + C_4 \mathscr{W}(\mathbf{q}_m, \mathbf{q}_f). \end{cases}$$

# 3 Two Examples of Gradient Discretizations of Hybrid Dimensional Models

In the spirit of [3], we consider generalized polyhedral meshes of $\Omega$. Let $\mathscr{M}$ be the set of cells that are disjoint open polyhedral subsets of $\Omega$ such that $\bigcup_{K \in \mathscr{M}} \overline{K} = \overline{\Omega}$. For all $K \in \mathscr{M}$, $\mathbf{x}_K$ denotes the so-called "centre" of the cell $K$ under the assumption that $K$ is star-shaped with respect to $\mathbf{x}_K$. We then denote by $\mathscr{F}_K$ the set of interfaces of non zero $d-1$ dimensional measure among the interior faces $\overline{K} \cap \overline{L}$, $L \in \mathscr{M}$, and the boundary interface $\overline{K} \cap \partial \Omega$, which possibly splits in several boundary faces. Let us denote by $\mathscr{F} = \bigcup_{K \in \mathscr{M}} \mathscr{F}_K$ the set of all faces of the mesh. The term "generalized polyhedral mesh" means that the faces are not assumed to be planar. For $\sigma \in \mathscr{F}$, let $\mathscr{E}_\sigma$ be the set of interfaces of non zero $d-2$ dimensional measure among the interfaces $\sigma \cap \sigma'$, $\sigma' \in \mathscr{F}$. Then, we denote by $\mathscr{E} = \bigcup_{\sigma \in \mathscr{F}} \mathscr{E}_\sigma$ the set of all edges of the mesh. Let $\mathscr{V}_\sigma = \bigcup_{e,e' \in \mathscr{E}_\sigma, e \neq e'} (e \cap e')$ be the set of vertices of $\sigma$, for each $K \in \mathscr{M}$ we define $\mathscr{V}_K = \bigcup_{\sigma \in \mathscr{F}_K} \mathscr{V}_\sigma$, and we also denote by $\mathscr{V} = \bigcup_{K \in \mathscr{M}} \mathscr{V}_K$ the set of all vertices of the mesh. It is then assumed that for each face $\sigma \in \mathscr{F}$, there exists a so-called "centre" of the face $\mathbf{x}_\sigma \in \sigma \setminus \bigcup_{e \in \mathscr{E}_\sigma} e$ such that $\mathbf{x}_\sigma = \sum_{\mathbf{s} \in \mathscr{V}_\sigma} \beta_{\sigma,\mathbf{s}} \mathbf{x}_\mathbf{s}$, with $\sum_{\mathbf{s} \in \mathscr{V}_\sigma} \beta_{\sigma,\mathbf{s}} = 1$, and $\beta_{\sigma,\mathbf{s}} \geq 0$ for all $\mathbf{s} \in \mathscr{V}_\sigma$; moreover the face $\sigma$ is assumed to match with the union of the triangles $T_{\sigma,e}$ defined by the face centre $\mathbf{x}_\sigma$ and each edge $e \in \mathscr{E}_\sigma$. The mesh is also supposed to be conforming w.r.t. the fracture network $\Gamma$ in the sense that for all $i \in I$ there exist the subsets $\mathscr{F}_{\Gamma_i}$ of $\mathscr{F}$ such that $\overline{\Gamma}_i = \bigcup_{\sigma \in \mathscr{F}_{\Gamma_i}} \sigma$. We will denote by $\mathscr{F}_\Gamma$ the set of fracture faces $\bigcup_{i \in I} \mathscr{F}_{\Gamma_i}$.

The discretization of the hybrid dimensional Darcy flow model with continuous pressures has been the object of several works such as [6] using a cell centred Multi-Point Flux Approximation scheme, [1] using a Mixed Finite Element (MFE) method, and [5] using a Control Volume Finite Element Method (CVFE). The MFE method, as well as some CVFE and MPFA schemes on e.g. tetrahedral meshes can be shown to be gradient discretizations. In the following we propose to apply the VAG and HFV schemes.

**Vertex Approximate Gradient Discretization**: The VAG discretization has been introduced in [3] for diffusive problems on heterogeneous anisotropic media. Its extension to the hybrid dimensional two-phase Darcy flow model is presented in [2]. The scheme is based on the following vector space of degrees of freedom:

$$X_{\mathscr{D}} = \{u_K, u_{\mathbf{s}}, u_\sigma \in \mathbb{R} \text{ for all } K \in \mathscr{M}, \mathbf{s} \in \mathscr{V}, \sigma \in \mathscr{F}_\Gamma\},$$

and its subspace with homogeneous Dirichlet boundary conditions on $\partial\Omega$: $X_{\mathscr{D}}^0 = \{u \in X_{\mathscr{D}} \,|\, u_{\mathbf{s}} = 0 \text{ for } \mathbf{s} \in \mathscr{V}_{ext}\}$ where $\mathscr{V}_{ext} = \mathscr{V} \cap \partial\Omega$ denotes the set of boundary vertices, and $\mathscr{V}_{int} = \mathscr{V} \setminus \mathscr{V}_{ext}$ denotes the set of interior vertices.

The discrete gradients in the matrix and in the fracture are defined as the usual gradient operators on the conforming space of continuous affine finite elements built upon a tetrahedral sub-mesh. In addition, the VAG discretization uses two non conforming piecewise constant reconstructions of functions from $X_{\mathscr{D}}$ into respectively $L^2(\Omega)$ and $L^2(\Gamma)$. In the matrix, it is such that $\pi_{\mathscr{D}_m} u(\mathbf{x})|_{\Omega_{m,\nu}} = u_\nu$ where the $\Omega_{m,\nu}$ for $\nu \in \mathscr{M} \cup \mathscr{V}_{int} \cup \mathscr{F}_\Gamma$ are neighbourhoods of $\mathbf{x}_\nu$ defining a partition of $\Omega$. In the fractures, it is such that $\pi_{\mathscr{D}_f} u(\mathbf{x})|_{\Omega_{f,\nu}} = u_\nu$ where the $\Omega_{f,\nu}$ for $\nu \in (\mathscr{V}_\Gamma \cap \mathscr{V}_{int}) \cup \mathscr{F}_\Gamma$ are neighbourhoods of $\mathbf{x}_\nu$ defining a partition of $\Gamma$.

**Hybrid Finite Volume Discretization**: The Hybrid Finite Volume (HFV) scheme introduced in [4] can be extended to the hybrid dimensional Darcy flow model as follows. The faces $\sigma \in \mathscr{F}$ are assumed to be planar and $\mathbf{x}_\sigma$ is assumed to be the centre of gravity of the face $\sigma$. We also denote by $\mathbf{x}_e$ the centre of the edge $e \in \mathscr{E}$. Let $\mathscr{F}_{int} \subset \mathscr{F}$ (resp. $\mathscr{E}_{int} \subset \mathscr{E}$) denote the subset of interior faces (resp. interior edges). The vector space of degrees of freedom $X_{\mathscr{D}}$ is defined by

$$X_{\mathscr{D}} = \{u_K, u_\sigma, u_e \in \mathbb{R} \text{ for all } K \in \mathscr{M}, \sigma \in \mathscr{F}, e \in \mathscr{E}_\Gamma\},$$

where $\mathscr{E}_\Gamma \subset \mathscr{E}$ denotes the subset of edges of $\Gamma$, and its subspace $X_{\mathscr{D}}^0$ is such that $u_\sigma = 0$ for all $\sigma \in \mathscr{F} \setminus \mathscr{F}_{int}$ and $u_e = 0$ for all $e \in \mathscr{E}_\Gamma \setminus \mathscr{E}_{int}$. For each cell $K$ and $u \in X_{\mathscr{D}}$, let us define $\nabla_K u = \frac{1}{|K|} \sum_{\sigma \in \mathscr{F}_K} |\sigma|(u_\sigma - u_K)\mathbf{n}_{K,\sigma}$, where $|K|$ is the volume of the cell $K$, $|\sigma|$ is the surface of the face $\sigma$, and $\mathbf{n}_{K,\sigma}$ is the unit normal vector of the face $\sigma \in \mathscr{F}_K$ outward to the cell $K$. The discrete gradient $\nabla_K u$ is stabilized using $\nabla_{K,\sigma} u = \nabla_K u + R_{K,\sigma}(u)\mathbf{n}_{K,\sigma}$, $\sigma \in \mathscr{F}_K$, with $R_{K,\sigma}(u) = \frac{\sqrt{d}}{d_{K,\sigma}}\left(u_\sigma - u_K - \nabla_K u \cdot (\mathbf{x}_K - \mathbf{x}_\sigma)\right)$, and $d_{K,\sigma} = \mathbf{n}_{K,\sigma} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K)$ which leads to the definition of the matrix discrete gradient $\nabla_{\mathscr{D}_m} u(\mathbf{x}) = \nabla_{K,\sigma} u$ on $K_\sigma$ for all $K \in \mathscr{M}, \sigma \in \mathscr{F}_K$, where $K_\sigma$ is the cone joining the face $\sigma$ to the cell centre $\mathbf{x}_K$. The fracture discrete gradient is defined similarly by $\nabla_{\mathscr{D}_f} u(\mathbf{x}) = \nabla_{\sigma,e} u$ on $\sigma_e$ for all $\sigma \in \mathscr{F}_\Gamma, e \in \mathscr{E}_\sigma$, with $\nabla_{\sigma,e} u = \nabla_\sigma u + R_{\sigma,e}(u)\mathbf{n}_{\sigma,e}$, and $\nabla_\sigma u = \frac{1}{|\sigma|} \sum_{e \in \mathscr{E}_\sigma} |e|(u_e - u_\sigma)\mathbf{n}_{\sigma,e}$, $R_{\sigma,e}(u) = \frac{\sqrt{d-1}}{d_{\sigma,e}}\left(u_e - u_\sigma - \nabla_\sigma u \cdot (\mathbf{x}_\sigma - \mathbf{x}_e)\right)$, where $\mathbf{n}_{\sigma,e}$ is the unit normal vector to the edge $e$ in the tangent plane of the face $\sigma$ and outward to the face $\sigma$, $d_{\sigma,e} = \mathbf{n}_{\sigma,e} \cdot (\mathbf{x}_e - \mathbf{x}_\sigma)$, and $\sigma_e$ is the triangle of base $e$ and vertex $\mathbf{x}_\sigma$. The function reconstruction operators are piecewise constant on a partition of the cells and of the fracture faces. These partitions are respectively denoted,

for all $K \in \mathcal{M}$, by $K = \Omega_K \cup \left( \bigcup_{\sigma \in \mathcal{F}_K \cap \mathcal{F}_{int}} \Omega_{K,\sigma} \right)$, and, for all $\sigma \in \mathcal{F}_\Gamma$, by $\sigma = \Sigma_\sigma \cup \left( \bigcup_{e \in \mathcal{E}_\sigma \cap \mathcal{E}_{int}} \Sigma_{\sigma,e} \right)$. Then, the function reconstruction operators are defined by $\Pi_{\mathcal{D}_m} u(\mathbf{x}) = \begin{cases} u_K & \text{for all } \mathbf{x} \in \Omega_K, \ K \in \mathcal{M}, \\ u_\sigma & \text{for all } \mathbf{x} \in \Omega_{K,\sigma}, \ \sigma \in \mathcal{F}_K \cap \mathcal{F}_{int}, \ K \in \mathcal{M}, \end{cases}$ and

$\Pi_{\mathcal{D}_f} u(\mathbf{x}) = \begin{cases} u_\sigma & \text{for all } \mathbf{x} \in \Sigma_\sigma, \ \sigma \in \mathcal{F}_\Gamma, \\ u_e & \text{for all } \mathbf{x} \in \Sigma_{\sigma,e}, \ e \in \mathcal{E}_\sigma \cap \mathcal{E}_{int}, \ \sigma \in \mathcal{F}_\Gamma. \end{cases}$

We can show the following proposition which can be proven using a lemma stating the density of smooth function subspaces in the spaces $V$, $V_0$, and $W(\Omega, \Gamma)$.

**Proposition 4** *Let us consider a family of meshes $\mathcal{M}^{(m)}$, $m \in \mathbb{N}$ as defined above. It is assumed that the family of tetrahedral submeshes of $\mathcal{M}^{(m)}$ is shape regular, that the cardinal of $\mathcal{V}_K$ is uniformly bounded for all $K \in \mathcal{M}^{(m)}$, and all $m \in \mathbb{N}$, and that the maximum diameter $h^{(m)}$ of the cells $K \in \mathcal{M}^{(m)}$ tends to zero with $m \to +\infty$. In addition, in the case of the HFV scheme, the faces are assumed to be planar. Then, the VAG and HFV discretizations are coercive, consistent, limit conforming and compact gradient discretizations of the hybrid dimensional Darcy flow model.*

# 4 Numerical Experiments

Let $\Omega = (0, 1)^3$ and consider the 2 planar fractures defined by $x = 0.5$ and $y = 0.5$ and splitting $\Omega$ into the four subdomains $\Omega_\alpha$, $\alpha = 1, \cdots, 4$ corresponding respectively to $\{x < 0.5, y < 0.5\}$, $\{x > 0.5, y < 0.5\}$, $\{x > 0.5, y > 0.5\}$ and $\{x < 0.5, y > 0.5\}$. In the fractures, we set $\Lambda_f(\mathbf{x}) = 100 I$ and $d_f(\mathbf{x}) = 0.01$. In the matrix, the permeability tensor $\Lambda_m(\mathbf{x})$ is fixed to $\Lambda_{m,\alpha}$ on each subdomain $\Omega_\alpha$, $\alpha = 1, \cdots, 4$ with two choices of the subdomain permeabilities. The first choice considers isotropic heterogeneous permeabilities setting $\Lambda_{m,\alpha} = \lambda_\alpha I$ with $\lambda_1 = 1$, $\lambda_2 = 0.1, \lambda_3 = 0.01, \lambda_4 = 10$. The second choice defines anisotropic heterogeneous permeabilities by

$$\Lambda_{m,1} = \begin{pmatrix} a_1 & b_1 & 0 \\ b_1 & c_1 & 0 \\ 0 & 0 & \lambda \end{pmatrix}, \ \Lambda_{m,2} = \begin{pmatrix} a_2 & 0 & b_2 \\ 0 & \lambda & 0 \\ b_2 & 0 & c_2 \end{pmatrix}, \ \Lambda_{m,3} = \begin{pmatrix} a_3 & b_3 & 0 \\ b_3 & c_3 & 0 \\ 0 & 0 & \lambda \end{pmatrix}, \ \Lambda_{m,4} = \begin{pmatrix} \lambda & 0 & 0 \\ 0 & a_4 & b_4 \\ 0 & b_4 & c_4 \end{pmatrix},$$

with $a_\alpha = \cos^2 \beta_\alpha + \Omega \sin^2 \beta_\alpha$, $b_\alpha = (1 - \Omega) \cos \beta_\alpha \sin \beta_\alpha$, $c_\alpha = \Omega \cos^2 \beta_\alpha + \sin^2 \beta_\alpha$, $\lambda = 0.01$, $\beta_1 = \frac{\pi}{6}$, $\beta_2 = -\frac{\pi}{6}$, $\beta_3 = 0$, $\beta_4 = \frac{\pi}{6}$ and $\Omega = 0.01$. For each subdomain let us define $t_1(\mathbf{x}) = y - x + z$, $t_2(\mathbf{x}) = x + y + z - 1$, $t_3(\mathbf{x}) = x - y + z$ and $t_4(\mathbf{x}) = 1 - x - y + z$. It can be checked that the function $u(\mathbf{x}) = e^{\cos(t_\alpha(\mathbf{x}))}$, $\mathbf{x} \in \Omega_\alpha$, $\alpha = 1, \cdots, 4$, belongs to $V$ and is such that $\mathbf{q}_m(\mathbf{x}) = -\Lambda_m \nabla u(\mathbf{x})$, $\mathbf{q}_f(\mathbf{x}) = -d_f \Lambda_f \nabla_\tau \gamma u(\mathbf{x})$ belongs to $W(\Omega, \Gamma)$. It will be used as exact solution of (1) with ad-hoc right hand sides and Dirichlet boundary conditions on $\partial \Omega$. For the numerical solutions, three different families of meshes are considered: uniform Cartesian meshes, a random perturbation of the previous Cartesian meshes, and

**Fig. 2** For the 3 families of meshes (*top* Cartesian meshes, *middle* randomly perturbated Cartesian meshes, and *bottom* tetrahedral meshes), and for the isotropic (*left*) and anisotropic (*right*) test cases: sum of the relative $L_2$ norm of the error in the matrix and in the fracture for the function and its gradients reconstructions and for both the VAG and HFV schemes function of the number of d.o.f. after elimination of the cell and Dirichlet unknowns

tetrahedral meshes generated by TetGen. To assess the error estimates of Proposition 3, we have computed the sum of the relative $L_2$ norms of the errors in the matrix and in the fractures, both for the function and for the gradient reconstructions. As exhibited in Fig. 2, the expected first orders of convergence are obtained both for the function reconstructions and the gradient reconstructions with observed superconvergence of order 2 for Cartesian meshes. We note that the HFV scheme seems to be less robust than the VAG scheme with respect to anisotropy. Also, as expected on tetrahedral meshes, the CPU time of the computation of the HFV solution is much larger of a factor around 10 than the CPU time obtained with the VAG scheme using for both schemes a GMRES solver preconditioned by ILUT.

# References

1. Alboin, C., Jaffré, J., Roberts, J., Serres, C.: Modeling fractures as interfaces for flow and transport in porous media. Contemp. Math. **295**, 13–24 (2002)
2. Brenner, K., Groza, M., Guichard, C., Masson, R.: Vertex approximate gradient scheme for hybrid dimensional two-phase darcy flows in fractured porous media. In: Proceedings of FVCA 7 (2014)
3. Eymard, R., Guichard, C., Herbin, R.: Small-stencil 3d schemes for diffusive flows in porous media. ESAIM: M2AN **46**, 265–290 (2010)
4. Eymard, R., Herbin, R., Gallouet, T.: Discretisation of heterogeneous and anisotropic diffusion problems on general nonconforming meshes. SUSHI: a scheme using stabilisation and hybrid interfaces. IMA J. Numer. Anal. **30**, 1009–1043 (2010). doi:10.1093/imanum/drn084
5. Reichenberger, V., Jakobs, H., Bastian, P., Helmig, R.: A mixed-dimensional finite volume method for multiphase flow in fractured porous media. Adv. Water Resources **29**, 1020–1036 (2006)
6. Tunc, X., Faille, I., Gallouet, T., Cacas, M.C., Havé, P.: A model for conductive faults with non matching grids. Comput. Geosci. **16**, 277–296 (2012)

# A Gradient Scheme for the Discretization of Richards Equation

Konstantin Brenner, Danielle Hilhorst and Huy Cuong Vu Do

**Abstract** We propose a finite volume method on general meshes for the discretization of Richards equation, an elliptic—parabolic equation modeling groundwater flow. The diffusion term, which can be anisotropic and heterogeneous, is discretized in a gradient scheme framework, which can be applied to a wide range of unstructured possibly non-matching polyhedral meshes in arbitrary space dimension. More precisely, we implement the SUSHI scheme which is also locally conservative. As is needed for Richards equation, the time discretization is fully implicit. We obtain a convergence result based upon energy-type estimates and the application of the Fréchet-Kolmogorov compactness theorem. We implement the scheme and present the results of a number of numerical tests.

## 1 Richards Equation

In this article, we study Richards equation using Kirchhoff transformation. Let $\Omega$ be an open bounded polygonal subset of $\mathbb{R}^d$ ($d = 1, 2$ or $3$) and let $T$ be a positive real number; Richards equation in the space-time domain $Q_T = \Omega \times (0, T)$ is given by

$$\partial_t \Big( \phi(\mathbf{x})\theta(p) \Big) - \mathrm{div}\Big( k_r(\theta(p))\mathbf{K}(\mathbf{x})\nabla(p + z) \Big) = 0, \tag{1}$$

K. Brenner
LJAD, University Nice Sophia-Antipolis, Nice, France

K. Brenner
Coffee Team Inria Sophia-Antipolis-Méditerranée, Valbonne, France
e-mail: konstantin.brenner@unice.fr

D. Hilhorst (✉)
Laboratoire de Mathématiques, CNRS et Université de Paris-Sud, Orsay, France
e-mail: Danielle.Hilhorst@math.u-psud.fr

H. C. Vu Do
Laboratoire de Mathématiques, Université de Paris-Sud, Orsay, France
e-mail: vdhuycuong@math.u-psud.fr

where $p(\mathbf{x}, t)$ is pressure head. The function $\theta(p)$ is the water saturation, $\phi(\mathbf{x})$ is the porosity, $\mathbf{K}(\mathbf{x})$ is the absolute permeability tensor and the scalar function $k_r(\theta)$ corresponds to the relative permeability, which depends on the water content. The space coordinates are defined by $\mathbf{x} = (x, z)$ in the case of space dimension 2 and $\mathbf{x} = (x, y, z)$ in the case of space dimension 3. Next we perform Kirchhoff's transformation. We set

$$F(s) := \int_0^s k_r(\theta(\tau))d\tau,$$

and suppose that the function $F$ is invertible. Then we set $u = F(p)$ in $Q_T$ and $c(u) = c(F(p)) = \theta(p)$. We remark that Kirchhoff's transformation leads to $\nabla u = k_r(\theta(p))\nabla p$. Thus, the Eq. (1) becomes

$$\partial_t\Big(\phi(\mathbf{x})c(u)\Big) - \mathrm{div}\Big(\mathbf{K}(\mathbf{x})\nabla u\Big) - \mathrm{div}\Big(k_r(c(u))\mathbf{K}(\mathbf{x})\nabla z\Big) = 0. \qquad (2)$$

Next, we consider the Eq. (2) together with the inhomogeneous Dirichlet boundary and the initial conditions

$$
\begin{aligned}
u(\mathbf{x}, t) &= \hat{u}(\mathbf{x}) \quad \text{a.e. on } \partial\Omega \times (0, T), \\
u(\mathbf{x}, 0) &= u_0(\mathbf{x}) \quad \text{a.e. in } \Omega.
\end{aligned}
\qquad (3)
$$

We make the following hypotheses:

$(H_1)$ $c$ is a continuous nondecreasing function such that there exist $\overline{\xi} > 0$ and $\underline{\xi} \geq 0$ satisfying $|c(u)| \leq \overline{\xi}(1 + |u|)$ for all $u \in \mathbb{R}$ and $|c(u) - c(v)| \geq \underline{\xi}|u - v|$ for all $u, v \in \mathbb{R}$.

$(H_2)$ $k_r$ is a continuous function such that $0 \leq k_r \leq \overline{k_r}$.

$(H_3)$ $\mathbf{K}$ is a bounded function from $\Omega$ to $\mathbb{M}_d(\mathbb{R})$, where $\mathbb{M}_d(\mathbb{R})$ denotes the set of real $d \times d$ matrices. Moreover for a.e. $\mathbf{x}$ in $\Omega$, $\mathbf{K}(\mathbf{x})$ is a symmetric positive definite matrix and there exist two positive constants $\overline{K}$ and $\underline{K}$ such that the eigenvalues of $\mathbf{K}(\mathbf{x})$ are included in $[\overline{K}, \underline{K}]$.

$(H_4)$ $u_0 \in L^2(\Omega)$, $\hat{u} \in H^1(\Omega)$ and $\phi \in L^\infty(\Omega)$ is such that $0 < \underline{\phi} \leq \phi(\mathbf{x}) \leq \overline{\phi}$ for a.e. $\mathbf{x} \in \Omega$.

**Definition** A function $u(\mathbf{x}, t)$ is said to be a weak solution of Problem (2)–(3) if:

(i)  $u(\mathbf{x}, t) - \hat{u}(\mathbf{x}) \in L^2(0, T; H_0^1(\Omega))$,

(ii)  $c(u) \in L^\infty(0, T; L^2(\Omega))$,

$$
\begin{aligned}
(iii) \quad &-\int_0^T \int_\Omega \phi(\mathbf{x})c(u(\mathbf{x}, t))\partial_t\varphi(\mathbf{x}, t)\, d\mathbf{x}dt - \int_\Omega \phi(\mathbf{x})c(u_0(\mathbf{x}))\varphi(\mathbf{x}, 0)\, d\mathbf{x} \\
&+ \int_0^T \int_\Omega \mathbf{K}(\mathbf{x})\nabla u(\mathbf{x}, t) \cdot \nabla\varphi(\mathbf{x}, t)\, d\mathbf{x}dt \qquad\qquad (4) \\
&+ \int_0^T \int_\Omega k_r(c(u(\mathbf{x}, t))\mathbf{K}(\mathbf{x})\nabla z \cdot \nabla\varphi(\mathbf{x}, t)\, d\mathbf{x}dt = 0,
\end{aligned}
$$

for all $\varphi \in L^2(0, T; H_0^1(\Omega))$ with $\varphi(\cdot, T) = 0$ and $\partial_t \varphi \in L^\infty(Q_T)$.

The discretization of Richards equation by means of gradient schemes has already been proposed by Eymard, Guichard, Herbin and Masson [3], where they consider Richards equation as a special case of two phase flow; however, they make the extra hypothesis that the relative permeability $k_r$ is bounded away from zero.

## 2 Gradient Discretization

Following [2] we define a gradient discretization $D$ of Problem (2)–(3) on a vector space $X_D$, or more precisely its subspace $X_D^0$ associated with the homogeneous Dirichlet boundary condition, and the two following linear operators:

- A gradient operator on the matrix domain: $\nabla_D : X_D \to L^2(\Omega)^d$.
- A function reconstruction operator on the matrix domain: $\pi_D : X_D \to L^2(\Omega)$.

**Coercivity**: We assume that $\|\nabla_D \cdot \|_{L^2(\Omega)^d}$ defines a norm on $X_D^0$. A gradient discretization $D$ is said to be coercive if there exists $C_D \geq 0$ such that for all $v \in X_D^0$ one has

$$\|\pi_D v\|_{L^2(\Omega)} \leq C_D \|\nabla_D v\|_{L^2(\Omega)^d}.$$

**Consistency**: Let $u \in H_0^1(\Omega)$, and let us define

$$S_D(u) = \inf_{v \in X_D^0} \left( \|\nabla_D v - \nabla u\|_{L^2(\Omega)^d} + \|\pi_D v - u\|_{L^2(\Omega)} \right).$$

Then, a sequence of gradient discretizations $(D^{(m)})_{m \in \mathbb{N}}$ is said to be consistent if for all $u \in H_0^1(\Omega)$, $\lim_{m \to +\infty} S_{D^{(m)}}(u) = 0$.

**Limit Conformity**: For all $\mathbf{q} \in H_{div}(\Omega)$, we define

$$W_D(\mathbf{q}) = \sup_{0 \neq v \in X_D^0} \frac{1}{\|\nabla_D v\|_{L^2(\Omega)^d}} \int_\Omega \nabla_D v \cdot \mathbf{q} + \pi_D v \mathrm{div}(\mathbf{q}) \, d\mathbf{x}. \tag{5}$$

Then, a sequence of gradient discretizations $(D^{(m)})_{m \in \mathbb{N}}$ is said to be limit conforming if for all $\mathbf{q} \in H_{div}(\Omega)$, $\lim_{m \to +\infty} W_{D^{(m)}}(\mathbf{q}) = 0$.

**Compactness**: A sequence of gradient discretizations $(D^{(m)})_{m \in \mathbb{N}}$ is said to be compact if for all sequences $v_m \in X_{D^{(m)}}^0$, $m \in \mathbb{N}$ such that there exists $C > 0$ with $\|\nabla_{D^{(m)}} v_m\|_{L^2(\Omega)^d} \leq C$ for all $m \in \mathbb{N}$, then there exist $\bar{v} \in L^2(\Omega)$ such that

$$\lim_{m \to +\infty} \|\pi_{D^{(m)}} v_m - \bar{v}\|_{L^2(\Omega)} = 0.$$

For $N \in \mathbb{N}^*$, let us consider the time discretization $t^0 = 0 < t^1 < \cdots < t^{n-1} < t^n \cdots < t^N = T$ of the time interval $[0, T]$. We denote the time steps by $\delta t^n = t^n - t^{n-1}$ for all $n \in \{1, \cdots, N\}$ while $\delta t$ stands for the whole sequence $(\delta t^n)_{n \in \{1,...,N\}}$. For all $v = \left(v^n \in X_D\right)_{n=1,\cdots,N}$ we set $\pi_{D,\delta t} v(\mathbf{x}, t) = \pi_D v^n(\mathbf{x})$ and $\nabla_{D,\delta t} v(\mathbf{x}, t) = \nabla_D v^n(\mathbf{x})$ for all $(\mathbf{x}, t) \in \Omega \times (t^{n-1}, t^n], n \in \{1, \ldots, N\}$.

**Discrete variational formulation**: For a given $u^0, \hat{u}_D \in X_D$ find $u = \left(u^n \in X_D\right)_{n \in \{1,...,N\}}$ such that for each $n \in \{1, \ldots, N\}$, $u^n - \hat{u}_D \in X_D^0$ and for all $v \in X_D^0$

$$\int_\Omega \phi \frac{c(\pi_D u^n) - c(\pi_D u^{n-1})}{\delta t^n} \pi_D v \, d\mathbf{x} + \int_\Omega \mathbf{K}(\nabla_D u^n + k_r(\pi_D u^n)\nabla z) \cdot \nabla_D v \, d\mathbf{x} = 0 \tag{6}$$

**Proposition 1** *There exists at least one solution of* (6)*; moreover there exists a positive* $C$ *only depending on* $\underline{\phi}, \overline{\phi}, \underline{\xi}, \overline{\xi}, \underline{K}, \overline{K}, \overline{k_r}, \Omega, T, u_0, \hat{u}$ *as well as on* $\|c(\pi_D u^0) - c(u_0)\|_{L^2(\Omega)}, \|\pi_D \hat{u}_D - \hat{u}\|_{L^2(\Omega)}$ *and* $\|\nabla_D \hat{u}_D - \nabla \hat{u}\|_{L^2(\Omega)}$ *such that*

$$\|c(\pi_{D,\delta t} u)\|_{L^\infty(0,T;L^2(\Omega))} + \|\nabla_{D,\delta t} u\|_{L^2(Q_T)^d} \leq C \tag{7}$$

*for any solution* $u$ *of* (6).

*Proof* In order to keep this presentation short, we only prove below the priori estimate (7), and only in the case of homogeneous Dirichlet boundary conditions; the adaptation to the inhomogeneous case is straightforward, and the existence of a discrete solution can be deduced using a standard argument based upon the topological degree. Let $u = (u^n)_{n \in \{1,...,N\}}$ be a solution of (6) and define

$$A_{D,\delta t}^n(v) = \int_\Omega \phi \frac{c(\pi_D u^n) - c(\pi_D u^{n-1})}{\delta t^n} \pi_D v \, d\mathbf{x},$$

$$B_{D,\delta t}^n(v) = \int_\Omega \mathbf{K} \nabla_D u^n \cdot \nabla_D v \, d\mathbf{x}, \quad C_{D,\delta t}^n(v) = \int_\Omega \mathbf{K} k_r(\pi_D u^n)\nabla z \cdot \nabla_D v \, d\mathbf{x}, \tag{8}$$

for all $n \in \{1, \ldots, N\}$ and $v \in X_D^0$. The terms defined above satisfy

$$A_{D,\delta t}^n(v) + B_{D,\delta t}^n(v) + C_{D,\delta t}^n(v) = 0 \text{ for all } v \in X_D^0. \tag{9}$$

Let us first estimate $\sum_{n=1}^m \delta t^n A_{D,\delta t}^n(u^n)$ for $m \in \{1, \ldots, N\}$; we define

$$\xi(u) = c(u)u - \int_0^u c(\tau) \, d\tau \quad \text{for all } u \in \mathbb{R}.$$

For all $a, b \in \mathbb{R}$, one has $\xi(a) - \xi(b) = (c(a) - c(b))a - \int_b^a (c(\tau) - c(b)) \, d\tau$ and since $c$ is nondecreasing we have that $\xi(a) - \xi(b) \leq (c(a) - c(b))a$. It implies that

$$\sum_{n=1}^{m} \delta t^n A_{D,\delta t}^n (u^n) \geq \int_{\Omega} \phi(\xi(\pi_D u^m) - \xi(\pi_D u^0)) \, d\mathbf{x}. \tag{10}$$

For all $a \in \mathbb{R}$ it holds $\frac{1}{2}\underline{\xi}a^2 \leq \xi(u) \leq c(a)a \leq \frac{(c(a))^2}{\underline{\xi}}$, therefore

$$\sum_{n=1}^{m} \Delta t^n A_{D,\Delta t}^n (u^n) \geq \frac{1}{2}\underline{\xi}\underline{\phi}\|\pi_D u^m\|_{L^2(\Omega)}^2 - \frac{1}{\underline{\xi}\overline{\phi}}\|c(\pi_D u^0)\|_{L^2(\Omega)}^2. \tag{11}$$

Using the assumptions $(H_2)$–$(H_3)$ we deduce that $B_{D,\delta t}^n (u^n) \geq \underline{K}\|\nabla_D u^n\|_{L^2(\Omega)^d}^2$ and that $C_{D,\delta t}^n (u^n) \leq \overline{k_r}\,\overline{K}|\Omega|^{1/2}\|\nabla_D u^n\|_{L^2(\Omega)^d}$ for all $n \in \{1, \ldots, N\}$. Combining these inequalities with (9) and (11) gives

$$\frac{1}{2}\underline{\xi}\,\underline{\phi}\|\pi_D u^m\|_{L^2(\Omega)}^2 + \underline{K}\sum_{n=1}^{m} \delta t^n \|\nabla_D u^n\|_{L^2(\Omega)^d}^2$$
$$\leq \frac{1}{\underline{\xi}\overline{\phi}}\|c(\pi_D u^0)\|_{L^2(\Omega)}^2 + \overline{k_r}\,\overline{K}|\Omega|^{1/2}\sum_{n=1}^{m} \delta t^n \|\nabla_D u^n\|_{L^2(\Omega)^d}.$$

Applying Young's inequality to the last term above, we obtain

$$\overline{k_r}\,\overline{K}|\Omega|^{1/2}\sum_{n=1}^{m} \delta t^n \|\nabla_D u^n\|_{L^2(\Omega)^d} \leq \frac{1}{2\varepsilon}\overline{k_r}^2\overline{K}T|\Omega| + \frac{\varepsilon}{2}\overline{K}\sum_{n=1}^{m} \delta t^n \|\nabla_D u^n\|_{L^2(\Omega)^d}^2.$$

This leads to

$$\frac{1}{2}\underline{\xi}\,\underline{\phi}\|(\pi_{D,\delta t}u)\|_{L^\infty(0,T;L^2(\Omega))}^2 + (\underline{K} - \frac{\varepsilon}{2}\overline{K}|)\|\nabla_{D,\delta t}u\|_{L^2(Q_T)^d}^2$$
$$\leq \frac{1}{\underline{\xi}\overline{\phi}}\|c(\pi_D u^0)\|_{L^2(\Omega)}^2 + \frac{1}{2\varepsilon}\overline{k_r}^2\overline{K}T|\Omega|. \tag{12}$$

One completes the proof of the estimate (7) by choosing $\varepsilon = \underline{K}/\overline{K}$ and using the assumptions $(H_1)$ and $(H_4)$.

The following result is rather standard and given without proof.

**Proposition 2** *Let $u$ be a solution to* (6). *There exists a positive constant $C$ only depending on $\underline{\phi}, \overline{\phi}, \underline{\xi}, \overline{\xi}, \underline{K}, \overline{K}, \overline{k_r}, \Omega, T, u_0, \hat{u}$ as well as on $\|c(\pi_D u^0) - c(u_0)\|_{L^2(\Omega)}$, $\|\pi_D \hat{u}_D - \hat{u}\|_{L^2(\Omega)}$ and $\|\nabla_D \hat{u}_D - \nabla \hat{u}\|_{L^2(\Omega)}$ such that for all $\tau \in (0, T)$, there holds*

$$\int_0^{T-\tau} \int_\Omega \left(\pi_{D,\delta t}u(\mathbf{x}, t+\tau) - \pi_{D,\delta t}u(\mathbf{x}, t)\right)^2 d\mathbf{x}dt \leq C\tau.$$

**Theorem 1** *Let $(D^{(m)}, \delta t^{(m)})_{m \in \mathbb{N}}$ be a family of discretizations, where $(D^{(m)})_{m \in \mathbb{N}}$ is assumed to be limit conforming, consistent, compact and uniformly coercive in the sense that there exists $C_1$ such that $C_{D^{(m)}} \leq C_1$ for all $m \in \mathbb{N}$; moreover we assume that $\|c(\pi_{D^{(m)}} u_m^0) - c(u_0)\|_{L^2(\Omega)}$, $\|\pi_{D^{(m)}} \hat{u}_{D^{(m)}} - \hat{u}\|_{L^2(\Omega)}$ and $\|\nabla_{D^{(m)}} \hat{u}_{D^{(m)}} - \nabla \hat{u}\|_{L^2(\Omega)}$, $\max_n \delta t^{(m),n}$ tend to 0 as $m \to \infty$. Let $u_m$ be a solution of (6) for all $m \in \mathbb{N}$. Then, up to a subsequence*

$$\pi_{D^{(m)}, \delta t^{(m)}} u_m \to \overline{u} \text{ in } L^2(Q_T),$$
$$\nabla_{D^{(m)}, \delta t^{(m)}} u_m \rightharpoonup \nabla \overline{u} \text{ in } L^2(Q_T)^d,$$

*where $\overline{u} \in L^2(0, T; H^1(\Omega))$ is a solution of (4).*

*Proof* Using the compactness and the uniform coercivity of the sequence $D^{(m)}$ as well as Propositions 1 and 2, we deduce from Fréchet-Kolmogorov theorem that the sequence $\{\pi_{D^{(m)}, \delta t^{(m)}} u_m - \pi_{D^{(m)}} \hat{u}_{D^{(m)}}\}$ is relatively compact in $L^2(Q_T)$. Therefore, we may extract a subsequence of $\{u_m\}$ (denoted again by $\{u_m\}$) such that $\pi_{D^{(m)}, \delta t^{(m)}} u_m$ converges to some $\overline{u} \in L^2(Q_T)$ strongly in $L^2(Q_T)$ and $\nabla_{D^{(m)}, \delta t^{(m)}} u_m$ is weakly convergent in $L^2(Q_T)$. It follows from Lemma 7.1 of [1] that the subsequence $u_m$ can also be chosen in such way that $c(\pi_{D^{(m)}, \delta t^{(m)}} u_m)$ and $k_r(c(\pi_{D^{(m)}, \delta t^{(m)}} u_m))$ converge strongly in $L^2(Q_T)$ to $c(\overline{u})$ and $k_r(c(\overline{u}))$ respectively; moreover one deduces from (7) that $c(\overline{u}) \in L^\infty(0, T; L^2(\Omega))$. Finally we deduce from the limit conformity of the scheme that $\overline{u} - \hat{u} \in L^2(0, T; H_0^1(\Omega))$ and that $\nabla_{D^{(m)}, \delta t^{(m)}} u_m \rightharpoonup \nabla \overline{u}$ in $L^2(Q_T)^d$ as $m \to +\infty$. Using again the limit conformity and consistency of the scheme we deduce that $\overline{u}$ is a weak solution of (4).

## 3 Numerical Tests

### 3.1 The Hornung-Messing Problem

The Hornung-Messing problem is a standard test (cf. for instance [5]). We consider a horizontal flow in a homogeneous ground $\Omega = [0, 1]^2$ and set $T = 1$. The problem after Kirchhoff's transformation is given by Problem (2) with

$$c(u) = \theta(p) = \begin{cases} \pi^2/2 - 2\arctan^2(\dfrac{u}{2-u}) & \text{if } p < 0, \\ \pi^2/2 & \text{otherwise}, \end{cases}$$

and suitable boundary and initial conditions. Let $s = x - z - t$, its solution is given:

$$u(x, z, t) = \begin{cases} \dfrac{2p(x, z, t)}{1 + p(x, z, t)} & \text{if } p < 0, \\ 2p(x, z, t) & \text{otherwise}, \end{cases} \quad p(x, z, t) = \begin{cases} -s/2 & \text{if } s < 0, \\ -\tan\left(\dfrac{e^s - 1}{e^s + 1}\right) & \text{otherwise}. \end{cases}$$
$$(13)$$

**Fig. 1** Saturation at $t = 0.1$ s and at $t = 0.4$ s. The medium is unsaturated on the right-hand side of the space domain where $\theta < 4.9348$ and fully saturated elsewhere

In this test, we apply the Sushi scheme [4] using an adaptive mesh driven by the variations of the saturation. We prescribe the Neumann boundary condition deduced from (13) on the line $x = 0$ and an inhomogeneous Dirichlet boundary condition elsewhere. We use an initially square mesh, which is such that each square can be decomposed again into four smaller square elements. Whereas the standard finite volume scheme is not suited to handle such a non-conforming adaptive mesh, the SUSHI scheme is compatible with these non-conforming volume elements (Fig. 1).

We introduce the relative error in $L^2(Q_T)$ between the exact and the numerical solution as well as the experimental order of convergence

$$err(u) = \frac{\|(u_{exact}(\mathbf{x}, t_n) - u_{D,\delta t}(\mathbf{x}, t_n))\|_{L^2(Q_T)}}{\|(u_{exact}(\mathbf{x}, t_n))\|_{L^2(Q_T)}}, \quad eoc = \frac{log(err(u_i)/err(u_{i+1}))}{log(h_{D_i}/h_{D_{i+1}})},$$

where $u_i$ is the solution corresponding to the space discretization $D_i$. Table 1 shows the error using a uniform square mesh with various mesh sizes and time steps in the four first lines. Note that the scheme is only first order accurate with respect to time; therefore in order to obtain second order convergence we choose $\delta t$ proportional to $h_D^2$. We also compare the error for the approximate saturation using a uniform mesh and an adaptive mesh with a similar number of unknowns. In both cases: about 300 unknowns (line 2–line 5) and 1,200 unknowns (line 3–line 6), the adaptive mesh compared to the fixed one provides slightly better results for the saturation $c(u)$. The observed computational gain is rather small (about 10–20%), which is due to the fact that the area of high gradients of $c$ is comparatively large.

## 3.2 The Haverkamp Problem

We consider the case of a sand ground represented by the space domain $\Omega = (0, 2) \times (0, 40)$ on the time interval $[0, 600]$. The parameters are given by [7]

**Table 1** Number of time steps $N$, mesh diameter $h_D$, number of unknown $N_{unk}$, the error on the solution $err(u)$, and on the saturation $err(c(u))$ and the experimental order of convergence $eoc$

| Mesh | $N$ | $h_D$ | $N_{unk}$ | $err(u)$ | $err(c(u))$ | $eoc(u)$ |
|------|-----|-------|-----------|----------|-------------|----------|
| Uniform | 25 | 0.2 | 85 | $2.40 \cdot 10^{-2}$ | $1.60 \cdot 10^{-5}$ | – |
| Uniform | 100 | 0.1 | 320 | $6.09 \cdot 10^{-3}$ | $4.13 \cdot 10^{-6}$ | 1.98 |
| Uniform | 400 | 0.05 | 1240 | $1.53 \cdot 10^{-3}$ | $2.90 \cdot 10^{-6}$ | 2.00 |
| Uniform | 1600 | 0.025 | 4880 | $3.76 \cdot 10^{-3}$ | $1.83 \cdot 10^{-6}$ | 2.02 |
| Adaptive | 200 | 0.143 | 302 | $5.62 \cdot 10^{-3}$ | $3.67 \cdot 10^{-6}$ | – |
| Adaptive | 800 | 0.071 | 1232 | $1.32 \cdot 10^{-3}$ | $2.19 \cdot 10^{-6}$ | – |



**Fig. 2** Time evolution of the pressure p and the adaptive mesh

$$\theta(p) = \begin{cases} \dfrac{\theta_s - \theta_r)}{1 + |\alpha p|^\beta} + \theta_r, & \text{if } p < 0, \\ \theta_s, & \text{otherwise,} \end{cases} \quad k_r(\theta(p)) = \begin{cases} \dfrac{K_s}{1 + |Ap|^\gamma}, & \text{if } p < 0, \\ K_s, & \text{otherwise,} \end{cases}$$

where $\theta_s = 0.287$, $\theta_r = 0.075$, $\alpha = 0.0271$ $\beta = 3.96$, $K_s = 9.44e - 3$, $A = 0.0524$ and $\gamma = 4.74$. From $\theta$ and $K$, we have tabulated suitable values for the functions $c$ and $K_c$. We have taken here the initial condition $p = -61.5$, a homogeneous Neumann boundary condition for $x = 0$ and $x = 1$, the Dirichlet boundary condition $p = -61.5$ for $z = 0$ and $p = -20.7$ for $z = 40$.

We use an adaptive mesh and the time step $\delta t = 1$ to perform the test. Figure 2-*left* represents the pressure profile at various times. In this test, no analytical solution is known. Therefore we compare our numerical solution with that of Pierre Sochala [8, Fig. 2.6, p. 35] which is obtained by means of a finite element method. Our results are quite similar to his. Figure 2-*right* shows the time evolution of the mesh at different times corresponding to the pressure profiles in Fig. 2-*left*.

# References

1. Angelini, O., Brenner, K., Hilhorst, D.: A finite volume method on general meshes for a degenerate parabolic convection-reaction-diffusion equation. Numer. Math. **123**, 219–257 (2013). doi:10.1007/s00211-012-0485-5

2. Droniou, J., Eymard, R., Gallouet, T., Herbin, R.: Gradient schemes: a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. Math. Models Meth. Appl. Sci. **23**(13), 2395–2432 (2013)

3. Eymard, R., Guichard, C., Herbin, R., Masson, R.: Gradient schemes for two-phase flow in heterogeneous porous media and Richards equation. ZAMM (2013). doi:10.1002/zamm.201200206

4. Eymard, R., Gallouët, T., Herbin, R.: Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes. Sushi: a scheme using stabilization and hybrid interfaces. IMA J. Numer. Anal. **30**(4), 1009–1043 (2010)

5. Eymard, R., Gutnic, M., Hilhorst, D.: The finite volume method for Richards equation. Comput. Geosci. **3**(3–4), 259–294 (2000)

6. Eymard, R., Hilhorst, D., Vohralik, M.: A combined finite volume scheme nonconforming/mixed-hybrid finite element scheme for degenerate parabolic problems. Numer. Math. **105**, 73–131 (2006)

7. Haverkamp, R., Vauclin, M., Touma, J., Wierenga, P., Vachaud, G.: A comparison of numerical simulation models for one-dimensional infiltration. Soil Sci. Soc. Am. J. **41**(2),285–294 (1977)

8. Sochala, P.: Méthodes numériques pour les écoulements souterrains et couplage avec le ruissellement. Thèse de doctorat, Ecole Nationale des Ponts et Chaussées (2008)

# Convergence of a Finite Volume Scheme for a Corrosion Model

**Claire Chainais-Hillairet, Pierre-Louis Colin and Ingrid Lacroix-Violet**

**Abstract** We consider a drift-diffusion system describing the corrosion of an iron based alloy in nuclear waste repository. In particular, we are interested in the convergence of a numerical scheme consisting in an implicit Euler scheme in time and a Scharfetter-Gummel finite volume scheme in space.

## 1 General Framework

The DPCM model, introduced by Bataillon et al. in [1] , is related to the corrosion of an iron based alloy in a nuclear waste repository. It describes the evolution of a dense oxide layer formed at the surface of the metal when it is in contact with claystone.

The system is made of drift-diffusion equations on the charge densities coupled with a Poisson equation on the electric potential. The boundary conditions induced by the electrochemical reactions at the interfaces are Robin boundary conditions. Moreover, the system includes moving boundary equations.

A numerical scheme for the DPCM model has been proposed and verified by numerical experiments in [2]. The proof of convergence of the scheme proposed in [2] for the full system is challenging. In this paper, we will focus on a simplified model with only two species and on a fixed domain. It permits to show how to deal with the boundary conditions in the proof of convergence.

C. Chainais-Hillairet · P.-L. Colin (✉) · I. Lacroix-Violet
Laboratoire Paul Painlevé, UMR CNRS 8524, Université Lille 1,
59655 Villeneuve d'Ascq Cedex, France
e-mail: pierre-louis.colin@math.univ-lille1.fr

C. Chainais-Hillairet
e-mail: Claire.Chainais@math.univ-lille1.fr

I. Lacroix-Violet
e-mail: ingrid.violet@math.univ-lille1.fr

## 2 Presentation of the Model and of the Hypotheses

The unknowns are the densities of electrons $N$, cations $P$ and the electric potential $\Psi$. The current densities of electrons and cations are respectively denoted $J_N$ and $J_P$; they contain both a drift part and a diffusion part. For $u = N$ or $P$, we also denote by $z_u$ the charge of the species ($z_P = 3$, $z_N = -1$) and $\varepsilon_u$ the ratio of diffusion coefficients arising in the scaling. The dimensionless system writes:

$$\varepsilon_u \partial_t u + \partial_x J_u = 0, \qquad J_u = -\partial_x u - z_u u \partial_x \Psi, \quad \text{in } (0, 1), \quad \text{for } u = P, N, \quad \text{(1a)}$$

$$- \lambda^2 \partial_{xx}^2 \Psi = 3P - N + \rho_{hl}, \quad \text{in } (0, 1), \tag{1b}$$

where $\lambda^2$ is a dimensionless parameter ($\lambda$ is the rescaled Debye length) and $\rho_{hl}$ is the net charge density of the ionic species in the host lattice which is supposed constant in the whole layer, with $\rho_{hl} = -5$. Charge carriers are created and consumed at both interfaces. The boundary conditions are prescribed by the kinetics of the electrochemical reactions at the interfaces. It leads to Robin boundary conditions, which are assumed to have the same form for electrons and cations:

$$- J_u(0) = \beta_u^0 (\Psi(0)) u(0) - \gamma_u^0 (\Psi(0)), \quad \text{on } x = 0, \quad \text{for } u = P, N, \tag{2a}$$

$$J_u(1) = \beta_u^1 (V - \Psi(1)) u(1) - \gamma_u^1 (V - \Psi(1)), \quad \text{on } x = 1, \quad \text{for } u = P, N, \tag{2b}$$

where $V$ is the applied voltage and $(\beta_u^i, \gamma_u^i)_{i=0,1}$ are continuous, nonnegative functions defined by:

$$\beta_u^i(x) = m_u^i e^{-z_u b_u^i x} + k_u^i e^{z_u a_u^i x}, \quad \text{for } u = P, N, \tag{3a}$$

$$\gamma_u^i(x) = m_u^i u^m e^{-z_u b_u^i x}, \quad \text{for } u = P, N. \tag{3b}$$

For the electric potential, the boundary conditions have the following form

$$\Psi(0) - \alpha_0 \partial_x \Psi(0) = \Delta \Psi_0^{pzc}, \quad \text{on } x = 0, \tag{4a}$$

$$\Psi(1) + \alpha_1 \partial_x \Psi(1) = V - \Delta \Psi_1^{pzc}, \quad \text{on } x = 1, \tag{4b}$$

where $\alpha_0$ and $\alpha_1$ are nonnegative dimensionless parameters and $(\Delta \Psi_i^{pzc})_{i=0,1}$ are the inner and outer voltages of zero charge. The system is supplemented with initial conditions, given in $L^\infty(0, 1)$

$$u(x, 0) = u^0(x), \quad \text{for } u = P, N. \tag{5}$$

Let us give a sense to the parameters in Eqs. (1a)–(4b) and some hypotheses:

- $(m_P^i, m_N^i, k_P^i, k_N^i)_{i=0,1}$ are the interface kinetic coefficients which are supposed constant and nonnegative.
- $P^m$ is the maximum occupancy for cations and $N^m$ is the electron densities of state in the metal. We assume that

$$z_P P^m + z_N N^m + \rho_{hl} = 0. \tag{6}$$

- $(a_P^i, a_N^i, b_P^i, b_N^i)_{i=0,1}$ are the nonnegative transfer coefficients which satisfy $a_u^0 + b_u^0 = a_u^1 + b_u^1 = 1, \quad u = P, N.$

Let us also consider some compatibility hypotheses on the data (see [2])

$$-\frac{1}{3a_P^0}\left(1 + \log\left(\alpha_0 a_P^0 k_P^0\right)\right) \leqslant \Delta\Psi_0^{pzc} \leqslant \frac{1}{a_N^0}\left(1 + \log\left(\alpha_0 a_N^0 k_N^0\right)\right), \tag{7a}$$

$$-\frac{1}{b_N^1}\left(1 + \log\left(\alpha_1 b_N^1 m_P^1\right)\right) \leqslant \Delta\Psi_1^{pzc} \leqslant \frac{1}{3b_P^1}\left(1 + \log\left(\alpha_1 b_P^1 m_P^1\right)\right). \tag{7b}$$

## 3 Numerical Scheme

We are interested in the convergence analysis of the fully implicit scheme introduced in [2]. It is an Euler implicit in time and finite volume in space scheme with a Scharfetter-Gummel approximation of the convection-diffusion fluxes.

Let us consider a mesh $\mathcal{T}$ for the domain $[0, 1]$, i.e a family of given points $(x_i)_{0 \leq i \leq I+1}$ satisfying: $0 = x_0 < x_1 < x_2 < \cdots < x_I < x_{I+1} = 1$. We define $x_{i+\frac{1}{2}} = \dfrac{x_i + x_{i+1}}{2}$, for $1 \leq i \leq I - 1$ and we set $x_{\frac{1}{2}} = x_0 = 0$, $x_{I+\frac{1}{2}} = x_{I+1} = 1$. Let us consider the mesh cells $[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]$, discretization steps:

$$h_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}, \ \forall 1 \leq i \leq I \text{ and } h_{i+\frac{1}{2}} = x_{i+1} - x_i, \ \forall 0 \leq i \leq I.$$

We define $h = \max\{h_i, 1 \leq i \leq I\}$ the mesh size. Let us denote by $\Delta t$ the time step given by $N_T \Delta t = T$ and consider the sequence $(t^n)_{0 \leqslant n \leqslant N_T}$ such that $t^n = n\Delta t$. Then, the scheme writes:

$$-\lambda^2\left(\mathrm{d}\Psi_{i+\frac{1}{2}}^{n+1} - \mathrm{d}\Psi_{i-\frac{1}{2}}^{n+1}\right) = h_i\left(3P_i^{n+1} - N_i^{n+1} + \rho_{hl}\right), \qquad 1 \leqslant i \leqslant I, \tag{8a}$$

$$\varepsilon_u h_i \frac{u_i^{n+1} - u_i^n}{\Delta t} + \mathcal{F}_{u,i+\frac{1}{2}}^{n+1} - \mathcal{F}_{u,i-\frac{1}{2}}^{n+1} = 0, \qquad 1 \leqslant i \leqslant I, \tag{8b}$$

with the numerical fluxes defined for $0 \le i \le I$ by:

$$d\Psi_{i+\frac{1}{2}}^{n+1} = \frac{\Psi_{i+1}^{n+1} - \Psi_i^{n+1}}{h_{i+\frac{1}{2}}}, \tag{9a}$$

$$\mathscr{F}_{u,i+\frac{1}{2}}^{n+1} = \frac{B\left(z_u h_{i+\frac{1}{2}} d\Psi_{i+\frac{1}{2}}^{n+1}\right) u_i^{n+1} - B\left(-z_u h_{i+\frac{1}{2}} d\Psi_{i+\frac{1}{2}}^{n+1}\right) u_{i+1}^{n+1}}{h_{i+\frac{1}{2}}}, \tag{9b}$$

where $B$ is the Bernoulli function, leading to Scharfetter-Gummel fluxes [7], i.e.:

$$B(x) = \frac{x}{e^x - 1} \quad \forall x \ne 0, \quad B(0) = 1.$$

We supplement the scheme with the discretization of the boundary conditions

$$\Psi_0^{n+1} - \alpha_0 d\Psi_{\frac{1}{2}}^{n+1} = \Delta\Psi_0^{pzc}, \tag{10a}$$

$$\Psi_{I+1}^{n+1} + \alpha_1 d\Psi_{I+\frac{1}{2}}^{n+1} = V - \Delta\Psi_1^{pzc}, \tag{10b}$$

$$-\mathscr{F}_{u,\frac{1}{2}}^{n+1} = \beta_u^0\left(\Psi_0^{n+1}\right) u_0^{n+1} - \gamma_u^0\left(\Psi_0^{n+1}\right), \quad \text{for } u = P, \, N, \tag{10c}$$

$$\mathscr{F}_{u,I+\frac{1}{2}}^{n+1} = \beta_u^1\left(V - \Psi_{I+1}^{n+1}\right) u_{I+1}^{n+1} - \gamma_u^1\left(V - \Psi_{I+1}^{n+1}\right), \quad \text{for } u = P, \, N, \tag{10d}$$

and of the initial conditions

$$u_i^0 = \frac{1}{h_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} u^0(x)\,dx, \quad 1 \le i \le I, \quad \text{for } u = P, \, N. \tag{11}$$

## 4 Main Results

The goal is to prove the convergence of a sequence of solutions to the numerical scheme (8a)–(11) to a solution of (1a)–(5).

**Proposition 1** *Let $\varepsilon_N, \varepsilon_P \ge 0$. Under the hypotheses on the data given in Sect. 2, there exists a solution $(P_i^{n+1}, N_i^{n+1}, \Psi_i^{n+1})_{0 \le i \le I+1, n \ge 0}$ to the fully implicit scheme (8a)–(11) satisfying the following stability property:*

$$0 \le P_i^n \le P^m \text{ and } 0 \le N_i^n \le N^m, \quad \forall 0 \le i \le I+1, \, \forall n \ge 0. \tag{12}$$

*Outline of the Proof*
This proposition is proved for $\varepsilon_N, \varepsilon_P > 0$ in [2]. Moreover, we can easily expand the result if $\varepsilon_N = 0$ or/and $\varepsilon_P = 0$. It is based on a linearization of the scheme and on a fixed point theorem.

As it is classical for finite volume scheme, we consider an approximate solution which is piecewise constant in time and space. For a sequence of meshes and time steps $(\mathcal{T}_m, \Delta t_m)_m$ such that $size(\mathcal{T}_m) \to 0$ and $\Delta t_m \to 0$ as $m \to +\infty$, we define a sequence of approximate solutions $(P_m, N_m, \Psi_m)_m$.

**Theorem 1** *Let $\varepsilon_N, \varepsilon_P > 0$. Under the hypotheses on the data given Sect. 2, up to subsequence, we have as $m \to +\infty$*

$$u_m \to u \quad strongly\ in\ \mathrm{L}^2(0, T, \mathrm{L}^2(0, 1)), \quad for\ u = P,\ N$$
$$\Psi_m \to \Psi \quad strongly\ in\ \mathrm{L}^2(0, T, \mathrm{L}^2(0, 1)),$$

*with additional weak convergence on the discrete gradients. Moreover $(P, N, \Psi)$ is a weak solution to the model (1a)–(5).*

*Outline of the Proof*
Let us consider the sequences $(\Psi_m)_m$ and $(u_m)_m$, for $u = P,\ N$. The proof splits into three steps: a priori estimates in the appropriate space (see Sect. 5), compactness of the sequences of approximate solutions and passage to the limit in the numerical scheme. In this paper, we focus on the proof of the a priori estimates satisfied by $N$, $P$ and $\Psi$. The main difficulty lies in the treatment of the boundary conditions.

## 5 A Priori Estimates

We introduce a discrete $\mathrm{H}^1$ norm in order to have a synthetical form for the estimates.

**Definition 1** Let $w = (w_i) \in \mathbb{R}^{I+2}$, $\|w\|_{\mathrm{H}^1(0,1)}^2 = \sum_{i=0}^{I} \frac{(w_{i+1} - w_i)^2}{h_{i+\frac{1}{2}}} + w_0^2 + w_{I+1}^2$.

Let us remark that the following discrete Poincaré estimate holds: $\forall w = (w_i) \in \mathbb{R}^{I+2}$, $\sum_{i=1}^{I} h_i(w_i)^2 \leq 2\|w\|_{\mathrm{H}^1(0,1)}^2$.

**Proposition 2** *Under the hypotheses of Theorem 1, there exists a constant $C$ depending only on the data and independent of $\Delta t$ and $h$, such that:*

$$\|\Psi^{n+1}\|_{\mathrm{H}^1(0,1)}^2 \leqslant C, \quad \forall 0 \leqslant n \leqslant N_T - 1, \tag{13}$$

$$\sum_{n=0}^{N_T-1} \Delta t \|\Psi^{n+1}\|_{\mathrm{H}^1(0,1)}^2 \leqslant CT. \tag{14}$$

*Outline of the proof*
We multiply (8a) by $\Psi_i^{n+1}$ and we sum over $i$. Then, using the boundary conditions (10a)–(10d), Young inequality and the discrete Poincaré estimate, we easily obtain estimate (13). Estimate (14) is a direct consequence of (13).

**Proposition 3** *Under the hypotheses of Theorem 1, there exists a constant C depending only on the data and independent of $\Delta t$ and h, such that:*

$$\sum_{n=0}^{N_T-1} \Delta t \|u^{n+1}\|_{H^1(0,1)}^2 \leqslant C. \tag{15}$$

*Proof* Although the proof is based on a classical method [5], the particular boundary conditions imply a new difficulty. Moreover, the Scharfetter-Gummel fluxes are treated as in [3].

Let us multiply (8b) with $\Delta t u_i^{n+1}$ and sum over $i$ and $n$, then:

$$0 = \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \varepsilon_u h_i u_i^{n+1} \left(u_i^{n+1} - u_i^n\right) + \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \Delta t u_i^{n+1} \left(\mathscr{F}_{u,i+\frac{1}{2}}^{n+1} - \mathscr{F}_{u,i-\frac{1}{2}}^{n+1}\right). \tag{16}$$

We have

$$\sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \varepsilon_u h_i u_i^{n+1} \left(u_i^{n+1} - u_i^n\right) \geqslant \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \frac{\varepsilon_u h_i}{2} \left(u_i^{n+1} - u_i^n\right)^2 - \sum_{i=1}^{I} \frac{\varepsilon_u h_i}{2} \left(u_i^0\right)^2. \tag{17}$$

Using as in [3] the following decomposition of fluxes:

$$\mathscr{F}_{u,i+\frac{1}{2}}^{n+1} = -z_u d\Psi_{i+\frac{1}{2}}^{n+1} \frac{u_i^{n+1} + u_{i+1}^{n+1}}{2}$$
$$+ \frac{z_u d\Psi_{i+\frac{1}{2}}^{n+1}}{2} \coth\left(\frac{-z_u h_{i+\frac{1}{2}} d\Psi_{i+\frac{1}{2}}^{n+1}}{2}\right) \left(u_{i+1}^{n+1} - u_i^{n+1}\right),$$

we have

$$\sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \Delta t u_i^{n+1} \left(\mathscr{F}_{u,i+\frac{1}{2}}^{n+1} - \mathscr{F}_{u,i-\frac{1}{2}}^{n+1}\right) = B_1 + B_2 + B_3,$$

with

$$B_1 = \sum_{n=0}^{N_T-1} \sum_{i=0}^{I} \frac{\Delta t z_u}{2} d\Psi_{i+\frac{1}{2}}^{n+1} \left(\left(u_{i+1}^{n+1}\right)^2 - \left(u_i^{n+1}\right)^2\right),$$

$$B_2 = -\sum_{n=0}^{N_T-1} \sum_{i=0}^{I} \frac{\Delta t z_u}{2} d\Psi_{i+\frac{1}{2}}^{n+1} \coth\left(\frac{-z_u h_{i+\frac{1}{2}} d\Psi_{i+\frac{1}{2}}^{n+1}}{2}\right) \left(u_{i+1}^{n+1} - u_i^{n+1}\right)^2,$$

$$B_3 = \sum_{n=0}^{N_T-1} \Delta t \left( u_{I+1}^{n+1} \mathscr{F}_{u,\,I+\frac{1}{2}}^{n+1} - u_0^{n+1} \mathscr{F}_{u,\,\frac{1}{2}}^{n+1} \right).$$

Using (6) and (8a), (8b) we get:

$$B_1 = \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \frac{\Delta t z_u h_i}{2} \frac{z_P \left( P_i^{n+1} - P^m \right) + z_N \left( N_i^{n+1} - N^m \right)}{\lambda^2} \left( u_i^{n+1} \right)^2$$

$$+ \sum_{n=0}^{N_T-1} \frac{\Delta t z_u}{2} \left( \mathrm{d}\Psi_{I+\frac{1}{2}}^{n+1} \left( u_{I+1}^{n+1} \right)^2 - \mathrm{d}\Psi_{\frac{1}{2}}^{n+1} \left( u_0^{n+1} \right)^2 \right),$$

and thanks to the $L^\infty$-estimates (12)

$$B_1 \geqslant \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \frac{\Delta t z_u^2 h_i}{2\lambda^2} \left( u_i^{n+1} - u^m \right) \left( u_i^{n+1} \right)^2$$

$$+ \sum_{n=0}^{N_T-1} \frac{\Delta t z_u}{2} \left( \mathrm{d}\Psi_{I+\frac{1}{2}}^{n+1} \left( u_{I+1}^{n+1} \right)^2 - \mathrm{d}\Psi_{\frac{1}{2}}^{n+1} \left( u_0^{n+1} \right)^2 \right).$$

As in [3], using $x \coth(x) \geqslant 1$ for all $x \in \mathbb{R}$, we obtain

$$B_2 \geqslant \sum_{n=0}^{N_T-1} \sum_{i=0}^{I} \frac{\Delta t}{h_{i+\frac{1}{2}}} \left( u_{i+1}^{n+1} - u_i^{n+1} \right)^2.$$

Then, we get

$$B_1 + B_2 + B_3 \geqslant \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} \frac{\Delta t z_u^2 h_i}{2\lambda^2} \left( u_i^{n+1} - u^m \right) \left( u_i^{n+1} \right)^2$$

$$+ \sum_{n=0}^{N_T-1} \sum_{i=0}^{I} \frac{\Delta t}{h_{i+\frac{1}{2}}} \left( u_{i+1}^{n+1} - u_i^{n+1} \right)^2 - (f_u^0 + f_u^1), \qquad (18)$$

with

$$f_u^0 = \left( u_0^{n+1} \right)^2 \left[ -\beta_u^0 \left( \psi_0^{n+1} \right) + \frac{z_u}{2} \frac{\psi_0^{n+1} - \Delta\psi_0^{pzc}}{\alpha_0} \right] + u_0^{n+1} \gamma_u^0 \left( \psi_0^{n+1} \right),$$

$$f_u^1 = \left( u_{I+1}^{n+1} \right)^2 \left[ -\beta_u^1 \left( V - \psi_{I+1}^{n+1} \right) - \frac{z_u}{2} \frac{V - \psi_{I+1}^{n+1} - \Delta\psi_1^{pzc}}{\alpha_1} \right] + u_{I+1}^{n+1} \gamma_u^1 \left( V - \psi_{I+1}^{n+1} \right).$$

Then, it remains to find an upper bound of $f_u^0 + f_u^1$. To this end, let us introduce:

$$\xi_u^0(x) = \gamma_u^0(x) - u^m \beta_u^0(x) + u^m \frac{z_u}{\alpha_0} \left(x - \Delta\psi_0^{pzc}\right), \qquad \forall x \in \mathbb{R},$$

$$\xi_u^1(x) = \gamma_u^1(x) - u^m \beta_u^1(x) - u^m \frac{z_u}{\alpha_1} \left(x - \Delta\psi_1^{pzc}\right), \qquad \forall x \in \mathbb{R}.$$

It permits to rewrite:

$$f_u^0 = \frac{\left(u_0^{n+1}\right)^2}{2u^m} \left[\xi_u^0\left(\psi_0^{n+1}\right) - u^m \beta_u^0\left(\psi_0^{n+1}\right) - \gamma_u^0\left(\psi_0^{n+1}\right)\right] + \gamma_u^0\left(\psi_0^{n+1}\right) u_0^{n+1},$$

$$f_u^1 = \frac{\left(u_{I+1}^{n+1}\right)^2}{2u^m} \left[\xi_u^1\left(V - \psi_{I+1}^{n+1}\right) - u^m \beta_u^1\left(V - \psi_{I+1}^{n+1}\right) - \gamma_u^1\left(V - \psi_{I+1}^{n+1}\right)\right]$$
$$+ \gamma_u^1\left(V - \psi_{I+1}^{n+1}\right) u_{I+1}^{n+1}.$$

As shown in [2], hypotheses (7a), (7b) ensure that $\xi_u^0$ and $\xi_u^1$ are nonpositive functions on $\mathbb{R}$. $\beta_u^0$, $\beta_u^1$, $\gamma_u^0$ and $\gamma_u^1$ are nonnegative functions. Then, using (12) and (13) and continuity of $\gamma_u^0$ and $\gamma_u^1$, we obtain

$$f_u^0 + f_u^1 \leqslant \gamma_u^0\left(\psi_0^{n+1}\right) u_0^{n+1} + \gamma_u^1\left(V - \psi_{I+1}^{n+1}\right) u_{I+1}^{n+1} \leqslant C,$$

It leads to:

$$\sum_{n=0}^{N_T-1} \sum_{i=0}^{I} \Delta t \frac{\left(u_{i+1}^{n+1} - u_i^{n+1}\right)^2}{h_{i+\frac{1}{2}}} + \frac{\varepsilon_u}{2} \sum_{n=0}^{N_T-1} \sum_{i=1}^{I} h_i \left(u_i^{n+1} - u_i^n\right)^2 \leqslant C, \qquad (19)$$

with $C$ depending on the data and independent of $\Delta t$ and $h$, which concludes the proof.

## 6 Conclusion

The a priori estimates (13)–(15) give us the compactness in space of the sequences of approximate solutions. The compactness in time is obtained for instance by discrete Aubin Simon compactness lemma (see [6]). We can already note that (19) also give an estimate in time for $P$ and $N$ (which holds only for $\varepsilon_N, \varepsilon_P > 0$). Then, passing to the limit in the numerical scheme, we obtain that $(P, N, \Psi)$ is a weak solution of (1a)–(5) (see [3]). In this step, we also have to pay attention to the boundary terms. This is not detailed in this paper but the integrality of the proof will be presented in [4].

In this paper, the convergence proof works for $\varepsilon_N$, $\varepsilon_P > 0$. But the dimensionless parameters in (1a) are in practice $\varepsilon_P = 1$ and $\varepsilon_N \ll 1$. It could be set to 0 in the model. Existence of solutions to this new model and convergence of a numerical scheme is still an open question. Future work will focus on it.

## References

1. Bataillon, C., Bouchon, F., Chainais-Hillairet, C., Desgranges, C., Hoarau, E., Martin, F., Perrin, S., Turpin, M., Talandier, J.: Corrosion modelling of iron based alloy in nuclear waste repository. Electrochim. Acta **55**(15), 4451–4467 (2010)
2. Bataillon, C., Bouchon, F., Chainais-Hillairet, C., Fuhrmann, J., Hoarau, E., Touzani, R.: Numerical methods for the simulation of a corrosion model with moving oxide layer. J. Comput. Phys. **231**(18), 6213–6231 (2012)
3. Bessemoulin-Chatard, M.: A finite volume scheme for convection-diffusion equations with non-linear diffusion derived from the Scharfetter-Gummel scheme. Numer. Math. **121**(4), 637–670 (2012)
4. Chainais-Hillairet, C., Colin, P.L., Lacroix-Violet, I.: Convergence of a finite volume scheme for a corrosion model. In preparation (2014)
5. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. Hand. Numer. Anal. **7**, 713–1018 (2000)
6. Gallouët, T., Latché, J.-C.: Compactness of discrete approximate solutions to parabolic pdes—application to a turbulence model. Commun. Pure Appl. Anal. **11**(6), 2371–2391 (2012)
7. Scharfetter, D., Gummel, H.: Large-signal analysis of a silicon read diode oscillator. Electron. Devices, IEEE Trans. **16**(1), 64–77 (1969)

# High Performance Computing Linear Algorithms for Two-Phase Flow in Porous Media

**Robert Eymard, Cindy Guichard and Roland Masson**

**Abstract**  We focus here on the difficult problem of linear solving, when considering implicit scheme for two-phase flow simulation in porous media. Indeed, this scheme leads to ill-conditioned linear systems, due to the different behaviors of the pressure unknown (which follows a diffusion equation) and the saturation unknown (mainly advected by the total volumic flow). This difficulty is enhanced by the parallel computing techniques, which reduce the choice of the possible preconditioners. We first present the framework of this study, and then we discuss different algorithms for linear solving. Finally, numerical results show the performances of these algorithms.

## 1 Introduction

We consider the flow of two immiscible compressible phases, the water phase (denoted $w$) and the gas phase (denoted $g$), in porous media; each phase is only composed of one component. In order to characterize the mathematical coupling of diffusion and advection, we consider the case where the capillary pressure effects

R. Eymard
Laboratoire d'Analyse et de Mathématiques Appliquées, CNRS, UPEM, UPEC,
5 boulevard Descartes, Champs-sur-Marne 77454, Marne-la-Vallée Cedex 2, France
e-mail: Robert.Eymard@u-pem.fr

C. Guichard
Laboratoire Jacques-Louis Lions, CNRS, UMR 7598, Sorbonne Universités,
UPMC Univ Paris 06, F-75005 Paris, France
e-mail: guichard@ljll.math.upmc.fr

R. Masson (✉)
Laboratoire de Mathématiques J.A. Dieudonné, UMR CNRS 7251 and team Coffee,
Université Nice Sophia Antipolis, CNRS and INRIA Sophia Antipolis Méditerranée,
Parc Valrose, 06108 Nice, France
e-mail: roland.masson@unice.fr

can be neglected in front of the high level of pressure gradients imposed by the production and injection wells. The mass conservation equations are therefore the following,

$$\phi \, \partial_t( \, \rho_\alpha(P) \, S_\alpha \, ) \, + \, \mathrm{div} \, ( \, \rho_\alpha(P) \, \mathbf{V}_\alpha \, ) \, = \, Q_\alpha, \quad \alpha = w, g \,, \tag{1}$$

together with the generalized Darcy law

$$\mathbf{V}_\alpha \, = \, - \, \frac{k_{r\alpha}(S_\alpha)}{\mu_\alpha(P)} \, \Lambda \, ( \, \nabla P \, - \, \rho_\alpha(P) \, \mathbf{g} \, ), \quad \alpha = w, g \,. \tag{2}$$

In Eqs. (1) and (2), the main unknowns are the pressure $P$ and one saturation, for example $S_w$, since the phase saturations are linked by $S_w + S_g = 1$. Additionally, $\phi$ is the porosity, $\Lambda$ is the absolute permeability tensor (these values only depending on the rock material), $\mathbf{g}$ is the gravity acceleration, and, for each phase $\alpha = w, g$, $\rho_\alpha$ represents the bulk density, $k_{r\alpha}$ is the relative permeability (nonnegative increasing function with respect to $S_\alpha$), $\mu_\alpha$ is the viscosity and $Q_\alpha$ is the source term that represents the contribution of the wells. The ratio $\frac{k_{r\alpha}}{\mu_\alpha}$ is called the mobility of the phase $\alpha$. These equations are considered in a time-space domain $(0, t_f) \times \Omega$, where $\Omega$ is a polygonal open bounded and connected subset of $\mathbb{R}^3$, and $t_f > 0$ is the time duration of the simulation. Finally, these equations are considered together with homogeneous Neumann conditions at the boundary of the domain $\Omega$, and initial conditions on the pressure and on the saturation.

The approximation of the solution to (1) and (2) in the industrial framework with large time and space scales, requires High Performance Computing techniques. This implies to handle the difficult problem of solving the linear systems which arise from fully coupled schemes and domain decomposition, using multi-threading algorithms: these schemes happen to be the only ones used for the approximation of (1) and (2) in the industrial framework. For this purpose, we consider here the use of PETSc [5] together with external preconditioners libraries like MUMPS [4] and HYPRE [2]. Note that other packages, like DUNE [1], are available. In the example of a gas storage case [10], a very good scalability has been observed using Boomer AMG [2] as a preconditioner on the full system, although AMG is usually not adapted to solve the full system but only the pressure elliptic block. But for more general situations of two-phase flow (such as the case considered in the numerical example of this paper), this strategy fails. This has led to the development of efficient Combinative-AMG preconditioners [11], combining typically an AMG preconditioner on the pressure block with an ILU preconditioner on the full system. This paper focuses on an alternative algorithm, based on the PETSc environment, for the resolution of the linear systems issued from a fully implicit scheme for the approximation of (1) and (2). Its main advantages are the following:

1. It makes a bridge between sequential and fully implicit schemes.
2. It leads to the sequential use of robust solvers suited for the nature of each unknown.

This paper is organized as follows. In Sect. 2, we present a discretization scheme and its parallel implementation. We then discuss in Sect. 3 the fix-point methods used for the approximation of the solution to the nonlinear systems of Eqs. (3) and (4). Some numerical results, in Sect. 4, illustrate our method.

## 2 Discretization and Parallel Implementation

In order to study the algorithms for solving the linear systems in a parallel framework, we have extended to the two-phase flow model a recent work (see [6]) done for a linear parabolic equation. For the implementation details, we use below the same notations as [6], thus we focus on the specific points regarding the discretization of two-phase flow. Hence, the continuous model (1) and (2) is discretized using an Euler fully implicit method in time, and the VAG scheme (Vertex Approximate Gradient scheme introduced in [7]) in space with up-winding of the mobilities according to the sign of the Darcy fluxes. We emphasize that the VAG scheme is a symmetric scheme based on a hybrid formulation, both in terms of vertices and cells unknowns, but in the resulting linear system the cell unknowns are algebraically eliminated without any fill-in. The VAG scheme involves linear fluxes between a cell and its vertices and its implementation matches with that of a standard Multi-Points Flux Approximation. We refer to [8, 9] for details on the VAG scheme for multiphase flow in porous media in the case of a sequential implementation.

**Parallel discretization.** We consider a mesh of the domain $\Omega$ (the elements of the mesh are called cells in the following). As in [6, Sect. 2.1], we denote the set of processes by $\mathscr{P}$, and we consider a partition of the mesh. For a given process $p \in \mathscr{P}$, we denote by $\mathscr{M}^p$ the set of its own cells (in practice selected by applying the Metis package [3]) and by $\overline{\mathscr{M}}^p$ the set of its overlapped cells which is defined as the set of all cells sharing a vertex with $\mathscr{M}^p$. Then we can define the overlapping decomposition of the set of vertices as follows:

$$\overline{\mathscr{V}}^p = \bigcup_{K \in \overline{\mathscr{M}}^p} \mathscr{V}_K, \quad p \in \mathscr{P},$$

where $\mathscr{V}_K$ is the set of the vertices of a given cell $K$. Finally, the set of the own vertices of a process $p \in \mathscr{P}$, denoted $\mathscr{V}^p$, is obtained by the application of a rule detailed in [6, Sect. 2.1]. We then discretize the continuous Eqs. (1) and (2) on each process $p$, for each phase $\alpha = w, g$, by writing

$$\frac{|\mathbf{s}|}{\delta t^{(n)}} \left( \rho_\alpha(P_\mathbf{s}^{(n+1)}) S_{\alpha,\mathbf{s}}^{(n+1)} - \rho_\alpha(P_\mathbf{s}^{(n)}) S_{\alpha,\mathbf{s}}^{(n)} \right) - \sum_{K \in \mathscr{M}_\mathbf{s}} M_{\alpha,K\mathbf{s}}^{\mathrm{up},(*)} V_{\alpha,K \to \mathbf{s}}^{(n+1)} = 0$$

$$\forall \mathbf{s} \in \mathscr{V}^p, \tag{3a}$$

$$\frac{|K|}{\delta t^{(n)}} \left( \rho_\alpha(P_K^{(n+1)}) S_{\alpha,K}^{(n+1)} - \rho_\alpha(P_K^{(n)}) S_{\alpha,s}^{(n)} \right) + \sum_{\mathbf{s} \in \mathcal{V}_K} M_{\alpha,K\mathbf{s}}^{\mathrm{up},(*)} V_{\alpha,K \to \mathbf{s}}^{(n+1)} = |K| Q_{\alpha,K}^{(n+1)}$$

$$\forall K \in \overline{\mathcal{M}}^p, \qquad (3b)$$

together with the Darcy fluxes (see [8, Sects. 3.1.2 and 3.2])

$$V_{\alpha,K \to \mathbf{s}}^{(n+1)} = \sum_{\mathbf{s}' \in \mathcal{V}_K} a_{K,\mathbf{s}}^{\mathbf{s}'} \left( P_K^{(n+1)} - P_{\mathbf{s}'}^{(n+1)} + \rho_\alpha(P_K^{(n)}) \mathbf{g} \cdot (x_K - x_{\mathbf{s}'}) \right). \qquad (4)$$

In (3) and (4), we use the following notations: $\mathcal{M}_\mathbf{s}$ is the set of cells $K$ such that $\mathbf{s} \in \mathcal{V}_K$, $|K|$ (resp. $|\mathbf{s}|$) is the porous volume associated to a cell $K$ (resp. to a vertex $\mathbf{s}$), computed from a redistribution of the total porous volume of the space domain with respect to the mesh and the rock type properties [8, 9], $x_K \in \mathbb{R}^3$ (resp. $x_\mathbf{s} \in \mathbb{R}^3$) denotes the coordinates of the center of the cell $K$ (resp. of the vertex $\mathbf{s}$). Note that the VAG scheme construction does not use the geometry of the control volumes $\mathbf{s} \in \mathcal{V}$ and $K \in \mathcal{M}$ but only their volumes. For $n \in \mathbb{N}$, $\delta t^{(n)} = t^{(n+1)} - t^{(n)}$ is the time step between times $t^{(n+1)}$ and $t^{(n)}$. For any control volume $I$ ($I = K$ or $I = \mathbf{s}$), $P_I^{(n)}$ (resp. $S_{\alpha,I}^{(n)}$) is an approximation of $P$ (resp. of $S_\alpha$) in $I$ at time $t^{(n)}$. $a_{K,\mathbf{s}}^{\mathbf{s}'}$ is computed with respect to the mesh and the permeability tensor $\Lambda$ [8, 9]. $M_{\alpha,K\mathbf{s}}^{\mathrm{up},(*)}$ denotes the upstream mobility of the phase $\alpha$ and is defined by

$$M_{\alpha,K\mathbf{s}}^{\mathrm{up},(*)} = \left( \rho_\alpha(P_{K\mathbf{s}}^{(n)}) \frac{k_{r\alpha}(S_{\alpha,K\mathbf{s}}^{(*)})}{\mu_\alpha(P_{K\mathbf{s}}^{(n)})} \right),$$

where $K\mathbf{s}$ denotes the cell $K$ if $V_{\alpha,K \to \mathbf{s}}^{(n+1)} \geq 0$, or the vertex $\mathbf{s}$ otherwise. The upper index $(*)$ stands for $(n)$ (ImPES scheme) or $(n+1)$ (fully implicit scheme); this point is reviewed in Sect. 3. Finally, $Q_{\alpha,K}^{(n+1)}$ is the possible source term if any well is open through cell $K$.

As usual, no special numerical treatment is needed for taking into account the homogeneous Neumann boundary conditions (see [8, 9]). The set of Eqs. (3) and (4) leads to a system of nonlinear equations at each time step. This system is solved by a fix-point algorithm based on the Newton-Raphson method (up to a possible under-relaxation in order to prevent from nonconvergence behaviors). Thus, the unknowns of the resulting discrete problem are, on each process $p \in \mathcal{P}$, the variations of $(P_I^{(n+1)})_{I \in \overline{\mathcal{V}}^p \cup \overline{\mathcal{M}}^p}$ and $(S_{w,I}^{(n+1)})_{I \in \overline{\mathcal{V}}^p \cup \overline{\mathcal{M}}^p}$ between two fix-point iterations. As in [6, Sect. 2.4], their values are obtained through the construction of rectangular linear systems on each process $p \in \mathcal{P}$ and, as mentioned above, a consequence of Eqs. (3) and (4) is that the cell unknowns can be eliminated by a Schur complement without fill-in, in order to reduce the linear system to the vertices unknowns. Thanks to our general definition of the overlap, the assembling step may be performed locally on each process without communication.

## 3 Fix-Point Methods

This section presents the fix-point method used in our parallel implementation, and its implementation thanks to open-source libraries. The variation of the vertices unknowns between two fix-point iterations is denoted as follows, omitting the time superscript $(n + 1)$ and the Newton iteration index for the sake of clarity,

$$U^p = \left( (\triangle P_\mathbf{s})_{\mathbf{s}\in\overline{\mathcal{V}}^p}, (\triangle S_{w,\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}}^p} \right), \quad \forall p \in \mathscr{P}.$$

If the fix-point method were exactly the Newton method, then $U^p$ would be the solution of a linear system over all processes $p$, under the form

$$\begin{pmatrix} A^p_{PP} & A^p_{PS} \\ A^p_{SP} & A^p_{SS} \end{pmatrix} U^p = \begin{pmatrix} B^p_P \\ B^p_S \end{pmatrix}, \tag{5}$$

where, for $u1, u2 = P, S$, the vectors right-hand-side $B^p_{u1}$ belong to $\mathbb{R}^{\mathcal{V}^p}$ and the sub-matrices $A^p_{u1,u2}$ belong to $\mathbb{R}^{\mathcal{V}^p} \otimes \mathbb{R}^{\overline{\mathcal{V}^p}}$.

Then, we define the diagonal blocks, of size $2 \times 2$, by

$$D^\mathbf{s} = \begin{pmatrix} A^p_{PP}(\mathbf{s}, \mathbf{s}) & A^p_{PS}(\mathbf{s}, \mathbf{s}) \\ A^p_{SP}(\mathbf{s}, \mathbf{s}) & A^p_{SS}(\mathbf{s}, \mathbf{s}) \end{pmatrix}, \quad \forall \mathbf{s} \in \mathcal{V}^p, \quad \forall p \in \mathscr{P},$$

where $A^p_{u1,u2}(\mathbf{s}, \mathbf{s})$ is the term associated to the equation on $\mathbf{s}$ and the unknown on $\mathbf{s}$ of the sub-matrix $A^p_{u1,u2}$. We then left-multiply the system (5) by the square matrix $[\mathrm{diag}(D^\mathbf{s}, \mathbf{s} \in \mathcal{V}^p)]^{-1} \in \mathbb{R}^{(\mathcal{V}^p)^2} \otimes \mathbb{R}^{(\mathcal{V}^p)^2}$. We then get the following linear system

$$\begin{pmatrix} \widehat{A}^p_{PP} & \widehat{A}^p_{PS} \\ \widehat{A}^p_{SP} & \widehat{A}^p_{SS} \end{pmatrix} U^p = \begin{pmatrix} \widehat{B}^p_P \\ \widehat{B}^p_S \end{pmatrix}, \tag{6}$$

where the new right-hand-side and sub-matrices have the same dimension as in (5) but now satisfy

$$\widehat{A}^p_{PP}(\mathbf{s}, \mathbf{s}) = 1, \ \widehat{A}^p_{PS}(\mathbf{s}, \mathbf{s}) = 0, \ \widehat{A}^p_{SP}(\mathbf{s}, \mathbf{s}) = 0, \ \widehat{A}^p_{SS}(\mathbf{s}, \mathbf{s}) = 1, \tag{7}$$

for all vertices $\mathbf{s} \in \mathcal{V}^p$, and for any process $p \in \mathscr{P}$.

Our implementation of the sequential fix-point scheme allows the choice between the three following schemes (in all cases, the solution at own and ghost cells is computed locally on each process $p \in \mathscr{P}$ by Schur complement).

1. The **ImPES scheme**, for Implicit in Pressure and Explicit in Saturation, is obtained by taking $(*) = (n)$ in (3). This means that the linear system, under the form (6), is such that $\widehat{A}^p_{PS} = 0$ and $\widehat{A}^p_{SS}(\mathbf{s}, \mathbf{s}) = 1$ if $\mathbf{s} \in \mathcal{V}^p$, and $\widehat{A}^p_{SS}(\cdot, \cdot) = 0$ otherwise. Hence the resolution of (6) first implies a full resolution on all vertices of the mesh for the equation

$$\widehat{A}^p_{PP}(\triangle P_{\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}} = \widehat{B}_p. \tag{8}$$

This sub-system is transferred line by line to PETSc which provides the solution vector on pressure variations $(\triangle P_{\mathbf{s}})_{\mathbf{s}\in\mathcal{V}^p}$ at own vertices for each process $p \in \mathcal{P}$. The pressure variations at own and ghost vertices $(\triangle P_{\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}}$, $p \in \mathcal{P}$ is obtained by a synchronization of the vertices. Then the saturation variations at own vertices are immediately obtained from (6) and the saturation variations at ghost vertices are obtained by a second synchronization of the vertices. Note that the matrix $A^p_{PP}$ has the main properties of a finite element matrix for a diffusion problem. Therefore, in PETSc, it is standard to use Algebraic Multi-Grid preconditioning for the resolution of these linear systems. Unfortunately, the ImPES scheme implies a limit on the time step, for standard stability reasons, which is generally not compatible with industrial requirements.

2. The **fully implicit scheme** is obtained by taking $(*) = (n + 1)$ in (3) (which implies that the matrices $\widehat{A}^p_{PS}$ do no longer vanish) and then to solve the coupled linear system (6) with a unique linear solver issued from the PETSc library. In this case, PETSc provides the solution on both pressure and saturation variations at own vertices for each process $p \in \mathcal{P}$. The variations at the unknowns at the ghost vertices are then obtained by synchronization. As discussed in the introduction of this paper, this strategy implies the implementation of Combinative preconditioners [11]. The following sequential scheme proposes a related but simpler to implement approach.

3. The **sequential scheme** consists in only approximately solving the linear systems (6), thanks to the following algorithm, based on the combination of adapted preconditioners on both the pressure and saturation blocks of the full system. For any process $p \in \mathcal{P}$, we consider the Gauss-Seidel type method,

$$\forall \mathbf{s} \in \overline{\mathcal{V}^p}, \ \triangle S^{\{0\}}_{w,\mathbf{s}} = 0, \tag{9a}$$

$$\widehat{A}^p_{PP}(\triangle P^{\{k+1\}}_{\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}} + \widehat{A}^p_{PS}(\triangle S^{\{k\}}_{w,\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}} = \widehat{B}_p \tag{9b}$$

$$\widehat{A}^p_{SP}(\triangle P^{\{k+1\}}_{\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}} + \widehat{A}_{SS}(\triangle S^{\{k+1\}}_{w,\mathbf{s}})_{\mathbf{s}\in\overline{\mathcal{V}^p}} = \widehat{B}_s, \tag{9c}$$

where in (9a)–(9c), the upper index $k = 0, \ldots, \mathbf{M}$ is corresponding to the sequential scheme iterations. Hence, if the integer $\mathbf{M}$ is large enough, this algorithm leads to a fix point method for solving the linear systems (6) issued from Newton's method applied to the fully implicit scheme. Nevertheless, we take in practice $\mathbf{M} \leq 5$ to ensure a reasonable wall clock time (denoted by WCT in the tables).

The resolution of (9a)–(9c), at each scheme iteration $k + 1$, requires to solve the first parallel linear system (9b) which has the same skeleton as (8) issued from the ImPES scheme. Its resolution by PETSc provides the solution vector on pressure variations $(\triangle P^{\{k+1\}}_{\mathbf{s}})_{\mathbf{s}\in\mathcal{V}^p}$ at own vertices for each process $p \in \mathcal{P}$. The pressure variations at own and ghost vertices is still obtained by a synchronization of the vertices. Then the saturation variations is obtained by solving (9c). This implies the assembling of a second parallel linear system, once again solved by PETSc. This provides the saturation variations at own vertices for each process $p \in \mathcal{P}$,

**Fig. 1** *Left* residual with respect to the Newton iteration at the first time step of $401 \times 401 \times 12$ mesh. *Right* random log normal heterogeneous permeability tensor on the same mesh

and a synchronization step concludes the iteration of the method. In terms of cost of communication, this sequential method is close to that resulting from the fully implicit scheme, if the total (sum on $k$) number of iterations of its two successive solvers is close to the number of iterations of the unique linear solver used for the fully implicit method. To achieve this, a very efficient strategy is then to specify the residual tolerance $\varepsilon(k)$ of these two resolutions with respect to the value of $k$. We have implemented the relation $\varepsilon(k+1) = \varepsilon(k)^2$. This leads to a very small number of linear solver iterations for the first values of $k$.

## 4 Numerical Results

We now consider a two-phase flow on a 3D "five-spot pattern", i.e. with 4 vertical injection wells (at each corner of the domain) and 1 vertical production well (at the center of the domain). The geometry and the permeability field of the test case is illustrated by Fig. 1 (*right side*). The cluster used is composed of 32 processors Intel® Xeon® CPU E5-4620 with frequency 2.20 GHz. Three successive meshes have been built (with resp. $101 \times 101 \times 12$, $201 \times 201 \times 12$ and $401 \times 401 \times 12$ cells), and we focus on the beginning of the simulation where the pressure and saturation variations are the highest. Referring to Sect. 3: firstly, the ImPES scheme is not efficient in front of the highest variations of the unknowns; secondly, we did not find any efficient preconditioner for the coupled system in the PETSc framework, involving a strong motivation for exploring the sequential algorithm which is simpler to implement than a Combinative-AMG preconditioner.

Let us first comment the results obtained with 32-processors runs. We present in Fig. 1 (*left side*) the residual in function of the Newton iterations, for different values of **M** (see Sect. 3). This figure shows that the convergence rate is only linear for **M** = 1, and is more and more quadratic as **M** increases, since the scheme becomes closer to the pure Newton method. Table 1 exhibits the wall clock times, the total number of Newton iterations, and of solver iterations for the pressure resolution and the saturation resolution. For the ill-conditioned pressure block, the preconditioner is

**Table 1** Results on the three meshes with 32 processors: "WCT init" denotes the wall clock time for initialization operations (in particular including mesh reading and partitioning), "WCT Newton" is the wall clock time for nonlinear iterations with total number ♯Newton, ♯iter/P (resp. /S) total number of pressure (resp. saturation) linear iterations

| WCT init (s) | 6 time steps on the 101 × 101 × 12 mesh | | | | | 7 time steps on the 201 × 201 × 12 mesh | | | | | 4 time steps on the 401 × 401 × 12 mesh | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | 8 | | | | | 27 | | | | | 110 | | | | |
| M | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 |
| WCT Newton (s) | 46 | 38 | 41 | 61 | 72 | 353 | 337 | 309 | 391 | 449 | 1831 | 1638 | 1631 | 1905 | 2252 |
| ♯Newton | 57 | 42 | 38 | 38 | 38 | 85 | 53 | 44 | 45 | 44 | 65 | 40 | 33 | 31 | 32 |
| ♯iter/P | 739 | 737 | 717 | 788 | 809 | 699 | 585 | 558 | 638 | 654 | 722 | 577 | 520 | 540 | 593 |
| ♯iter/S | 143 | 161 | 168 | 202 | 239 | 255 | 250 | 238 | 275 | 308 | 239 | 228 | 210 | 210 | 245 |

**Table 2** Results for 7 time steps on the $201 \times 201 \times 12$ mesh and $\mathbf{M} = 3$ (the total Wall Clock Time includes the initialization and Newton iteration times)

| ♯proc. | 1 | 2 | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|---|
| Total WCT (s) | 3214 | 1743 | 1039 | 655 | 460 | 336 |
| ♯Newton | 45 | 45 | 45 | 45 | 45 | 44 |
| ♯iter/P | 569 | 560 | 570 | 562 | 562 | 558 |
| ♯iter/S | 243 | 244 | 243 | 244 | 243 | 238 |

1 V-cycle of boomer AMG of HYPRE with Gauss Seidel relaxation, whereas for the much better conditionned saturation block, we selected the Jacobi preconditioner. The final residual reduction specified for each linear resolution is $\varepsilon(\mathbf{M}) = 10^{-5}$. For these numerical tests we imposed a reduction of Newton residual equal to $10^{-6}$. This criterion has been selected in order to impose a high precision in the nonlinear resolution, hence indicating the robustness of the algorithm with respect to severe convergence requirements. We observe that an optimum, both with respect to the wall clock time and the number of Newton iterations, is obtained with $\mathbf{M} = 3$.

The scalability results presented in Table 2 are similar to those of [10]. The parallel efficiency is reduced from 16 to 32 processors due to a too small number of unknowns per processor in the AMG preconditioner for this problem size (see [10]). They also show a very good stability of the linear algorithms with respect to the increase of the number of processors.

# References

1. DUNE Distributed and Unified Numerics Environment. http://www.dune-project.org/
2. Hypre Parallel High Performance Preconditioners. http://acts.nersc.gov/hypre/
3. Metis Serial Graph Partitioning and Fill-reducing Matrix Ordering. http://glaros.dtc.umn.edu/gkhome/views/metis
4. MUMPS MUltifrontal Massively Parallel sparse direct Solver. http://mumps.enseeiht.fr/
5. PETSc Portable, Extensible Toolkit for Scientific Computation. http://www.mcs.anl.gov/petsc
6. Dalissier, E., Guichard, C., Have, P., Masson, R., Yang, C.: ComPASS : a tool for distributed parallel finite volume discretizations on general unstructured polyhedral meshes. ESAIM: Proc. **43**, 147–163 (2013)
7. Eymard, R., Guichard, C., Herbin, R.: Small-stencil 3d schemes for diffusive flows in porous media. ESAIM. Math. Model. Numer. Anal. **46**, 265–290 (2012)
8. Eymard, R., Guichard, C., Herbin, R., Masson, R.: Vertex-centred discretization of multiphase compositional Darcy flows on general meshes. Comput. Geosci. **16**(4), 987–1005 (2012)
9. Eymard, R., Guichard, C., Herbin, R., Masson, R.: Vertex centred discretization of two-phase Darcy flows on general meshes. ESAIM: Proc. **35**, 59–78 (2012)
10. Eymard, R., Guichard, C., Masson, R.: Simulation of two-phase Darcy flow using the VAG scheme on parallel architecture. poster at MoMaS Multiphase Seminar Days (oct. 2013), available online: http://math.unice.fr/massonr/films/Compass.pdf
11. Scheichl, R., Masson, R., Johannes, W.: Decoupling and block preconditioning for sedimentary basin simulations. Comput. Geosci. **7**, 295–318 (2003)

# Numerical Solution of Fluid-Structure Interaction by the Space-Time Discontinuous Galerkin Method

**Miloslav Feistauer, Martin Hadrava, Jaromír Horáček and Adam Kosík**

**Abstract** This paper is devoted to the numerical solution of the interaction of compressible viscous flow with elastic structures. The flow in a time-dependent domain is described by the compressible Navier-Stokes equations written in the ALE formulation and the deformation of elastic structures is described by the dynamic linear elasticity system. For each individual problem we employ the discretization by the space-time discontinuous Galerkin finite element method (ST-DGM). The flow and elasticity problems are coupled via transmission conditions. The developed method is tested by numerical experiments.

## 1 Formulation of the Problem

### 1.1 Flow Problem

We are concerned with the problem of compressible flow in a time-dependent bounded domain $\Omega_t \subset I\!\!R^2$ with $t \in [0, T]$. The boundary of $\Omega_t$ is formed by

M. Feistauer (✉) · M. Hadrava · A. Kosík
Faculty of Mathematics and Physics, Charles University in Prague, Sokolovská
83,186 75 Praha 8, Czech Republic
e-mail: feist@karlin.mff.cuni.cz

M. Hadrava
e-mail: martin@hadrava.eu

A. Kosík
e-mail: adam.kosik@atlas.cz

J. Horáček
Institute of Thermomechanics, The Academy of Sciences of the Czech Republic,
v. v. i., Dolejškova 1402/5,182 00 Praha 8, Czech Republic
e-mail: jaromirh@it.cas.cz

three disjoint parts: $\partial\Omega_t = \Gamma_I \cup \Gamma_O \cup \Gamma_{W_t}$, where $\Gamma_I$ is the inlet, $\Gamma_O$ is the outlet and $\Gamma_{W_t}$ represents impermeable time-dependent walls.

The time dependence of the domain $\Omega_t$ is taken into account with the aid of the *Arbitrary Lagrangian-Eulerian* (ALE) method (see, e.g., [4]). It is based on a regular one-to-one ALE mapping of the reference configuration $\Omega_0$ onto the current configuration $\Omega_t : \mathcal{A}_t : \bar{\Omega}_0 \longrightarrow \bar{\Omega}_t$, i.e. $X \in \bar{\Omega}_0 \longmapsto x = x(X, t) = \mathcal{A}_t(X) \in \bar{\Omega}_t$. Further, we define the domain velocity $\tilde{z}(X, t) = \frac{\partial}{\partial t}\mathcal{A}_t(X)$, $t \in [0, T]$, $X \in \Omega_0$, $z(x, t) = \tilde{z}(\mathcal{A}_t^{-1}(x), t)$, $t \in [0, T]$, $x \in \Omega_t$ and the ALE derivative of the state vector function $w = w(x, t)$ defined for $x \in \Omega_t$ and $t \in [0, T]$: $\frac{D^{\mathcal{A}}}{Dt}w(x, t) = \frac{\partial\tilde{w}}{\partial t}(X, t)$, $\tilde{w}(X, t) = w(\mathcal{A}_t(X), t)$, $X \in \Omega_0$, $x = \mathcal{A}_t(X)$. Then the continuity equation, the Navier-Stokes equations and the energy equation can be written in the ALE form

$$\frac{D^{\mathcal{A}}w}{Dt} + \sum_{s=1}^{2}\frac{\partial g_s(w)}{\partial x_s} + w\,\mathrm{div}z = \sum_{s=1}^{2}\frac{\partial R_s(w, \nabla w)}{\partial x_s}, \tag{1}$$

where $w = (\rho, \rho v_1, \rho v_2, E)^T \in I\!\!R^4$, $g_s(w) = f(w)_s - z_s w$, $f_s = (\rho v_s, \rho v_1 v_s + \delta_{1s}p, \rho v_2 v_s + \delta_{2s}p, (E + p)v_s)^T$, $R_s(w, \nabla w) = (0, \tau_{s1}^V, \tau_{s2}^V, \tau_{s1}^V v_1 + \tau_{s2}^V v_2 + k\frac{\partial\theta}{\partial x_s})^T$, $s = 1, 2$, $\tau_{ij}^V = \lambda\delta_{ij}\mathrm{div}v + 2\mu d_{ij}(v)$, $d_{ij}(v) = \frac{1}{2}\left(\frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i}\right)$, $i, j = 1, 2$. We have $R_s(w, \nabla w) = \sum_{k=1}^{2}\mathbb{K}_{s,k}(w)\frac{\partial w}{\partial x_k}$, where $\mathbb{K}_{s,k}(w)$ are $4 \times 4$ matrices depending on $w$, and $f_s(w) = \mathbb{A}(w)w$ with $\mathbb{A}(w) = Df_s(w)/Dw$.

The following notation is used: $\rho$—fluid density, $p$—pressure, $E$—total energy, $v = (v_1, v_2)$—velocity vector, $\theta$—absolute temperature, $c_v > 0$—specific heat at constant volume, $\gamma > 1$—Poisson adiabatic constant, $\mu > 0, \lambda = -2\mu/3$—viscosity coefficients, $k > 0$—heat conduction coefficient, $\tau_{ij}^V$—components of the viscous part of the stress tensor. System (1) is completed by the thermodynamical relations $p = (\gamma - 1)\left(E - \rho\frac{|v|^2}{2}\right)$, $\theta = \frac{1}{c_v}\left(\frac{E}{\rho} - \frac{|v|^2}{2}\right)$ and equipped with the initial condition $w(x, 0) = w^0(x)$, $x \in \Omega_0$ and the boundary conditions:
$\rho = \rho_D$, $v = v_D$, $\sum_{j=1}^{2}\left(\sum_{i=1}^{2}\tau_{ij}^V n_i\right)v_j + k\frac{\partial\theta}{\partial n} = 0$ on the inlet $\Gamma_I$,
$v = z_D(t) =$ velocity of a moving wall, $\frac{\partial\theta}{\partial n} = 0$, on the moving wall $\Gamma_{W_t}$,
$\sum_{j=1}^{2}\tau_{ij}^V n_j = 0$, $\frac{\partial\theta}{\partial n} = 0$, $i = 1, 2$, on the outlet $\Gamma_O$,
with prescribed data $\rho_D$, $v_D$, $z_D$. By $n$ we denote the unit outer normal.

## 1.2 Elasticity Problem

We consider an elastic body $\Omega^b \subset I\!\!R^2$, which has a common boundary $\Gamma_N^b$ with the reference domain $\Omega_0$ occupied by the fluid at the initial time. Further, the boundary of $\Omega^b$ is formed by two disjoint parts $\partial\Omega^b = \Gamma_N^b \cup \Gamma_D^b$, $\Gamma_N^b \subset \Gamma_{W_0}$ and $\Gamma_D^b$ is a fixed part of the boundary. Using the notation of the displacement of the body $u = u(X, t)$, $X \in \Omega^b$, $t \in (0, T)$ we can write the equations describing the defor-

mation of the elastic body $\Omega^b$ in the form

$$\rho^b \frac{\partial^2 \boldsymbol{u}}{\partial t^2} + c_M \rho^b \frac{\partial \boldsymbol{u}}{\partial t} - div\, \boldsymbol{\sigma}(\boldsymbol{u}) - c_K \frac{\partial}{\partial t} div\, \boldsymbol{\sigma}(\boldsymbol{u}) = \boldsymbol{f} \quad \text{in } \Omega^b \times (0, T), \quad (2)$$

$$\boldsymbol{u} = \boldsymbol{u}_D \quad \text{in } \Gamma^b_D \times (0, T), \quad \boldsymbol{\sigma}(\boldsymbol{u}) \cdot \boldsymbol{n} = \boldsymbol{g}_N \quad \text{in } \Gamma^b_N \times (0, T), \quad (3)$$

$$\boldsymbol{u}(x, 0) = \boldsymbol{u}_0(x), \quad x \in \Omega^b, \quad \frac{\partial \boldsymbol{u}}{\partial t}(x, 0) = z_0(x), \quad x \in \Omega^b. \quad (4)$$

Here $\boldsymbol{\sigma}(\boldsymbol{u}) = \{\sigma_{ij}\}^2_{i,j=1}$, $\sigma_{ij} = \lambda^b div \boldsymbol{u} \delta_{ij} + 2\mu^b e^b_{ij}(\boldsymbol{u})$ with $e^b_{ij}(\boldsymbol{u}) = (\partial u_i/\partial x_j + \partial u_j/\partial x_i)/2$. Further, $\boldsymbol{f} : \Omega^b \times (0, T) \rightarrow \mathbb{R}^2$—outer volume force, $\boldsymbol{u}_D : \Gamma^b_D \times (0, T) \rightarrow \mathbb{R}^2$—boundary displacement, $\boldsymbol{g}_N : \Gamma^b_N \times (0, T) \rightarrow \mathbb{R}^2$—boundary normal stress, $\boldsymbol{u}_0 : \Omega^b \rightarrow \mathbb{R}^2$—initial displacement, $z_0 : \Omega^b \rightarrow \mathbb{R}^2$—initial deformation velocity and $\rho^b > 0$—material density are given functions. The expressions $c_M \rho^b \frac{\partial \boldsymbol{u}}{\partial t}$ and $c_K \frac{\partial}{\partial t} div\, \boldsymbol{\sigma}(\boldsymbol{u})$ represent the damping terms, with $c_M, c_K \geq 0$.

The flow and structural problems are coupled by the transmission conditions

$$\boldsymbol{v} = \frac{\partial \boldsymbol{u}}{\partial t}, \quad \sum^2_{j=1} \sigma_{ij}(\boldsymbol{X}, t) n_j(\boldsymbol{X}) = -\sum^2_{j=1} \tau^f_{ij}(\boldsymbol{x}, t) n_j(\boldsymbol{X}), \quad i = 1, 2, \quad (5)$$

$$\boldsymbol{X} \in \Gamma^b_N, \quad \boldsymbol{x} = \boldsymbol{X} + \boldsymbol{u}(\boldsymbol{X}, t), \quad \tau^f_{ij} = -p\, \delta_{ij} + \tau^V_{ij}.$$

## 2 Discrete Problem

### 2.1 Discretization of the Flow Problem

The problem will be discretized by the space-time discontinuous Galerkin method (ST-DGM). We construct a polygonal approximation $\Omega_{ht}$ of the domain $\Omega_t$. By $\mathcal{T}_{ht}$ we denote a partition of the closure $\overline{\Omega}_{ht}$ of the domain $\Omega_t$ into a finite number of closed triangles $K$ with mutually disjoint interiors such that $\overline{\Omega}_{ht} = \bigcup_{K \in \mathcal{T}_{ht}} K$.

By $\mathcal{F}_h, \mathcal{F}^B_h, \mathcal{F}^I_h$ we denote the systems of all faces of all elements $K \in \mathcal{T}_{ht}$, boundary faces and inner faces, respectively. Further, we introduce the set of "Dirichlet" boundary faces $\mathcal{F}^D_h = \{\Gamma \in \mathcal{F}^B_h; \text{ a Dirichlet condition is prescribed on } \Gamma\}$. Each face $\Gamma$ is associated with a unit normal $\boldsymbol{n}_\Gamma$, which has the same orientation as the outer normal on $\Gamma \in \mathcal{F}^B_h$. We set $h_\Gamma = \text{length of } \Gamma \in \mathcal{F}_h$.

We introduce the space of piecewise polynomial functions $\boldsymbol{S}^r_{ht} = \{v; v|_K \in P_r(K) \,\forall\, K \in \mathcal{T}_{ht}\}^4$, where $r > 0$ is an integer and $P_r(K)$ denotes the space of all polynomials on $K$ of degree $\leq r$. A function $\varphi \in \boldsymbol{S}^r_{ht}$ is, in general, discontinuous on interfaces $\Gamma \in \mathcal{F}^I_h$. By $\varphi^{(L)}_\Gamma$ and $\varphi^{(R)}_\Gamma$ we denote the values of $\varphi \in \boldsymbol{S}^r_{ht}$ on $\Gamma$ from the side of the element $K^{(L)}_\Gamma$ and $K^{(R)}_\Gamma$ adjacent to $\Gamma$ lying in the opposite direction to $\boldsymbol{n}_\Gamma$ and in the direction of $\boldsymbol{n}_\Gamma$, respectively. Then we set $\langle \varphi \rangle_\Gamma = (\varphi^{(R)}_\Gamma + \varphi^{(L)}_\Gamma)/2$ and $[\varphi]_\Gamma = \varphi^{(L)}_\Gamma - \varphi^{(R)}_\Gamma$.

The discrete problem is derived in the following way: We multiply system (1) by a test function $\varphi_h \in S_{ht}^r$, integrate over $K \in \mathcal{T}_{ht}$, apply Green's theorem, sum over all elements $K \in \mathcal{T}_{ht}$, use the concept of the numerical flux and introduce suitable terms mutually vanishing for a regular exact solution and linearize the resulting forms (see, e.g. [1, 3]). In this way we get the following forms:

$$\hat{a}_h(\overline{w}_h, w_h, \varphi_h, t) = \sum_{K \in \mathcal{T}_{ht}} \int_K \sum_{s=1}^{2} \sum_{k=1}^{2} \mathbb{K}_{s,k}(\overline{w}_h) \frac{\partial w_h}{\partial x_k} \cdot \frac{\partial \varphi_h}{\partial x_s} \, d\boldsymbol{x} \qquad (6)$$

$$- \sum_{\Gamma \in \mathcal{F}_{ht}^I} \int_\Gamma \sum_{s=1}^{2} \left\langle \sum_{k=1}^{2} \mathbb{K}_{s,k}(\overline{w}_h) \frac{\partial w_h}{\partial x_k} \right\rangle (\boldsymbol{n}_\Gamma)_s \cdot [\varphi_h] \, dS$$

$$- \sum_{\Gamma \in \mathcal{F}_{ht}^D} \int_\Gamma \sum_{s=1}^{2} \sum_{k=1}^{2} \mathbb{K}_{s,k}(\overline{w}_h) \frac{\partial w_h}{\partial x_k} (\boldsymbol{n}_\Gamma)_s \cdot \varphi_h \, dS$$

$$- \Theta \sum_{\Gamma \in \mathcal{F}_{ht}^I} \int_\Gamma \sum_{s=1}^{2} \left\langle \sum_{k=1}^{2} \mathbb{K}_{k,s}^T(\overline{w}_h) \frac{\partial \varphi_h}{\partial x_k} \right\rangle (\boldsymbol{n}_\Gamma)_s \cdot [w_h] \, dS$$

$$- \Theta \sum_{\Gamma \in \mathcal{F}_{ht}^D} \int_\Gamma \sum_{s=1}^{2} \sum_{k=1}^{2} \mathbb{K}_{k,s}^T(\overline{w}_h) \frac{\partial \varphi_h}{\partial x_k} (\boldsymbol{n}_\Gamma)_s \cdot w_h \, dS,$$

$$d_h(w_h, \varphi_h, t) = \sum_{K \in \mathcal{T}_{ht}} \int_K (w_h \cdot \varphi_h) \operatorname{div} z \, dx, \qquad (7)$$

$$J_h(w_h, \varphi_h, t) = \sum_{\Gamma \in \mathcal{F}_{ht}^I} \int_\Gamma \frac{\mu C_W}{h_\Gamma} [w_h] \cdot [\varphi_h] \, dS + \sum_{\Gamma \in \mathcal{F}_{ht}^D} \int_\Gamma \frac{\mu C_W}{h_\Gamma} w_h \cdot \varphi_h \, dS,$$
$$(8)$$

$$\ell_h(w_h, \varphi_h, t) = \sum_{\Gamma \in \mathcal{F}_{ht}^D} \int_\Gamma \frac{\mu C_W}{h_\Gamma} w_B \cdot \varphi_h \, dS \qquad (9)$$

$$- \Theta \sum_{\Gamma \in \mathcal{F}_{ht}^D} \int_\Gamma \sum_{k=1}^{2} \mathbb{K}_{k,s}^T(\overline{w}_h) \frac{\partial \varphi_h}{\partial x_k} (\boldsymbol{n}_\Gamma)_s \cdot w_B \, dS,$$

$$\hat{b}_h(\overline{w}_h, w_h, \varphi_h, t) = \qquad (10)$$

$$- \sum_{K \in \mathcal{T}_{ht_{k+1}}} \int_K \sum_{s=1}^{2} ((\mathbb{A}_s(\overline{\boldsymbol{w}}_h(x)) - z_s(x)\mathbb{I})\boldsymbol{w}_h(x)) \cdot \frac{\partial \varphi_h(x)}{\partial x_s} dx$$

$$+ \sum_{\Gamma \in \mathcal{F}_{ht}^I} \int_\Gamma \left( \mathbb{P}_g^+ (\langle \overline{\boldsymbol{w}}_h \rangle_\Gamma, \boldsymbol{n}_\Gamma) \boldsymbol{w}_h^{(L)} + \mathbb{P}_g^- (\langle \overline{\boldsymbol{w}}_h \rangle_\Gamma, \boldsymbol{n}_\Gamma) \boldsymbol{w}_h^{(R)} \right) \cdot [\varphi_h] \, dS$$

$$+ \sum_{\Gamma \in \mathcal{F}_{ht}^B} \int_\Gamma \left( \mathbb{P}_g^+ (\langle \overline{\boldsymbol{w}}_h \rangle_\Gamma, \boldsymbol{n}_\Gamma) \boldsymbol{w}_h^{(L)} + \mathbb{P}_g^- (\langle \overline{\boldsymbol{w}}_h \rangle_\Gamma, \boldsymbol{n}_\Gamma) \overline{\boldsymbol{w}}_h^{(R)} \right) \cdot \varphi_h \, dS,$$

$C_W > 0$ is a sufficiently large constant. We set $\Theta = 1$ or $\Theta = 0$ or $\Theta = -1$ and get the so-called symmetric version (SIPG) or incomplete version (IIPG) or nonsymmetric version (NIPG), respectively, of the discretization of viscous terms. The symbols $\mathbb{P}_g^+(\boldsymbol{w}, \boldsymbol{n})$ and $\mathbb{P}_g^-(\boldsymbol{w}, \boldsymbol{n})$ denote the "positive" and "negative" parts of the matrix $\mathbb{P}_g(\boldsymbol{w}, \boldsymbol{n}) = \sum_{s=1}^{2}(\mathbb{A}_s(\boldsymbol{w}) - z_s\mathbb{I})n_s$ defined, e.g., in [2]. The boundary state $\boldsymbol{w}_B$ is defined on the basis of the prescribed Dirichlet boundary conditions and extrapolation.

For the space-time discretization we consider a partition $0 = t_0 < t_1 < \ldots < t_M = T$ of the time interval $[0, T]$ and denote $I_m = (t_{m-1}, t_m)$, $\tau_m = t_m - t_{m-1}$, for $m = 1, \ldots, M$. We define the space $\mathbf{S}_{h\tau}^{rq} = \{\phi \,;\, \phi|_{I_m} = \sum_{i=0}^{q} \zeta_i \phi_i, \text{ where } \phi_i \in S_{ht}^r, \zeta_i \in P^q(I_m)\}^2$ with integers $r, q \geq 1$. $P^q(I_m)$ denotes the space of all polynomials in $t$ on $I_m$ of degree $\leq q$. For $\varphi \in \mathbf{S}_{h\tau}^{rq}$ we set $\varphi_m^\pm = \varphi(t_m^\pm) = \lim_{t \to t_{m\pm}} \varphi(t)$, $\{\varphi\}_m = \varphi_m^+ - \varphi_m^-$. The initial state $\boldsymbol{w}_{h\tau}(0-) \in \mathbf{S}_{h0}^p$ is defined as the $L^2(\Omega_{h0})$-projection of $\boldsymbol{w}^0$ on $\mathbf{S}_{h0}^r$. Moreover, we introduce the prolongation $\overline{\boldsymbol{w}}_{h\tau}(t)$ of $\boldsymbol{w}_{h\tau}|_{I_{m-1}}$ on the interval $I_m$. By $(\cdot, \cdot)_t$ we denote the $L^2(\Omega_{ht})$-scalar product.

Now the *space-time DG approximate solution* is defined as a function $\boldsymbol{w}_{h\tau} \in \mathbf{S}_{h\tau}^{rq}$ satisfying the following relation for $m = 1, \ldots, M$:

$$\int_{I_m} \left( \left( \frac{D^\mathcal{A} \boldsymbol{w}_{h\tau}}{Dt}(t), \varphi_{h\tau} \right)_t + \hat{a}_h(\overline{\boldsymbol{w}}_{h\tau}, \boldsymbol{w}_{h\tau}, \varphi_{h\tau}, t) \right) \, dt \qquad (11)$$

$$+ \int_{I_m} \left( \hat{b}_h(\overline{\boldsymbol{w}}_{h\tau}, \boldsymbol{w}_{h\tau}, \varphi_{h\tau}, t) + \int_{I_m} J_h(\boldsymbol{w}_{h\tau}, \varphi_{h\tau}, t) \right) \, dt$$

$$+ (\{\boldsymbol{w}_{h\tau}\}_{m-1}, \varphi_{h\tau}(t_{m-1}+)) = \int_{I_m} \ell_h(\boldsymbol{w}_{hD}, \varphi_{h\tau}, t) \, dt, \quad \forall \varphi_{h\tau} \in \mathbf{S}_{h\tau}^{rq}.$$

## 2.2 Discretization of the Elasticity Problem

The elasticity problem will also be discretized by the ST-DGM. To this end, the problem is reformulated as a couple of equations of the first order in time: find functions $\boldsymbol{u}$ and $z : \Omega^b \times [0, T] \to \mathbb{R}^2$ such that

$$\rho^b \frac{\partial z}{\partial t} + c\rho^b z - div\, \boldsymbol{\sigma}(\boldsymbol{u}) = \boldsymbol{f} \quad \text{in } \Omega^b \times (0, T), \tag{12}$$

$$\frac{\partial \boldsymbol{u}}{\partial t} - z = 0 \quad \text{in } \Omega^b \times (0, T), \tag{13}$$

$$\boldsymbol{u} = \boldsymbol{u}_D \quad \text{in } \Gamma_D^b \times (0, T), \quad \boldsymbol{\sigma}(\boldsymbol{u}) \cdot \boldsymbol{n} = \boldsymbol{g}_N \quad \text{in } \Gamma_N^b \times (0, T), \tag{14}$$

$$\boldsymbol{u}(x, 0) = \boldsymbol{u}_0(x), \quad z(x, 0) = z_0(x), \quad x \in \Omega^b. \tag{15}$$

Now we proceed in a similar way as in Sect. 2.1. By $\Omega_h^b$ we denote a polygonal approximation of the domain $\Omega^b$. The sets $\Gamma_{Dh}^b, \Gamma_{Nh}^b \subset \partial \Omega_h^b$ will approximate $\Gamma_D^b$ and $\Gamma_N^b$. Let $\mathcal{T}_h^b$ be a partition of the closure $\overline{\Omega}_h^b$. We define the finite dimensional space $\boldsymbol{S}_{hs}^b = \{v \in L^2(\Omega_h^b); v|_K \in P_s(K), K \in \mathcal{T}_h^b\}^2$, where $s > 0$ is an integer. By $\mathcal{F}_h^b, \mathcal{F}_h^{bD}, \mathcal{F}_h^{bN}, \mathcal{F}_h^{bI}$ we denote the system of all faces of all elements $K \in \mathcal{T}_h^b$, boundary Dirichlet, Neumann faces and inner faces. If we introduce the forms

$$a_h^b(\boldsymbol{u}, \boldsymbol{v}) = \sum_{K \in \mathcal{T}_h^b} \int_K \boldsymbol{\sigma}(\boldsymbol{u}) : \boldsymbol{e}(\boldsymbol{v}) \, dx - \sum_{\Gamma \in \mathcal{F}_h^{bI}} \int_\Gamma (\langle \boldsymbol{\sigma}(\boldsymbol{u}) \rangle \cdot \boldsymbol{n}) \cdot [\boldsymbol{v}] \, dS \tag{16}$$

$$- \sum_{\Gamma \in \mathcal{F}_h^{bD}} \int_\Gamma (\boldsymbol{\sigma}(\boldsymbol{u}) \cdot \boldsymbol{n}) \cdot \boldsymbol{v} \, dS - \Theta \sum_{\Gamma \in \mathcal{F}_h^{bI}} \int_\Gamma (\langle \boldsymbol{\sigma}(\boldsymbol{v}) \rangle \cdot \boldsymbol{n}) \cdot [\boldsymbol{u}] \, dS$$

$$- \Theta \sum_{\Gamma \in \mathcal{F}_h^{bD}} \int_\Gamma (\boldsymbol{\sigma}(\boldsymbol{v}) \cdot \boldsymbol{n}) \cdot \boldsymbol{u} \, dS,$$

$$J_h^b(\boldsymbol{u}, \boldsymbol{v}) = \sum_{\Gamma \in \mathcal{F}_h^{bI}} \int_\Gamma \frac{C_W^b}{h_\Gamma} [\boldsymbol{u}] \cdot [\boldsymbol{v}] \, dS + \sum_{\Gamma \in \mathcal{F}_h^{bD}} \int_\Gamma \frac{C_W^b}{h_\Gamma} \boldsymbol{u} \cdot \boldsymbol{v} \, dS, \tag{17}$$

$$\ell_h^b(\boldsymbol{v})(t) = \sum_{K \in \mathcal{T}_h^b} \int_K \boldsymbol{f}(t) \cdot \boldsymbol{v} \, dx + \sum_{\Gamma \in \mathcal{F}_h^{bN}} \int_\Gamma \boldsymbol{g}_N(t) \cdot \boldsymbol{v} \, dS \tag{18}$$

$$- \Theta \sum_{\Gamma \in \mathcal{F}_h^{bD}} \int_\Gamma (\boldsymbol{\sigma}(\boldsymbol{v}) \cdot \boldsymbol{n}) \cdot \boldsymbol{u}_D(t) \, dS + \sum_{\Gamma \in \mathcal{F}_h^{bD}} \int_\Gamma \frac{C_W^b}{h_\Gamma} \boldsymbol{u}_D(t) \cdot \boldsymbol{v} \, dS,$$

$$(\boldsymbol{u}, \boldsymbol{v})_{\Omega_h^b} = \int_{\Omega_h^b} \boldsymbol{u} \cdot \boldsymbol{v} \, dx = \sum_{K \in \mathcal{T}_h^b} \int_K \boldsymbol{u} \cdot \boldsymbol{v} \, dx, \tag{19}$$

where $C_W^b > 0$ is a sufficiently large constant, $\Theta = 1$, $\Theta = 0$ or $\Theta = -1$ and $\boldsymbol{S}_{h\tau}^{b,sq} = \{v \in L^2(\Omega_h^b \times (0, T); v|_{I_m} = \sum_{i=0}^q t^i \varphi_i \text{ with } \varphi_i \in S_{hs}^b, m = 1, \dots, M\}^2$, the ST-DG approximate solution can be defined as a couple $\boldsymbol{u}_{h\tau}, z_{h\tau} \in \boldsymbol{S}_{h\tau}^{b,sq}$ such that

(a) $\displaystyle\int_{I_m}\left(\rho^b\left(\frac{\partial z_{h\tau}}{\partial t}, \boldsymbol{v}_{h\tau}\right)_{\Omega_h^b} + C\left(\rho^b z_{h\tau}, \boldsymbol{v}_{h\tau}\right)_{\Omega_h^b} + a_h^b(\boldsymbol{u}_{h\tau}, \boldsymbol{v}_{h\tau})\right.$ (20)

$\displaystyle\left. + J_h^b(\boldsymbol{u}_{h\tau}, \boldsymbol{v}_{h\tau})\right) dt + (\{\boldsymbol{u}_{h\tau}\}_{m-1}, \boldsymbol{v}_{h\tau}(t_{m-1}+))_{\Omega_h^b}$

$\displaystyle= \int_{I_m} \ell(\boldsymbol{v}_{h\tau})\, dt \quad \forall \boldsymbol{v}_{h\tau} \in \boldsymbol{S}_{h\tau}^{b,sq},$

(b) $\displaystyle\int_{I_m}\left(\left(\frac{\partial \boldsymbol{u}_{h\tau}}{\partial t}, \boldsymbol{w}_{h\tau}\right)_{\Omega_h^b} - (z_{h\tau}, \boldsymbol{w}_{h\tau})_{\Omega_h^b}\right) dt$

$\displaystyle+ (\{\boldsymbol{u}_{h\tau}\}_{m-1}, \boldsymbol{w}_{h\tau}(t_{m-1}+))_{\Omega_h^b} = 0 \quad \forall \boldsymbol{w}_{h\tau} \in \boldsymbol{S}_{h\tau}^{b,sq},$

$m = 1, \ldots, M.$

The initial states $\boldsymbol{u}_h(0-), z_h(0-) \in \boldsymbol{S}_{hs}^b$ are defined by $(\boldsymbol{u}_h(0-), \boldsymbol{v}_h)_{\Omega_h^b} = (\boldsymbol{u}^0, \boldsymbol{v}_h)_{\Omega_h^b}$, $(z_h(0-), \boldsymbol{v}_h)_{\Omega_h^b} = (z^0, \boldsymbol{v}_h)_{\Omega_h^b}$ for all $\boldsymbol{v}_h \in \boldsymbol{S}_{hs}^b$.

In the FSI problem the coupling of the discrete flow problem (11) and structural problem (20) are realized via the discrete version of transmission conditions (5). The coupled problem is solved with the aid of the following coupling procedure.

1. Assume that the approximate solution of the flow problem on the time level $t_k$ is known as well as the deformation of the structure $\boldsymbol{u}_{h,k}$.
2. Set $\boldsymbol{u}_{h,k+1}^0 := \boldsymbol{u}_{h,k}$, $l := 1$ and apply the iterative process:

   a. Compute the stress tensor $\tau_{ij}^f$ and the aerodynamical force acting on the structure and transform it to the interface $\Gamma_{Nh}^b$.
   b. Solve the elasticity problem, compute the deformation $\boldsymbol{u}_{h,k+1}^l$ at time $t_{k+1}$ and approximate the domain $\Omega_{ht_{k+1}}^l$.
   c. Determine the ALE mapping $\mathcal{A}_{t_{k+1}h}^l$ and approximate the domain velocity $z_{h,k+1}^l$.
   d. Solve the flow problem on the approximation of $\Omega_{ht_{k+1}}^l$.
   e. If the variation of the displacement $\boldsymbol{u}_{h,k+1}^l$ and $\boldsymbol{u}_{h,k+1}^{l-1}$ is larger than the prescribed tolerance, go to (a) and $l := l + 1$. Else $k := k + 1$ and goto (2).

This represents the so-called strong coupling. If in the step (e) we set $k := k + 1$ and go to (2) already in the case when $l = 1$, then we get the weak (loose) coupling.

## 3 Numerical Results

We consider a 2D model of gas flow past an elastic airfoil. For testing our method we assume that the material of the airfoil is very soft. It is characterized by the Lamè parametres $\lambda^b = 2 \cdot 10^7$ Pa and $\mu^b = 5 \cdot 10^6$ Pa. The structural damping coefficients

**Fig. 1** Triangulation at time $t = 0$ used for the computation of fluid flow and triangulation for the elasticity problem



**Fig. 2** Visualization of velocity vectors and of the deformed elastic airfoil at time $t = 0.15$ s

are chosen as $c_M = 0.1\,\mathrm{s}^{-1}$ and $c_K = 0.1\,\mathrm{s}$ and the material density is given by $\rho^b = 10^4\,\mathrm{kg}\,\mathrm{m}^{-3}$.

The fluid flow simulation was carried out using the following data: $\mu = 1.72 \cdot 10^{-5}\,\mathrm{kg}\,\mathrm{m}^{-1}.\mathrm{s}$, far-field pressure $p = 101250\,\mathrm{Pa}$, far-field density $\rho = 1.225\,\mathrm{kg}\,\mathrm{m}^{-3}$, Poisson adiabatic constant $\gamma = 1.4$, specific heat $c_v = 721.428\,\mathrm{m}^2\,\mathrm{s}^{-2}\,\mathrm{K}^{-1}$, heat conduction coefficient $k = 2.428 \cdot 10^{-2}\,\mathrm{kg}\,\mathrm{m}\,.\,\mathrm{s}^{-2}\,\mathrm{K}^{-1}$. The far-field velocity was $40\,\mathrm{m}\,\mathrm{s}^{-1}$. Figure 1 shows the triangulation at the initial time $t = 0$.

Fluid flow is solved by the ST-DGM with quadratic polynomials in space and linear polynomials in time. For the elasticity problem we also used the ST-DGM, but with linear polynomials in space and constant polynomials in time. For both problems the non-symmetric version (NIPG) was used. For flow problem we set $C_W = 1000$ on the interior elements and $C_W = 10000$ on the boundary elements in order to keep the prescribed Dirichlet boundary conditions, particularly in the boundary layer. For elasticity we set $C_W^b = 10^{10}$ in order to match the magnitude of the Lamè parametres. We used the time step $\tau = 2.25 \cdot 10^{-6}$ s. The strong coupling was used for the FSI

process. The accuracy $10^{-6}$ was achieved with at most 5 iteration on each time level. Figure 2 shows the visualization of the deformed airfoil and the velocity vectors.

# References

1. Česenek, J., Feistauer, M., Kosík, A.: DGFEM for the analysis of airfoil vibrations induced by compressible flow. Z. Angew. Math. Mech. **93**(6–7), 387–402 (2013)
2. Feistauer, M., Felcman, J., Straškraba, I.: Mathematical and Computational Methods for Compressible Flow. Clarendon Press, Oxford (2003)
3. Feistauer, M., Horáček, J., Kučera, V., Prokopová, J.: On numerical solution of compressible flow in time-dependent domains. Mathematica Bohemica **137**, 1–16 (2011)
4. Nomura, T., Hughes, T.J.R.: An arbitrary Lagrangian-Eulerian finite element method for interaction of flow and a rigid body. Comput. Meth. Appl. Mech. Eng. **95**, 115–138 (1992)

# An Anisotropic Diffusion Finite Volume Algorithm Using a Small Stencil

**Martin Ferrand, Jacques Fontaine and Ophélie Angelini**

**Abstract**   This article presents a finite volume algorithm to solve anisotropic heterogeneous diffusion equations within the open source CFD software *Code_Saturne*. This algorithm has the advantage to use a small stencil composed of face neighbouring cells only, which makes it easy to parallelize. The resolution is performed through an iterative process (fixed point Picard algorithm). Second order convergence in space is numerically obtained on various analytical test-cases and mesh sequences of the FVCA6 benchmark and the results are compared to the barycentric version of the SUSHI scheme [3].

## 1 Introduction

Several discretization schemes for anisotropic, heterogeneous diffusion problems are presented in the literature, especially with the finite volume method (see [2], or the 2D [7] and 3D [5] benchmarks). In this proceeding, we propose a two-point flux approximation (TPFA) scheme for non-cartesian grids within an industrial code called *Code_Saturne*  (see [6] for more information), applied to various steps in the core solver (diffusion part of scalar transport equation with Generalized Gradient Diffusion Hypothesis (GGDH) [1], projection step of the predictor-corrector Navier-Stokes solver in presence of head-losses, etc.). Therefore, the proposed scheme is aimed to be highly parallelized, with a ghost-cells technique, and should use the smallest stencil possible to be run on large meshes.

In this paper, the space discretization is presented in Sect. 2, then the scheme is derived from the flux continuity property, the resolution is performed through an iterative process. The obtained results are compared to the barycentric version of the SUSHI scheme [3] and provided in Sect. 3.

M. Ferrand (✉) · J. Fontaine · O. Angelini
Fluid Mechanics, Energy and Environment, EDF R&D,
6, quai Watier, 78400 Chatou, France
e-mail: martin.ferrand@edf.fr

## 2 Space Discretization of the Anisotropic Heterogeneous Diffusion Equation

In this article, first-order tensors (identified to vector fields) are underlined once, and second-order tensors (identified to matrix fields) are twice underlined.

The studied Poisson equation reads:

$$\begin{cases} -\text{div}\left(\underline{\underline{K}} \cdot \underline{\nabla} Y\right) = f & \text{on } \Omega \\ \qquad\qquad\qquad Y = Y_d & \text{on } \mathscr{D} \subset \partial\Omega \\ \left(\underline{\underline{K}} \cdot \underline{\nabla} Y\right) \cdot \underline{n} = Q_d & \text{on } \mathscr{N} = \partial\Omega \setminus \mathscr{D} \end{cases} \tag{1}$$

where we denote by $\partial\Omega = \overline{\Omega} \setminus \Omega$ the boundary of the domain $\Omega$, an open bounded connected polyhedral subset of $\mathbb{R}^d$ ($d$ is the space dimension), $Y$ is the scalar field defined on the domain $\Omega$ and belongs to $H^1(\Omega)$, $\underline{\underline{K}}$ is the tensor diffusivity field, assumed to be symmetric positive definite on the whole domain, continuous by part and limited, and $f \in L^2(\Omega)$ a source term. The boundary conditions on the field $Y$ are composed of Dirichlet conditions on $\mathscr{D}$, and of Neumann conditions on $\mathscr{N}$. $\mathscr{D}$ and $\mathscr{N}$ are partitions of the boundary $\partial\Omega$.

### 2.1 Space Discretization

The domain is discretized into cells $\Omega_i$ on which $K$ is supposed to be constant. The barycentre of a cell $i$ (respectively of a cell $j$) is denoted by $I$ (respectively by $J$). Two cells $i$ and $j$ are said to be neighbours if they share a face noted $f_{ij}$ of centre $F$. The intersection between $f_{ij}$ and the vector $\underline{IJ}$ is $O$. The unit normal vector to the face $f_{ij}$ oriented from $i$ to $j$ is $\underline{n}_{ij}$, the surface of $f_{ij}$ is $S_{f_{ij}}$, and we define $\underline{S}_{ij} = S_{f_{ij}}\underline{n}_{ij}$. The set of all interior faces of $i$ is denoted $\mathscr{F}_i^{int}$. Finally, $I'$ (respectively $J'$) is the orthogonal projection of $I$ (respectively of $J$) with respect to the face $f_{ij}$.

Faces shared by one and only one cell $i$ are said to be boundary faces and denoted $f_b$. $\underline{n}_{ib}$ is the outward normal to the face $f_b$, and $\underline{S}_{ib}$ is the normal vector to the face which norm is the surface. The set of boundary faces of $i$ is denoted by $\mathscr{F}_i^{ext}$, whereas $\mathscr{F}_i$ is the set of all faces of $i$. Finally, $I'$ is the orthogonal projection of $I$ with respect to face normal of $f_b$. All the geometric definitions are recalled on Fig. 1a and on Fig. 1b. The cell mean of $Y$ in $i$ is defined by $Y_i \equiv \dfrac{1}{|\Omega_i|} \displaystyle\int_{\Omega_i} Y \, d\Omega$.

*Code_Saturne* uses a finite volume scheme where the solved variables are stored at cell centres: i.e. $Y_I = Y_i$. The discretized field $Y$ is assumed to be affine on each cell. Therefore, for every point $I''$ of a cell $\Omega_i$, the following relationship reads:

$$Y_{I''} = Y_i + \underline{\nabla}_i Y \cdot \underline{II''} \tag{2}$$

**Fig. 1** Sketch displaying interior and exterior faces with reconstruction points $I''$ and $J''$, **a** interior face, **b** boundary face

where $\underline{\nabla}_i Y$ is the $Y$ gradient, constant within each cell. Face quantities are defined by: $Y_f = \dfrac{1}{S_f} \displaystyle\int_{S_f} Y \, dS$. Thanks to the affinity of $Y$, we have $Y_f = Y_F$. Integrating Equation (1) over a cell $\Omega_i$ gives:

$$- \sum_{f \in \mathscr{F}_i} \left( \underline{\underline{K}} \cdot \underline{\nabla} Y \right)_f \cdot \underline{S}_f = |\Omega_i| \, f_i \tag{3}$$

## 2.2 Two-Point Flux Approximation Scheme

The aim of this section is to write a two-point flux as for isotropic diffusion problem on orthogonal meshes (see [4] for a description of TPFA schemes). One can notice that for the face $f_{ij}$, the diffusive flux seen by cell $i$ can be written as:

$$\boxed{\left( \underline{\underline{K}}_i \cdot \underline{\nabla} Y \right) \cdot \underline{n}_{ij} = \left( \underline{\nabla} Y \right) \cdot \left( \underline{\underline{K}}_i^T \cdot \underline{n}_{ij} \right)} \tag{4}$$

The tensor $\underline{\underline{K}}$ is symmetric, therefore $\underline{\underline{K}}^T = \underline{\underline{K}}$. The Eq. (4) indicates that the direction $\underline{\underline{K}}_i \cdot \underline{n}_{ij}$ is optimal to discretize the flux. Let $I''$ (resp. $J''$) be a point on the line passing through $F$ with direction vector $\underline{\underline{K}}_i \cdot \underline{n}_{ij}$ (resp. $\underline{\underline{K}}_j \cdot \underline{n}_{ij}$). The diffusive flux seen from cell $i$ is then approximated by:

$$\left( \underline{\underline{K}} \cdot \underline{\nabla} Y \right)_F \cdot \underline{n}_{ij} = \left( \underline{\nabla} Y \right)_F \cdot \left( \underline{\underline{K}}_i \cdot \underline{n}_{ij} \right) \simeq \frac{\left\| \underline{\underline{K}}_i \cdot \underline{n}_{ij} \right\|}{I''F} \left( Y_F - Y_{I''} \right) \tag{5}$$

where $\overline{I''F}$ is the algebraic distance between $I''$ and $F$. The flux seen from cell $j$ is approximated by:

$$\left(\underline{\underline{K}} \cdot \underline{\nabla} Y\right)_F \cdot \underline{n}_{ij} = \left(\underline{\nabla} Y\right)_F \cdot \left(\underline{\underline{K}} \cdot \underline{n}_{ij}\right) \simeq \frac{\left\|\underline{\underline{K}}_j \cdot \underline{n}_{ij}\right\|}{\overline{FJ''}} (Y_{J''} - Y_F) \qquad (6)$$

where $\overline{FJ''}$ is the algebraic distance between $F$ and $J''$. Enforcing the continuity of the fluxes (5) and (6) yields:

$$Y_F = \left(\frac{\overline{I''F}}{\left\|\underline{\underline{K}}_i \cdot \underline{n}_{ij}\right\|} + \frac{\overline{FJ''}}{\left\|\underline{\underline{K}}_j \cdot \underline{n}_{ij}\right\|}\right)^{-1} \left(Y_{I''} \frac{\overline{FJ''}}{\left\|\underline{\underline{K}}_j \cdot \underline{n}_{ij}\right\|} + Y_{J''} \frac{\overline{I''F}}{\left\|\underline{\underline{K}}_i \cdot \underline{n}_{ij}\right\|}\right) \qquad (7)$$

Therefore $Y_F$ is the weighted harmonic mean depending on $\underline{\underline{K}} \cdot \underline{n}$ in cells $i$ and $j$. Injecting (7) in (5) gives a two-point formula for the flux through face $f_{ij}$:

$$- \left(\underline{\underline{K}} \cdot \underline{\nabla} Y\right)_{f_{ij}} \cdot \underline{n}_{ij} = \frac{K_{f_{ij}}}{\overline{I''J''}} (Y_{I''} - Y_{J''})$$

where the equivalent scalar face diffusivity $K_{f_{ij}}$ is defined by:

$$\frac{K_{f_{ij}}}{\overline{I''J''}} = \left(\frac{\overline{I''F}}{\left\|\underline{\underline{K}}_i \cdot \underline{n}_{ij}\right\|} + \frac{\overline{FJ''}}{\left\|\underline{\underline{K}}_j \cdot \underline{n}_{ij}\right\|}\right)^{-1} \qquad (8)$$

Besides, boundary terms are discretized as:

$$- \left(\underline{\underline{K}} \cdot \underline{\nabla} Y\right)_{f_b} \cdot \underline{n}_{ib} = A_Y^f + B_Y^f Y_{I'}$$

## 2.3 Non-Orthogonalities Reconstruction and Iterative Solving

As the field values of $Y$ at $I''$ and $J''$ are not degrees of freedom of the discretized field, they are written in terms of $Y_i$ and $Y_j$ using the cell gradient through (2):

$$\begin{aligned} Y_{I''} &= Y_I + \underline{\nabla}_i Y \cdot \underline{II''} \\ Y_{J''} &= Y_J + \underline{\nabla}_j Y \cdot \underline{JJ''} \end{aligned} \qquad (9)$$

One should note that the choices for $I''$ and $J''$ are still arbitrary, the only requirement is that $\underline{FI''}$ should be collinear to $\underline{\underline{K}}_i \cdot \underline{n}_{ij}$ and $\underline{FJ''}$ should be collinear to $\underline{\underline{K}}_j \cdot \underline{n}_{ij}$. Eventually, inserting (9) in Equation (3) reads:

$$
\begin{array}{|l}
\displaystyle\sum_{f_{ij}\in\mathscr{F}_i^{int}} \frac{K_{f_{ij}}}{I''J''}\left(Y_i - Y_j + \left(\underline{\nabla}_i Y \cdot \underline{II''} - \underline{\nabla}_j Y \cdot \underline{JJ''}\right)\right) \\[3mm]
+ \displaystyle\sum_{f_b\in\mathscr{F}_i^{ext}}\left(A_Y^f + B_Y^f\left(Y_i + \underline{\nabla}_i Y \cdot \underline{II''}\right)\right) - |\Omega_i|\, f_i \qquad = E_i\left(\mathbf{Y}\right) = 0
\end{array}
\tag{10}
$$

where $\mathbf{E}\left(.\right)$ is the linear operator to be solved and $\mathbf{Y}$ be the vector containing the degree of freedom (vector of size $N_{cel}$) of the discretized field $Y$: $\mathbf{Y} = \left[Y_1, \cdots, Y_i, \cdots, Y_{N_{cel}}\right]$.

Equation (10) is solved using a Picard fix point. This approach is widely used in industrial CFD codes because it gives a linear system with good invertibility properties and with low memory consumption due to the small stencil. The problem (10) can be rewritten as: find $\mathbf{Y}$ so that $\mathbf{E}\left(\mathbf{Y}\right) = \mathbf{0}$. Let $\left(\mathbf{Y}^k\right)_{k\in\mathbb{N}}$ be a series initialized at $\mathbf{0}$ defined by:

$$
\begin{cases}
\mathbf{EM}\left(\delta\mathbf{Y}^{k+1}\right) = \mathbf{E}\left(\mathbf{Y}^k\right) \\
\mathbf{Y}^{k+1} = \mathbf{Y}^k + \delta\mathbf{Y}^{k+1}
\end{cases}
\tag{11}
$$

where $\mathbf{EM}\left(\cdot\right)$ is the linear operator built with non-reconstructed fluxes:

$$
EM_i\left(\delta\mathbf{Y}^{k+1}\right) = \sum_{f_{ij}\in\mathscr{F}_i^{int}} \frac{K_{f_{ij}}}{I''J''}\left(\delta Y_i - \delta Y_j\right) + \sum_{f_b\in\mathscr{F}_i^{ext}} B_Y^f \delta Y_i
\tag{12}
$$

The incremental vector is computed at each sub-iteration solving a linear system of size $N_{cel} \times N_{cel}$. The sparse matrix $\mathbf{EM}$ is composed of a diagonal of size $N_{cel}$ and an extra-diagonal of size $N_{fac}$ (number of interior faces), and is a $M-$matrix, thus, is invertible and its inverse is positive definite. In all the numerical test cases, a conjugate gradient combined with an algebraic multi-grid algorithm is used.

If $I''$ and $J''$ were chosen as orthogonal projections on the lines passing through $F$ with respective direction vectors $\underline{\underline{K}}_i \cdot n_{ij}$ and $\underline{\underline{K}}_j \cdot n_{ij}$, it would not ensure the positivity of the face viscosity $K_{f_{ij}}$ defined by (8) and the resulting matrix would loose its invertibility property. Hence, we choose to maintain $I''$ and $J''$ on the same side of the

face $f_{ij}$ as $I$ and $J$, by setting: $\dfrac{\overline{I''F}}{\left\|\underline{\underline{K}}_i \cdot n_{ij}\right\|} = \max\left[\dfrac{\underline{IF} \cdot \underline{\underline{K}}_i \cdot n_{ij}}{\left\|\underline{\underline{K}}_i \cdot n_{ij}\right\|^2}, \varepsilon\,\dfrac{\overline{I'F}}{\left\|\underline{\underline{K}}_i \cdot n_{ij}\right\|}\right]$,

with $\varepsilon$ set to 0.1. That tends to drive away the built matrix from the one defined in (10), and thus requires more reconstruction sweeps, but keeps the matrix as a $M-$matrix. Finally, the reconstruction gradient used in Equation (10), is the standard cell-gradient displayed in the theory documentation of *Code_Saturne*[6].

Fig. 2 An overview of mesh sequences. **a** HEX mesh. **b** TET mesh. **c** BLS mesh. **d** DBLS mesh

**Table 1** Order of convergence for the two finer meshes of hexahedral (HEX), tetrahedral (TET), prism with triangle bases (BLS) and prism with general bases (DBLS) mesh sequences for cases (13), (14) and (15)

| (a) Case 1. | | | (b) Case 2. | | | (c) Case 3. | | |
|---|---|---|---|---|---|---|---|---|
| Meshes | Ratiol2 | Ratiograd | Meshes | Ratiol2 | Ratiograd | Meshes | Ratiol2 | Ratiograd |
| HEX | 1.985 | 1.265 | HEX | 1.968 | 1.447 | HEX | 2.002 | 2.001 |
| TET | 2.136 | 1.029 | TET | 2.090 | 0.927 | TET | 2.321 | 1.691 |
| BLS | 1.888 | 1.016 | BLS | 1.962 | 0.922 | BLS | 1.702 | 1.012 |
| DBLS | 1.727 | 1.009 | DBLS | 1.728 | 1.467 | DBLS | 2.042 | 1.593 |

# 3 Verification Test Cases

In this section, numerical results obtained on the FVCA6 benchmark test cases [5] are presented and compared to the barycentric version of SUSHI scheme [3]. The number of unknowns of the linear system is denoted by nu. We deliberately have not chosen exactly the mandatory meshes of the benchmark, but all the presented cases are run on hexahedral meshes (HEX), tetrahedral meshes (TET), prism meshes with triangle bases (BLS) and prism meshes with general bases (DBLS) of the benchmark mesh database (see Fig. 2).

Then the following orders of convergence are defined:

$$
\texttt{ratiol2(i)} = -3\frac{\log\left(\frac{\texttt{erl2(i)}}{\texttt{erl2(i-1)}}\right)}{\log\left(\frac{\texttt{nu(i)}}{\texttt{nu(i-1)}}\right)}, \quad \texttt{ratiograd(i)} = -3\frac{\log\left(\frac{\texttt{ergrad(i)}}{\texttt{ergrad(i-1)}}\right)}{\log\left(\frac{\texttt{nu(i)}}{\texttt{nu(i-1)}}\right)}
$$

where i is the number of the mesh (from the coarser to the finer), erl2(i) is the normalized discrete $L^2$-error on the solution of the mesh number i, ergrad(i) the normalized discrete $L^2$-error on the gradient. In all cases bellow, the value of the analytical solution is imposed as a Dirichlet. The first analytical test-case is defined by:

$$
\underline{\underline{K}} = \begin{pmatrix} 1 & 0.5 & 0 \\ 0.5 & 1 & 0.5 \\ 0 & 0.5 & 1 \end{pmatrix}
$$

$$
Y_{ana} = 1 + \sin(\pi x)\sin\left(\pi\left(y + \tfrac{1}{2}\right)\right)\sin\left(\pi\left(z + \tfrac{1}{3}\right)\right)
$$

(13)

**Fig. 3** $L^2$ error norm for $Y$ and its gradient. **a** $L^2$ Error norm for $Y$ for case 1. **b** $L^2$ Error norm for $\Delta Y$ for case 1. **c** $L^2$ Error norm for $Y$ for case 2. **d** $L^2$ Error norm for $\Delta Y$ for case 2. **e** $L^2$ Error norm for $Y$ for case 3. **f** $L^2$ Error norm for $\Delta Y$ for case 3.

The second analytical test case is defined as follows:

$$\underline{\underline{K}} = \begin{pmatrix} 1 + y^2 + z^2 & -xy & -xz \\ -xy & 1 + x^2 + z^2 & -yz \\ -xz & -yz & 1 + x^2 + y^2 \end{pmatrix} \tag{14}$$

$$Y_{ana} = x^3 y^2 z + x \sin(2\pi xy) \sin(2\pi xz) \sin(2\pi z)$$

This case is representative of the applications covered by the algorithm (GGDH on the diffusive term of a scalar transport equation) with an heterogeneous and anisotropic tensor $\underline{\underline{K}}$, but not discontinuous. The last analytical test-case presented here is:

$$\underline{\underline{K}} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1000 \end{pmatrix}, \ Y_{ana} = \sin(2\pi x) \sin(2\pi y) \sin(2\pi z) \qquad (15)$$

The Picard fixed point algorithm is considered converged when the normed residual becomes smaller than $10^{-10}$. Note that, by construction, no iteration is needed for cases with hexahedral meshes and orthotropic diffusion coefficients (e.g. case (15) and HEX meshes). On presented verification cases about 30 iterations are needed to obtain the required precision (up to 177 sweeps for (15) case with the finest TET mesh). A summary of convergence ratios is displayed in Table 1 for the two finer meshes of each mesh sequence. We can notice that the order of convergence is around 2, even greater for the two finer tetrahedral meshes. This is due to the fact that the penultimate mesh is of poorer quality than the finest one as displayed on Fig. 3. One can remark that the numerical results have a precision of the order of the barycentric version of the SUSHI scheme, but are more sensitive to the mesh quality. The worst accuracy is obtained with the tetrahedral meshes. We also must admit that the results on kershaw-type meshes are not shown, since the present scheme is sensitive to the mesh quality criterion, which are worst as the meshes are finer in the kershaw-type sequence.

## 4 Conclusion

A new algorithm which solves anisotropic heterogeneous diffusion problem is presented. It gives satisfactory results on various types of meshes with an approximate second order of convergence in space in $L^2$-norm and with a precision close to the one obtained with the SUSHI scheme. Its main advantage is the small stencil needed, which allows an highly parallelization of the algorithm, within the industrial CFD code *Code_Saturne*.

## References

1. Dehoux, F., Benhamadouche, S., Manceau, R.: Modelling turbulent heat fluxes using the elliptic blending approach for natural convection. In: Proceedings of 7th International Symposium Turbo Shear Flow Phenomena, Ottawa, Canada (2011)
2. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. Handb. Numer. Anal. **7**, 713–1018 (2000)
3. Eymard, R., Gallouët, T., Herbin, R.: Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes sushi: a scheme using stabilization and hybrid

interfaces. IMA j. Numer. Anal. **30**(4), 1009–1043 (2010)
4. Eymard, R., Gallouët, T., Herbin, R., Masson, R.: Tp or not tp, that is the question. Hal-00801648 2 **2**, 1–19 (2013)
5. Eymard, R., Henry, G., Herbin, R., Hubert, F., Klöfkorn, R., Manzini, G.: 3d benchmark on discretization schemes for anisotropic diffusion problems on general grids, pp. 95–130. Springer, berlin (2011)
6. EDF R&D: (2014). http://www.code-saturne.org
7. Svyatskiy, D.: Benchmark on anisotropic problems.nonlinear monotone finite volume method. In: Finite Volumes for Complex Applications V, pp. 935–947 (2008)

# Coupling of Fluid Flow and Solute Transport Using a Divergence-Free Reconstruction of the Crouzeix-Raviart Element

Jürgen Fuhrmann, Alexander Linke and Christian Merdon

**Abstract**  The nonconforming Crouzeix-Raviart finite element discretization for the Navier-Stokes equations allows for a divergence-free reconstruction of the discrete velocity field in the Raviart-Thomas finite element space. Integration over the faces of the control volumes of an admissible finite volume subdivision of the normal components of this reconstructed velocity field allows the coupling to the two-point flux based exponential fitting finite volume method for mass transport. The main advantage of this scheme is that it preserves positivity and maximum principles for the concentration. In comparison to previously introduced coupling schemes based on divergence-free finite element ansatzes for the fluid flow, the new method uses a significantly smaller number of degrees of freedom. The paper introduces the coupling method, demonstrates the preservation of the qualitative properties of the discrete concentration field and, based on numerical experiments, establishes the hypothesis that the coupled scheme is convergent.

## 1 Introduction

This paper concerns the coupling between a fluid with the velocity $\mathbf{u}$, pressure $p$ and viscosity $\eta$ which satisfies the steady, incompressible Navier-Stokes equation

$$(\mathbf{u} \cdot \nabla)\mathbf{u} + \nabla p - \eta \Delta \mathbf{u} = \mathbf{f}, \qquad \nabla \cdot \mathbf{u} = 0 \quad \text{and} \quad \mathbf{u} = \mathbf{u}_D \text{ along } \Gamma_D \qquad (1)$$

J. Fuhrmann (✉) · A. Linke · C. Merdon
Weierstrass Institute for Applied Analysis and Stochastics (WIAS),
Berlin, Germany
e-mail: juergen.fuhrmann@wias-berlin.de

A. Linke
e-mail: alexander.linke@wias-berlin.de

C. Merdon
e-mail: christian.merdon@wias-berlin.de

with the steady transport of a species dissolved in the fluid with a concentration $c$ and diffusion coefficient $D$ that satisfies

$$\nabla \cdot (-D\nabla c + \nabla \cdot c\mathbf{u}) = s \tag{2}$$

plus further boundary conditions.

Similar to the method proposed in [9], the Navier-Stokes equation is solved numerically with a modified nonconforming Crouzeix-Raviart finite element method. Employing a Raviart-Thomas velocity reconstruction operator, the modified Crouzeix-Raviart element stabilises the discrete Navier-Stokes solution whenever the exterior force term or the nonlinear convection term contains a large irrotational part in the sense of the Helmholtz decomposition.

Similar to [2], one can apply this reconstruction operator to the discrete velocity solution to obtain a divergence-free velocity field. This property is employed in the coupling to a transport equation which is solved by a two-point exponential fitting [1, 8, 10] finite volume method on an admissible finite volume mesh [5]. As a consequence, the maximum principle for the solution of the transport equation is preserved.

A similar, coupled FEM-FVM approach using the Scott-Vogelius finite element has been shown to be convergent [6]. It has been successfully applied to the numerical investigation of the limiting current problem in an electrochemical thin layer flow cell [7]. If successful, the advantage of the new approach introduced in the present paper would be the significantly reduced number of degrees of freedom for the FEM solution of the Navier-Stokes equations, while retaining the advantageous qualitative properties of the finite volume solution of the transport equation.

For the practical realization of the admissible finite volume mesh, Voronoi boxes have been used. It seems to be possible to realise a simpler coupling scheme of the same quality using a cell-centred finite volume method on acute simplex meshes. It is however not straightforward to generate these meshes in complex domains, or to relax this condition on the meshes. On the other hand, mesh generators which deliver boundary conforming Delaunay meshes for large classes of geometries are available [11, 12], allowing the implementation of Voronoi box based schemes [13].

## 2 Divergence-Free Reconstruction of the Velocity Field in the Navier-Stokes Equations

Let $\mathcal{T}$ denote a regular triangulation into triangles (2D) or tetrahedra (3D) in the sense of Ciarlet with nodes $\mathcal{N}$ and edges (2D) resp. faces (3D) $\mathcal{E}$. Let $P_k(\mathcal{T})$ be the space piece-wise polynomials of order $k$, i.e., every $v \in P_k(\mathcal{T})$ is a polynomial $v \in P_k(T)$ of order $k$ on every element $T \in \mathcal{T}$. The space of vector-valued nonconforming Crouzeix-Raviart functions reads

$$\mathrm{CR}(\mathscr{T}; \mathbb{R}^d) := \left\{ \mathbf{v}_h \in P_1(\mathscr{T})^d : \text{for all } T \in \mathscr{T}, \ [\mathbf{v}_h](\mathrm{mid}(E)) = \mathbf{0} \text{ for all } E \in \mathscr{E}(\Omega) \right\},$$

$$\mathrm{CR}_0(\mathscr{T}; \mathbb{R}^d) := \left\{ \mathbf{v}_h \in \mathrm{CR}(\mathscr{T}; \mathbb{R}^d) : \mathbf{v}_h(\mathrm{mid}(E)) = \mathbf{0} \text{ for all } E \in \mathscr{E}(\partial\Omega) \right\}.$$

The related nonconforming interpolation $\Pi_{\mathrm{CR}} : H^1(\Omega; \mathbb{R}^d) \to \mathrm{CR}(\mathscr{T}; \mathbb{R}^d)$ is defined by

$$(\Pi_{\mathrm{CR}}\mathbf{v})(\mathrm{mid}(E)) = \frac{1}{|E|} \int_E \mathbf{v} \, ds \quad \text{for all } E \in \mathscr{E}.$$

The unmodified Crouzeix-Raviart nonconforming finite element method for (1) seeks $\mathbf{u}_h \in \Pi_{\mathrm{CR}}\mathbf{u}_D + \mathrm{CR}_0(\mathscr{T}; \mathbb{R}^d)$ and $p_h \in Q(\mathscr{T}) := \{q_h \in P_0(\mathscr{T}) : \int_\Omega p_h \, dx = 0\}$ with

$$\int_\Omega (\mathbf{u}_h \cdot \nabla)\mathbf{u}_h \cdot \mathbf{v}_h \, dx - \int_\Omega p_h \nabla \cdot \mathbf{v}_h \, dx + \int_\Omega \eta \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, dx = \int_\Omega \mathbf{f} \cdot \mathbf{v}_h \, dx,$$

$$- \int_\Omega q_h \nabla \cdot \mathbf{u}_h \, dx = 0 \quad \text{for all } (\mathbf{v}_h, q_h) \in \mathrm{CR}_0(\mathscr{T}; \mathbb{R}^d) \times Q(\mathscr{T}).$$

$$(3)$$

The modified Crouzeix-Raviart finite element method for the Navier-Stokes equation employs the Fortin interpolation operator $\Pi_{\mathrm{RT}} : H^1(\Omega; \mathbb{R}^d) \cup \mathrm{CR}(\mathscr{T}; \mathbb{R}^d) \to \mathrm{RT}(\mathscr{T})$ into the space of Raviart-Thomas finite elements

$$\mathrm{RT}(\mathscr{T}) = \{\mathbf{v} \in H(\mathrm{div}, \Omega) \,|\, \forall T \in \mathscr{T}, \exists a \in P_0(T; \mathbb{R}^d), b \in P_0(T), \ \mathbf{v}(x)|_T = a + bx\},$$

that is uniquely defined by

$$(\Pi_{\mathrm{RT}}\mathbf{v}) \cdot \mathbf{n}_E := \frac{1}{|E|} \int_E \mathbf{v} \cdot \mathbf{n}_E \, ds \quad \text{for all } E \in \mathscr{E}.$$

The modified method seeks $\mathbf{u}_h \in \Pi_{\mathrm{CR}}\mathbf{u}_D + \mathrm{CR}_0(\mathscr{T}; \mathbb{R}^d)$ and $p_h \in Q(\mathscr{T})$ with

$$\int_\Omega (\Pi_{\mathrm{RT}}\mathbf{u}_h \cdot \nabla)\mathbf{u}_h \cdot \Pi_{\mathrm{RT}}\mathbf{v}_h \, dx - \int_\Omega p_h \nabla \cdot \mathbf{v_h} \, dx + \int_\Omega \eta \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, dx = \int_\Omega \mathbf{f} \cdot \Pi_{\mathrm{RT}}\mathbf{v}_h \, dx,$$

$$- \int_\Omega q_h \nabla \cdot \mathbf{u}_h \, dx = 0 \quad \text{for all } (\mathbf{v}_h, q_h) \in \mathrm{CR}_0(\mathscr{T}; \mathbb{R}^d) \times Q(\mathscr{T}).$$

$$(4)$$

Note, that $\Pi_{\mathrm{RT}}$ is well-defined for Crouzeix-Raviart functions $\mathbf{v}_h$, since $\int_E \mathbf{v}_h \cdot \mathbf{n}_E \, ds = |E| \mathbf{v}_h(\mathrm{mid}(E)) \cdot \mathbf{n}_E$ and $\mathbf{v}_h$ is continuous in $\mathrm{mid}(E)$. If $\mathbf{v}_h$ is piecewise divergence-free, i.e. $\nabla \cdot \mathbf{v}_h|_T = 0$ on every $T \in \mathscr{T}$, it follows that $\Pi_{\mathrm{RT}}\mathbf{v}_h \in H(\mathrm{div}, \Omega) \cap P_0(\mathscr{T}; \mathbb{R}^d)$ is divergence-free. Moreover, there is the Fortin interpolation estimate

$$\|\mathbf{v} - \Pi_{\mathrm{RT}}\mathbf{v}\|_{L^2(\Omega)} \leq Ch \|\nabla_h \mathbf{v}\|_{L^2(\Omega)} \quad \text{for all } \mathbf{v} \in H^1(\Omega; \mathbb{R}^d) \cup \mathrm{CR}(\mathscr{T}; \mathbb{R}^d).$$

A proof can be found in [3, 4] and is also valid for Crouzeix-Raviart functions. This divergence-free reconstruction of the velocity enters the finite volume method for the transport equation (2) as described below.

## 3 Finite Volume Scheme for the Transport Equation

Let $\mathscr{P}$ denote a set of points and let $\mathscr{K}$ denote the associated set of Voronoi cells with facets $\mathscr{F}$. For a point $\mathbf{x}_K \in \mathscr{P}$, the Voronoi cell $K \in \mathscr{K}$ around $\mathbf{x}_K$ is defined as the set of points $\mathbf{x} \in \Omega$ which are closer to $\mathbf{x}_K$ than to any other point in $\mathscr{P} \setminus \{\mathbf{x}_K\}$. The set of Dirichlet control volumes is denoted with $\mathscr{K}_D$ and constructed such that $x_K \in \Gamma_D$ for any $K \in \mathscr{K}_D$. Let $\mathscr{K}_0 := \mathscr{K} \setminus \mathscr{K}_D$. The set $\mathscr{K}(K)$ consists of all neighbouring cells $L \in \mathscr{K}$ such that the surface measure $|\sigma_{KL}|$ of their shared facet $\sigma_{KL} := \partial K \cap \partial L \in \mathscr{F}$ is positive. The set of facets of $K$ with nonempty intersection with the Neumann boundary $\Gamma_N$ is denoted by $\mathscr{F}_N(K)$. In practice, the Voronoi cells are constructed as the dual of a Delaunay triangulation of the domain. Therefore, due to the progress in Delaunay mesh generation algorithms [11, 12], this procedure is a constructive way to yield an admissible finite volume subdivision in the sense of [5].

The transmission coefficient along the facet $\sigma_{KL}$ of two neighbouring cells $K \in \mathscr{K}$ and $L \in \mathscr{K}(K)$ is given by

$$\tau_{\sigma_{KL}} := |\sigma_{KL}| / |\mathbf{h}_{\sigma_{KL}}| \quad \text{with} \quad \mathbf{h}_{\sigma_{KL}} := \mathbf{x}_L - \mathbf{x}_K. \tag{5}$$

The transmission coefficient for a facet $\sigma \in \mathscr{F}_N$ on the boundary reads $\tau_\sigma := |\sigma|$ and $\mathbf{h}_\sigma = \mathbf{n}_{\Gamma_N}$ equals the outer unit normal vector. The integral mean of a function $c \in L^2(K)$ over a control volume $K \in \mathscr{K}$ reads $c_K := \int_K c\,dx / |K|$. Given a divergence-free vector field $\mathbf{v} \in H(\mathrm{div}, \Omega)$ that satisfies the boundary conditions of (1) we define its scaled flux projection $v_h \in P_0(\mathscr{F})$ on the facets of the Voronoi cells by

$$v_h|_\sigma = v_\sigma := \frac{1}{|\sigma|} \int_\sigma \mathbf{v} \cdot \mathbf{h}_\sigma \, ds \quad \text{for all } \sigma \in \mathscr{F}. \tag{6}$$

This flux projection is discretely divergence-free in the following finite volume sense that

$$\sum_{L \in \mathscr{K}(K)} \tau_{\sigma_{KL}} v_{\sigma_{KL}} = \int_{\partial K} \mathbf{v} \cdot \mathbf{n}_K \, dx = \int_K \nabla \cdot \mathbf{v} \, dx = 0 \quad \text{for all } K \in \mathscr{K}. \tag{7}$$

Given $v_{\sigma_{KL}}$ from (6) and the Bernoulli function $B(z) = z/(1 - e^{-z})$, we define the exponentially fitted flux approximation [1, 8, 10]

$$g(c_K, c_L, v_{\sigma_{KL}}) := D \left( B \left( \frac{v_{\sigma_{KL}}}{D} \right) c_K - B \left( -\frac{v_{\sigma_{KL}}}{D} \right) c_L \right) \text{ for } K \in \mathscr{K}, L \in \mathscr{K}(K). \tag{8}$$

Then, the finite volume scheme seeks $c_h \in P_0(\mathcal{K})$ with $c_K = c_h|_K = c_D(\mathbf{x}_K)$ for all $K \in \mathcal{K}_D$ and

$$\sum_{L \in \mathcal{K}(K)} \tau_{\sigma_{KL}} g(c_K, c_L, v_{\sigma_{KL}}) + \sum_{\sigma \in \mathcal{F}_N(K)} \tau_\sigma g(c_K, c_K, v_\sigma) = |K| s_K \; \forall K \in \mathcal{K}_0.$$
(9)

There is an equivalent description of (9). The sum of (9) over all control volumes $K \in \mathcal{K}$ and the multiplication of a piece-wise constant test function $\lambda_h \in P_0(\mathcal{K})$ with $\lambda_K = 0$ for $K \in \mathcal{K}_D$ leads to

$$\sum_{K \in \mathcal{K}_0} \lambda_K \left( \sum_{L \in \mathcal{K}(K)} \tau_{\sigma_{KL}} g(c_K, c_L, v_{\sigma_{KL}}) + \sum_{\sigma \in \mathcal{F}_N(K)} \tau_\sigma g(c_K, c_K, v_\sigma) \right) = \sum_{K \in \mathcal{K}_0} \lambda_K |K| s_K.$$

By (8), it holds $g(c_K, c_L, v_{\sigma_{KL}}) = -g(c_L, c_K, -v_{\sigma_{KL}})$ and therefore,

$$a_h(c_h, \lambda_h) := \sum_{\sigma_{KL} \in \mathcal{F}_0} \tau_{\sigma_{KL}} g(c_K, c_L, v_{\sigma_{KL}})(\lambda_K - \lambda_L)$$

$$= \sum_{K \in \mathcal{K}_0} \lambda_K |K| s_K = \int_\Omega \lambda_h s \, dx =: F(\lambda_h).$$

Hence, an equivalent description of (9) is to search for $c_h \in P_0(\mathcal{K})$ with $c_K = c_h|_K = c_D(\mathbf{x}_K)$ for all $K \in \mathcal{K}_D$ and

$$a_h(c_h, \lambda_h) = F(\lambda_h) \quad \text{for all } \lambda_h \in P_0(\mathcal{K}) \text{ with } \lambda_K = 0 \text{ for } K \in \mathcal{K}_D.$$

Existence and uniqueness of the solutions depend on the coercivity of the bilinear form $a_h$ and have been shown e.g. in [6]. Furthermore the corresponding discretization matrix has the M-Property. If $s = 0$, for a discrete solution $(c_K)_{K \in \mathcal{K}}$, a discrete maximum principle is valid, which bounds the value $c_K$ by the values $c_L$ for $L \in \mathcal{K}(K)$, and the values at the boundary [6]. We note that the derivation of the discrete maximum principle explicitly uses the fact that the discrete velocity field $v_h$ is discretely divergence-free in the sense of (7).

## 4 The Coupling Method

Similarly to the method proposed in [6], in order to realise the coupling, we obtain a discrete solution of (1) and use the interpolation $\Pi_{\mathrm{RT}} \mathbf{u}_h$ as the velocity field $\mathbf{v}$ in (6). The numerical results obtained so far allow to put forward the hypothesis that— similar to the result in [6] which was based on the Scott-Vogelius finite element for

the Navier-Stokes equations—the coupled method is convergent. A possible convergence proof can exploit the fact that the distance between the reconstructed and the unreconstructed solution is $O(h)$ and and that the discrete velocity converges to a continuous function in $H^1(\Omega)$.

The velocity projections (6) onto the control volume faces are calculated simplex by simplex using a second order quadrature rule. For this somewhat cumbersome (especially in 3D) procedure, see [6]. We note however, that the complexity of this operation is proportional to the number of degrees of freedom if the point set $\mathscr{P}$ coincides with the set of simplex vertices, and that it can be parallelized in a straightforward manner. In the general case, with a properly implemented procedure to find the element containing a given point, optimal complexity seems to be in reach as well.

## 5 Numerical Examples

The method has been implemented within the framework of the numerical tool box pdelib2 [13] developed and maintained at WIAS.

**Convergence study.** In order to assess the convergence properties of the coupling scheme, we perform a numerical convergence study for the following coupled 2D problem taken from [6]. It is given in the unit square $\Omega$ with homogeneous Dirichlet boundary conditions. It has the exact solution

$$\mathbf{v} = \begin{pmatrix} 2(-1+x)^2x^2(-1+y)y(-1+2y) \\ 2(1-2x)(-1+x)x(-1+y)^2y^2 \end{pmatrix}$$
$$c = x^2(x-1)y(y-1).$$

The right hand sides have been chosen in such a way that they provide the indicated exact solutions.

We investigate the convection diffusion problem on two series of grids. The first series of triangular grids consisting of right angled triangles is created from a square mesh of $n \times n$ points by subdividing each square into two triangles. The second series of genuinely triangular grids is created using the mesh generator triangle [11] such that no alignment of grid points occurs. We characterise both meshes by their respective minimal edge length $h$, called mesh width.

The calculation of the flux integrals (6) is performed by a second order quadrature, the right hand side $s_K$ is calculated by point evaluation and multiplication by the size of the control volume.

The results are depicted in Fig. 1. We observe similar asymptotic behaviour for the case of exact flux and the case of a flux calculated numerically by the proposed modified Crouzeix-Raviart mixed finite element method. This allows to conjecture a similar convergence result for the coupled problem as obtained in [6]. We note that the numerical experiments for the unreconstructed method result in nearly the same

**Fig. 1** Convergence of the finite volume solution of the convection-diffusion problem on meshes of right angled triangles ("sq") resp. general triangles ("tr") for diffusion coefficient $D = 1$ (*left column*) and $D = 10^{-5}$ (*right column*), and velocities taken from the exact ("exact") resp. Raviart-Thomas reconstruction of the Crouzeix-Raviart finite element solution ("crrt") of the Stokes problem

convergence data. For the sake of readbility, these results have been omitted from the plots.

Therefore, within the context of the proposed coupling scheme, the advantage of the modified Crouzeix-Raviart mixed finite element method method lies in the existence of a discrete maximum principle for the solution of the transport problem, as will be demonstrated in the next example.

**Influence of reconstruction on discrete maximum principle.** The influence of the reconstruction can be observed in Fig. 2. The flow has been calculated using corresponding Hagen-Poiseuille velocity profiles for the inlet and the outlet boundary conditions, with maximum velocity $4.2 \times 10^9$, respectively. The U-shaped domain fits into a $5 \times 10$ rectangle, the width of the pipe is 1. At the other boundaries, no-slip boundary conditions are fulfilled. The concentration has been set to 1 at the inlet.

**Fig. 2** Stationary concentration in the longitudinal section of a *U shaped pipe*: Crouzeix-Raviart velocity without reconstruction (*left*) and with Raviart-Thomas reconstruction (*right*). Inlet and outlet are marked by the arrows. Both isolines (on the base) and elevation graph of the concentration are shown

At the outlet, an outflow boundary condition [6] is applied. At all other boundaries, impermeability conditions are imposed. The diffusion coefficient is $D = 1$.

The correct physical solution of this stationary problem with strong convection dominance is a constant concentration $c = 1$. Without reconstruction, the maximum principle for the concentration is significantly violated. The divergence-free reconstruction allows to keep the maximum principle in the numerical solution.

# References

1. Allen, D.N., Southwell, R.V.: Relaxation methods applied to determine the motion, in two dimensions, of a viscous fluid past a fixed cylinder. Quart. J. Mech. Appl. Math. **8**, 129–145 (1955)
2. Araya, R., Barrenechea, G.R., Poza, A.H., Valentin, F.: Convergence analysis of a residual local projection finite element method for the navier-stokes equations. SIAM J. Numer. Anal. **50**(2), 669–699 (2012)
3. Braess, D.: Finite Elements—Theory, Fast Solvers, and Applications in Solid Mechanics. Cambridge University Press, New York (2007)
4. Carstensen, C., Gedicke, J., Rim, D.: Explicit error estimates for Courant, Crouzeix-Raviart and Raviart-Thomas finite element methods. J. Comput. Math. **30**(4), 337–353 (2012)
5. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Handbook of Numerical Analysis, Vol. VII, pp. 713–1020. North-Holland, Amsterdam (2000)
6. Fuhrmann, J., Linke, A., Langmach, H.: A numerical method for mass conservative coupling between fluid flow and solute transport. Appl. Num. Math. **61**(4), 530–553 (2011)
7. Fuhrmann, J., Linke, A., Langmach, H., Baltruschat, H.: Numerical calculation of the limiting current for a cylindrical thin layer flow cell. Electrochim. Acta **55**, 430–438 (2009)

8. Il'in, A.M.: A difference scheme for a differential equation with a small parameter multiplying the second derivative. Mat. zametki **6**, 237–248 (1969)
9. Linke, A.: On the role of the Helmholtz decomposition in mixed methods for incompressible flows and a new variational crime. Comp. Methods Appl. Mech. Engrg. **268**, 782–800 (2014)
10. Scharfetter, D.L., Gummel, H.K.: Large signal analysis of a silicon Read diode. IEEE Trans. Electron. Dev. **16**, 64–77 (1969)
11. Shewchuk, J.R.: Triangle version 1.6. http://www.cs.cmu.edu/quake/triangle.html (2007). Accessed 31 Jan 2014
12. Si, H.: TetGen version 1.5. http://tetgen.org/ (2014). Accessed 31 Jan 2014
13. Streckenbach, T., Fuhrmann, J., et al.: Pdelib—a software toolbox for numerical computations. http://pdelib.org (2014). Accessed 31 Jan 2014

# Activity Based Finite Volume Methods for Generalised Nernst-Planck-Poisson Systems

**Jürgen Fuhrmann**

**Abstract** The paper shortly introduces models which improve the Nernst-Planck-Poisson system to obtain more realistic ion concentrations near electrode surfaces in comparison to classical models. The resulting equations are reformulated using activities as basic variables describing the species amounts. This reformulation allows to introduce a straightforward generalisation of the Scharfetter-Gummel scheme for drift-diffusion equations. Numerical examples demonstrate the improved physical correctness of the generalised model, the thermodynamic consistency in the sense of the decay of the free energy, and the usefulness in nanofluidic problems.

## 1 Ion Transport in a Fluid in Mechanical Equilibrium

We regard an incompressible isothermal mixture of $N$ components characterised by molar densities (in the sequel called concentrations) $c_\alpha$, chemical potentials $\mu_\alpha$, diffusive fluxes $\mathbf{N}_\alpha$, charge numbers $z_\alpha$, molar masses $M_\alpha$, diffusion coefficients $D_\alpha$. The component $N$ is regarded as an electroneutral ($z_N = 0$) solvent with dielectric permittivity $\varepsilon$. The evolution of the concentrations and the electrostatic potential $\phi$ is described by the Nernst-Planck-Poisson system [7, 10]:

$$-\nabla \cdot \varepsilon \nabla \phi = q = F \sum_{\alpha=1}^{N} z_\alpha c_\alpha \tag{1a}$$

$$\partial_t c_\alpha + \nabla \cdot (c_\alpha \mathbf{v} + \mathbf{N}_\alpha) = 0 \qquad (\alpha = 1 \ldots N-1) \tag{1b}$$

$$\mathbf{N}_\alpha = -\frac{D_\alpha}{RT} c_\alpha \left( \nabla \tilde{\mu}_\alpha + z_\alpha F \nabla \phi \right). \quad (\alpha = 1 \ldots N-1) \tag{1c}$$

J. Fuhrmann (✉)
Weierstrass Institute for Applied Analysis and Stochastics (WIAS),
Berlin, Germany
e-mail: juergen.fuhrmann@wias-berlin.de

The remaining notations are $T$—temperature, $R$—gas constant, $F$—Faraday constant. The effective chemical potentials $\tilde{\mu}_\alpha$ relate to the chemical potentials as

$$\tilde{\mu}_\alpha = \mu_\alpha - \frac{M_\alpha}{M_N}\mu_N. \qquad (\alpha = 1 \ldots N - 1) \qquad (2)$$

The barycentric velocity $\mathbf{v}$ of the mixture follows the incompressible Navier-Stokes equations. The assumption of mechanical equilibrium allows to set $\mathbf{v} = 0$ throughout the paper. The momentum equation results in the force balance [7, 10]:

$$\nabla p = -q\nabla\phi. \qquad (3)$$

The solvent concentration $c_N = \bar{c} - \sum_{\alpha=1}^{N-1} c_\alpha$ is the difference between the constant (due to incompressibility) concentration $\bar{c}$ of the mixture and the sum of the concentrations of the dissolved species. The solvent flux can be obtained from the condition $\sum_{\alpha=1}^{N} M_\alpha \mathbf{N}_\alpha = 0$.

This ansatz differs from the classical treatment of ion drift diffusion (see e.g. [11]) by explicitly taking into account the chemical potential of the solvent [7].

## 2 Constitutive Relationships for Chemical Potential

To close system (1a–1c), it is necessary to introduce constitutive relationships between the chemical potentials $\mu_1 \ldots \mu_N$ and the other quantities describing the system.

**Ideal dilute solution**. Here, the motion of the solvent is not influenced by the motion of the dissolved species, and with given constant reference chemical potentials $\mu_\alpha^\circ$, the chemical potential can be set to [1]:

$$\mu_N = 0, \quad \mu_\alpha = \mu_\alpha^\circ + RT \ln \frac{c_\alpha}{\bar{c}} \ (\alpha = 1 \ldots N - 1), \qquad (4)$$

corresponding to a free energy density which does not take into account the solvent:

$$\psi = \frac{1}{2}\varepsilon|\nabla\phi|^2 + RT \sum_{\alpha=1}^{N-1} c_\alpha \left(\ln \frac{c_\alpha}{\bar{c}} - 1\right). \qquad (5)$$

This ansatz regards ions as point charges and misses the fact that the finite size of real ions limits the maximum possible species concentrations $c_\alpha$.

**Bikerman model**. The introduction of an additional term in the chemical potential which takes into account the increase of the free energy resulting from the movement of a molecule into a volume already crowded with other molecules is the subject of a significant number of papers, see e.g. the reviews [3, 5].

The summary volume fraction of dissolved species amounts to $\Phi = \sum_{\alpha=1}^{N-1} v_\alpha c_\alpha$, where $v_\alpha$ is the partial molar volume accommodating 1 mole of species $\alpha$ together with the hydration shells [5]. The Bikerman model [6] assumes that all molecules are placed on a given lattice with lattice constant $a$ and that $v_\alpha = v = a^3$. A reasonable choice is $v = \frac{1}{\bar{c}}$, resulting in results in $1 - \Phi = \frac{c_N}{\bar{c}}$, the mole fraction of the solvent. A common way to incorporate the volume constraint is to set $\mu_N = 0$ and

$$\mu_\alpha = \tilde{\mu}_\alpha = \mu_\alpha^\circ + RT \ln \frac{c_\alpha}{\bar{c}} - RT \ln \frac{c_N}{\bar{c}} \qquad (\alpha = 1 \ldots N-1) \quad (6)$$

$$= \mu_\alpha^\circ + RT \ln \frac{c_\alpha}{\bar{c}} - RT \ln \left( 1 - \sum_{\alpha=1}^{N-1} \frac{c_\alpha}{\bar{c}} \right). \quad (7)$$

This ansatz introduces a nonlinear coupling between the species and corresponds to a free energy density

$$\psi = \frac{1}{2}\varepsilon|\nabla\phi|^2 + RT \sum_{\alpha=1}^{N} c_\alpha \left( \ln \frac{c_\alpha}{\bar{c}} - 1 \right). \quad (8)$$

**Dreyer et al. model.** The authors of [7] propose

$$\mu_\alpha = \mu_\alpha^\circ + \frac{1}{\bar{c}}(p - p^\circ) + RT \ln \frac{c_\alpha}{\bar{c}}, \qquad (\alpha = 1 \ldots N) \quad (9)$$

using consistent expressions for all species including the solvent. Here, $p^\circ$ is a constant reference pressure. The effective chemical potential is

$$\begin{aligned} \tilde{\mu}_\alpha = &\mu_\alpha^\circ + RT \ln \frac{c_\alpha}{\bar{c}} + \left( 1 - \frac{M_\alpha}{M_N} \right) \frac{(p - p^\circ)}{\bar{c}} \\ &- \frac{M_\alpha}{M_N} RT \ln \frac{c_N}{\bar{c}} \qquad (\alpha = 1 \ldots N-1). \end{aligned}$$

Introducing the simplifying assumption of equal molar masses $M_\alpha$ of all species including the solvent, we arrive at (6). Therefore the model of [7] appears as a consistent generalisation of the Bikerman model (6).

We just remark that it is possible to treat the general case (9) by taking the divergence on both sides of (3) arriving at a second order equation for the pressure. Due to space restrictions, all subsequent considerations will be made for the models (4) and (6). They can be readily generalized to the case (9).

## 3 Activity Based Formulation

Chemical potentials as primary variables have the disadvantage that it is hard to handle small concentrations. A common choice for the basic variables are the concentrations $c_\alpha$. The presence of $c_N$ in all species fluxes results in a coupling between the concentration gradients which is quite inconvenient to handle numerically. We make an argument in favour of an activity based re-formulation of the system. This formulation as well has its drawbacks as, for large voltage differences, the domain of values of activities may exceed the standard range of floating point implementations. In electrochemistry, these, however do not occur.

The idea is to start with the expression $\tilde{\mu}_\alpha = \tilde{\mu}_\alpha^\circ + RT \ln a_\alpha$ for the effective chemical potential which has a similar form as for a dilute solution. The activity coefficient $\gamma_\alpha$ defined by $a_\alpha = \gamma_\alpha \frac{c_\alpha}{\bar{c}}$, allows to express $c_\alpha$ through $a_\alpha$:

$$c_\alpha = \bar{c} \frac{a_\alpha}{\gamma_\alpha} = \bar{c} \beta_\alpha a_\alpha, \tag{10}$$

where $\beta_\alpha = \frac{1}{\gamma_\alpha}$ denotes the inverse activity coefficient. After a straightforward calculation, the Nernst-Planck-Poisson system (1a–1c) becomes

$$-\nabla \cdot \varepsilon \nabla \phi = q = F\bar{c} \sum_{\alpha=1}^{N-1} z_\alpha \beta_\alpha a_\alpha \tag{11a}$$

$$\partial_t (\bar{c} \beta_\alpha a_\alpha) = -\nabla \cdot \mathbf{N}_\alpha \qquad (\alpha = 1 \ldots N-1) \tag{11b}$$

$$\mathbf{N}_\alpha = -D_\alpha \bar{c} \beta_\alpha \left( \nabla a_\alpha + a_\alpha z_\alpha \frac{F}{RT} \nabla \phi \right). \quad (\alpha = 1 \ldots N-1) \tag{11c}$$

The expressions in the activities under the time derivative and the divergence operator are close to the dilute solution case. No gradient coupling is introduced.

For the dilute solution case (4), one obtains $\beta_\alpha = 1$, and the activity $a_\alpha$ is identical to the mole fraction $y_\alpha = \frac{c_\alpha}{\bar{c}}$.

For the Bikerman model (6) one obtains $\beta_\alpha = 1 - \sum_{i=1}^{N-1} \frac{c_i}{\bar{c}} =: \beta$ which is the same for all species. Expressing $c_i$ yields $\beta = 1 - \sum_{i=1}^{N-1} \beta a_i$ and

$$\beta = \frac{1}{1 + \sum_{i=1}^{N-1} a_i}, \tag{12}$$

introducing a nonlinear coupling of the species fluxes.

## 4 Equilibrium Case: Nonlinear Poisson System

Assuming $\mu_\alpha^\circ = 0$ and zero flux due to thermodynamic equilibrium, one arrives at

$$\nabla \tilde{\mu}_\alpha = -z_\alpha F \nabla \phi \qquad (\alpha = 1 \dots N - 1). \qquad (13)$$

To fulfil (13), we introduce the constant quasi-Fermi (electrochemical) potential $\psi_\alpha$ and set $\tilde{\mu}_\alpha = z_\alpha F (\psi_\alpha - \phi)$ leading to the nonlinear Poisson equation

$$-\nabla \cdot \varepsilon \nabla \phi = F \bar{c} \sum_{\alpha=1}^{N-1} z_\alpha \beta \exp\left(\frac{z_\alpha F}{RT}(\psi_\alpha - \phi)\right) \qquad (14)$$

which in the case $\beta = 1$ is exactly the Poisson-Boltzmann equation leading to the classical Gouy-Chapman theory of the electric double layer [2] .

For the Bikerman model (6), we arrive at

$$-\nabla \cdot \varepsilon \nabla \phi = F \bar{c} \frac{\sum_{\alpha=1}^{N-1} z_\alpha \exp\left(\frac{z_\alpha F}{RT}(\psi_\alpha - \phi)\right)}{1 + \sum_{\alpha=1}^{N-1} \exp\left(\frac{z_\alpha F}{RT}(\psi_\alpha - \phi)\right)}. \qquad (15)$$

For a binary 1:1 electrolyte ($N = 3, z_1 = 1, z_2 = -1$), this equation has been introduced in [6] by statistical mechanics considerations.

## 5 Finite Volume Scheme Consistent with Equilibrium

We discuss a numerical scheme for the activity based formulation of the generalised Nernst-Planck-Poisson system (11) which is consistent with the the equilibrium problem (15), i.e. the zero flux condition is consistent with the expression

$$c_\alpha = \bar{c} \beta_\alpha a_\alpha = \bar{c} \beta_\alpha \exp\left(\frac{z_\alpha F}{RT}(\psi_\alpha - \phi)\right) \qquad (16)$$

where $\phi$ is a given value of the electrostatic potential and $\psi_\alpha$ is the constant quasi-Fermi potential. In this section, we set $Z_\alpha = \frac{z_\alpha F}{RT}$ and omit the index $\alpha$.

For the concentration $c = c_\alpha$, we define a time implicit two-point flux finite volume scheme on an admissible mesh [8] (e.g. on Voronoi boxes) of control volumes $K$ containing the collocation points $\mathbf{x}_K$ in a given domain $\Omega$:

$$|K| \frac{c_K^n - c_K^{n-1}}{t^n - t^{n-1}} + \sum_{L \text{ neighbour of } K} |\partial K \cap \partial L| N_{KL}^n = 0 \qquad (17)$$

In a similar manner, one obtains for the Poisson equation

$$\sum_{L \text{ neighbour of } K} |\partial K \cap \partial L| E^n_{KL} = |K| q^n_K \tag{18}$$

Here, $c^n_K$, $q^n_K$ are the values of the concentration and the charge in the collocation points $\mathbf{x}_K$ at moment $t^n$, respectively. $N^n_{KL}$ and $E^n_{KL}$ are the respective averaged projections of the molar flux and the electric field onto the normal directions of the control volume faces $\partial K \cap \partial L$. $N^n_{KL}$ and $E^n_{KL}$ can be expressed consistently as functions of the unknown values in the control volumes $K$ and $L$. Aiming at the unconditional stability of the scheme, we chose these unknown values solely from the moment $t^n$.

In a straightforward manner, the electric field projection is expressed by

$$E^n_{KL} = \varepsilon \frac{\phi^n_K - \phi^n_L}{|\mathbf{x}_K - \mathbf{x}_L|}. \tag{19}$$

In the equilibrium, one obtains the value of $q_K$ by inserting $\phi = \phi_K$ into the different right hand side expression of (14). A correct approximation $N^n_{KL}$ should be consistent with this choice [4].

**Case $\beta = 1$.** In equilibrium, we have $c_= = \bar{c} \exp(Z(\psi - \phi))$ with a constant quasi-Fermi potential $\psi$ and a position dependent electrostatic potential $\phi$. Correspondingly, $c_{=,K} = \bar{c} \exp(Z(\psi - \phi_K))$. Consistency with equilibrium means that for such values of $c_{=,K}$, the resulting numerical flux $N_{KL}$ is zero. Using $B(\xi) = \frac{\xi}{\exp(\xi)-1}$, the Scharfetter-Gummel scheme [12]

$$N_{KL} = D \frac{B\left(Z(\phi_L - \phi_K)\right) c_K - B\left(Z(\phi_K - \phi_L)\right) c_L}{|\mathbf{x}_K - \mathbf{x}_L|} \tag{20}$$

is consistent with equilibrium: For any given constant $\psi$, assuming $N_{KL} = 0$ implies

$$\begin{aligned}
\frac{c_K}{c_L} &= \frac{B(Z(\phi_K - \phi_L))}{B(Z(\phi_L - \phi_K))} = -\frac{\exp(Z(\phi_L - \phi_K)) - 1}{\exp(Z(\phi_K - \phi_L)) - 1} \\
&= -\frac{\exp(Z\phi_L)}{\exp(Z\phi_K)} \cdot \frac{\exp(-Z\phi_K) - \exp(-Z\phi_L)}{\exp(-Z\phi_L) - \exp(-Z\phi_K)} \\
&= \frac{\exp(Z\phi_L)}{\exp(Z\phi_K)} = \frac{\exp(Z(\psi - \phi_K))}{\exp(Z(\psi - \phi_L))} = \frac{c_{=,K}}{c_{=,L}}.
\end{aligned}$$

**General $\beta$.** The flux in the activity based formulation up to the prefactor $\bar{c}\beta$ has the same structure as in the case $\beta = 1$, so we propose the ansatz

$$N_{KL} = \bar{c} D \bar{\beta} (B(Z(\phi_L - \phi_K)) a_K - B(Z(\phi_K - \phi_L)) a_L), \tag{21}$$

**Fig. 1** Negative ion concentration (*left*) and potential profile (*right*) at an ideally polarizable electrode with bulk ion concentration $c_\infty^\pm = 0.01$ mol/dm$^3$ and applied voltage of 0.5 V

**Fig. 2** Decay of free energy to equilibrium value during discharge of double layer



where $\overline{\beta}$ is some average of $\beta$ on $[a_K, a_L]$. In equilibrium, we get

$$\frac{a_K}{a_L} = \frac{\exp(Z(\psi - \phi_K))}{\exp(Z(\psi - \phi_L))}. \tag{22}$$

This is consistent with the expression for $a$ in (16) resulting in a similarly consistent expression for $c = \bar{c}\beta a$.

## 6 Numerical Examples

The examples are solved numerically using the described finite volume method implemented within pdelib [14]. The nonlinear systems are solved using Newton's method. The linear systems are solved using the direct solver Pardiso [13].

We regard an aqueous binary 1:1 electrolyte with a given molarity of the bulk solution $c_{\alpha,\infty} = c_\infty$ for $\alpha = 1, 2$. The summary concentration $\bar{c}$ is set to the molarity of water at standard conditions $\bar{c} = 55.508$ mol/dm$^3$.

**Ideally polarizable electrode.** Regard system (14) in the domain $\Omega = (0, L)$. Assume the boundary conditions $\phi|_{x=0} = \phi_0$, $\phi|_{x=L} = 0$. The quasi-Fermi potentials are obtained from given concentration values $c_\alpha|_{x=L} = c_{\alpha,\infty} << \bar{c}$ such that $q|_{x=L} = 0$. Figure 1 demonstrates the most important difference between the

$$\varepsilon\partial_{\mathbf{n}}\phi = \sigma \qquad\qquad \varepsilon\partial_{\mathbf{n}}\phi = -\sigma$$

$$c_{1,2} = c_{\text{bulk}} \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad c_{1,2} = c_{\text{bulk}}$$
$$\phi = 0 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \phi = \phi_{\text{bias}}$$

**Fig. 3** Electrolytic diode in a nanopore. If not stated otherwise: $\partial_{\mathbf{n}}\phi = 0$, $\partial_{\mathbf{n}}p = 0$, $\mathbf{N}_{1,2}\cdot\mathbf{n} = 0$, $\sigma = 250\,\mu\text{A s/m}^2$, $c_{\text{bulk}} = 2\,\text{mol/dm}^3$. Pore length = 100 nm, pore width = 2, 4, 8 nm



**Fig. 4** *Left*: IV-Curves for the different models (pore width = 2 nm). *Right*: discrepancy between the models for different pore widths

models. The dilute solution model (4) with $\beta = 1$—also called Gouy-Chapman model—overestimates the concentration close to the electrode, beyond the physical limit given by the concentration of the water molecules. The Bikermann model (6), (12) with $\beta = \frac{1}{1+a_1+a_2}$ is more realistic by limiting the solute concentration by the solvent concentration.

**Discharge of double layer.** Let $\Omega = (0, 2L)$ and solve the time dependent problem using the flux (21) starting with an initial solution obtained by charging the double layers by applying a potential $\phi|_{x=0} = \phi_0$, $\phi|_{x=2L} = -\phi_0$. We confirm that the discrete equilibrium solution from the time dependent scheme is identical to that given by the nonlinear Poisson system (14) for both cases $\beta = 1$ and $\beta = \frac{1}{1+a_1+a_2}$. After that, remove the applied potential (apply homogeneous Neumann boundary conditions) for the Poisson equation. As the potential is asymmetric with respect to $x = L$, its value in $x = L$ is always zero, and we fix this value in order to obtain uniqueness of the solution of the Poisson equation. In Fig. 2, we observe the evolution of the free energy during the approach to the equilibrium. For both models it decays monotonically, thus calling for an analysis of the scheme similar to [9].

**Electrolytic diode.** We apply the scheme in a 2D situation (Fig. 3). Figure 4 (left) demonstrates the behaviour of a diode. In difference to a semiconductor diode, the current is an ionic current, and the fixed surface charge $\pm\sigma$ plays the role of the doping. Figure 4 demonstrates the influence of the model discrepancy on the IV curve.

# References

1. Atkins, P., de Paula, J.: Atkins Physical Chemistry. Oxford University Press, Oxford (2006)
2. Bard, A.J., Faulkner, L.R.: Electrochemical Methods. Wiley, New York (1980)
3. Bazant, M.Z., Kilic, M.S., Storey, B.D., Ajdari, A.: Towards an understanding of induced-charge electrokinetics at large applied voltages in concentrated solutions. Adv. Coll. Interface Sci. **152**(1), 48–88 (2009)
4. Bessemoulin-Chatard, M.: A finite volume scheme for convection-diffusion equations with nonlinear diffusion derived from the Scharfetter-Gummel scheme. Numer. Math. **121**(4), 637–670 (2012)
5. Biesheuvel, P., Van Soestbergen, M.: Counterion volume effects in mixed electrical double layers. J. Coll. Interface Sci. **316**(2), 490–499 (2007)
6. Bikerman, J.J.: Structure and capacity of electrical double layer. Philos. Mag. **33**(220), 384–397 (1942)
7. Dreyer, W., Guhlke, C., Müller, R.: Overcoming the shortcomings of the Nernst-Planck model. Phys. Chem. Chem. Phys. **15**, 7075–7086 (2013)
8. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Handbook of Numerical Analysis, vol. VII, pp. 713–1020. Elsevier, Netherlands (2000)
9. Glitzky, A., Gärtner, K.: Energy estimates for continuous and discretized electro-reaction-diffusion systems. Nonlinear Anal. **70**(2), 788–805 (2009)
10. de Groot, S.R., Mazur, P.O.: Non-Equilibrium Thermodynamics. Dover Publications, New York (1962)
11. Newman, J., Thomas-Alyea, K.E.: Electrochemical Systems. Wiley, New York (2012)
12. Scharfetter, D.L., Gummel, H.K.: Large signal analysis of a silicon Read diode. IEEE Trans. Electron Dev. **16**, 64–77 (1969)
13. Schenk, O., Gärtner, K., Karypis, G., Röllin, S., Hagemann, M.: PARDISO solver project. http://www.pardiso-project.org (2014). Accessed 15 Jan 2014
14. Streckenbach, T., Fuhrmann, J., et al.: Pdelib—a software toolbox for numerical computations. http://www.wias-berlin.de/software/pdelib/ (2014). Accessed 15 Jan 2014

# Suitable Formulations of Lagrange Remap Finite Volume Schemes for Manycore/GPU Architectures

**Thibault Gasc and Florian De Vuyst**

**Abstract**  This paper is dedicated to Lagrange-Remap schemes (also referred to as Lagrange-Euler schemes) and their suitable formulations for manycore/GPU architectures. High performance computing efficiency requires a suitable balance between floating point operations and memory accesses, uniform compactly supported stencils, memory alignment, SIMD-based instructions and minimal dereferencing into memory. We provide various formulations, from the basis geometrical remapping to remap by flux balances and operator splitting variant approach. We present numerical experiments of two-dimensional Euler hydrodynamics on Cartesian grids up to $2048^2$ cells and provide performance results.

## 1 Introduction

For multimaterial hydrodynamics, Lagrangian methods are considered as the most accurate methods because they inherently follow material interfaces and the convection is solved by means of the moving mesh. Unfortunately, they often show a lack of robustness due to the possible cell degeneracy. Arbitrary Lagrangian-Eulerian (ALE) methods try to keep most of the Lagrangian accuracy while smoothing the mesh (if needed) by some remap processing. The limit case is the so-called Lagrange-

T. Gasc (✉)
Maison de la Simulation, USR 3441, CEA - CNRS - INRIA - University Paris - Sud - University Versailles, 91191 Gif-sur-Yvette, France
e-mail: thibault.gasc@cea.fr

T. Gasc
CEA DIF Bruyère-le-Châtel, 91297 Arpajon, France

F. De Vuyst
Centre de Mathématiques et de Leurs Applications, CMLA UMR 8536, ÉNS CACHAN, 61 avenue du président Wilson, 94235 Cachan, France
e-mail: devuyst@cmla.ens-cachan.fr

Remap (or Lagrange-Euler) solver, pioneered by Von Neumann and Richtmyer in 1950 [5], where a remap step is performed after each Lagrangian time advance. Lagrange-Remap schemes are known to have attractive mathematical properties, and are still subject to developments and analyses [6]. However, the portability of these methods on large scale parallel computers is still subject to ongoing research [2].

Compressible Euler equations written under a Lagrangian integral form over a volume $\mathscr{V}^t$ moving at velocity $\mathbf{u}$ are written as follows:

$$\frac{d\mathscr{V}^t}{dt} = \int_{\mathscr{S}^t} \mathbf{u} \cdot \mathbf{n} , \quad \frac{d}{dt}(\int_{\mathscr{V}^t} \rho) = 0 , \tag{1}$$

$$\frac{d}{dt}(\int_{\mathscr{V}^t} \rho\mathbf{u}) = -\int_{\mathscr{S}^t} P\mathbf{n} , \quad \frac{d}{dt}(\int_{\mathscr{V}^t} \rho E) = -\int_{\mathscr{S}^t} P\mathbf{u} \cdot \mathbf{n}, \tag{2}$$

with $\mathscr{S}^t = \partial\mathscr{V}^t$. The Lagrangian step simply consists in discretizing these equations by choosing some moving control volumes. The remap step is a geometrical projection over fixed (Euler) control volumes:

$$\int_{\mathscr{V}_i} q = \sum_j \int_{\mathscr{V}_i \cap \mathscr{V}_j^{Lag}} q_j^{Lag}, \tag{3}$$

where $\mathscr{V}_i$ is the volume of cell $i$ in the fixed Eulerian grid, $\mathscr{V}_j^{Lag}$ the volume of the cell $j$ in the Lagrangian grid obtained from the Lagrangian step (1) and $q_j^{Lag}$ the value of the quantity which is remapped ($\rho$, $\rho u$ or $\rho E$) on the Lagrangian cell $j$; this value are obtained from the previous Lagrangian step (1–2). For the sake of simplicity, we here assume a perfect gas equation of state $P = (\gamma - 1)\rho e$, $\gamma \in (1, 3]$, where the speed of sound $c$ is given by $c = \sqrt{\frac{\gamma P}{\rho}}$.

The Lagrangian step does not introduce specific choices nor difficulties. But building an accurate and computing efficient remap can be challenging. Indeed one has to compute many volume intersections and accuracy of this operation strongly impacts the order of the scheme. In order to avoid computing complex exact volume intersections, various techniques can be used such as Alternating Direction methods.

## 2 Lagrange Remap in Finite Volume Formalism

To avoid complex geometrical computations, Lagrange Remap schemes can be approximated by using Finite Volume formalism. Specific fluxes can be defined to approximate both steps. For example, the energy equation integrated in time between $n$ and $n + 1$ can be written as follows:

**Fig. 1** Remap operation interpreted as a flux balance: Eulerian fixed cell $(i, j)$ with material velocity at the nodes, Lagrangian deformed cell $(i, j)$ and definition of areas corresponding to the fluxes through interfaces of the Eulerian cell



$$\int_{\mathscr{V}^{n+1}} \rho E - \int_{\mathscr{V}^n} \rho E = -\int_{t^n}^{t^{n+1}} \int_{\mathscr{S}(t)} P\mathbf{u} \cdot \mathbf{n},$$

$$(mE)^L - (mE)^n \approx -\Delta t \sum_{S \subset \mathscr{S}} |S| P_s \mathbf{u}_s \cdot \mathbf{n}_s,$$

where $\mathscr{S}$ is a mean border of $\mathscr{V}$ between $t$ and $t^{n+1}$ and $\{S\}$ define a partition of $\mathscr{S}$. The latest form is a pure flux formulation. The difference between the flux and the pure Lagrangian step rely on the non trivial integration $\int_{t^n}^{t^{n+1}} \int_{\mathscr{S}(t)}$. The mass and momentum equations can be also written with the same formalism.

The Remap step can also be rewritten as flux balance. In order to only use fluxes to define an approximate Remap step, it should be noticed that a flux can only be defined between two volumes (or area in 2D) sharing a face (or edge in 2D). The Remap operation has been firstly described as a pure geometrical process (3). It can also be understood as a redistribution process from a deformed grid to a fixed one. Any quantity that has traveled too much during the Lagrangian step should be assigned to a new location, which means a new volume. Quantities that have left the fixed volume during the Lagrangian step have crossed an interface. This allows us to define the remap operation by using fluxes. Since we construct this fluxes with Lagrangian quantities which cross interface, Lagrangian upwind values of the corresponding quantities are used to define the fluxes. The flux should match with the area swept by the edge during the Lagrangian step $\int_{t^n}^{t^{n+1}} \int_{\mathscr{S}_i(t)} dS_i dt$ .

Any quantity that crosses multiple interface will be counted once for each crossed interface. For example, in Fig. 1, the quantity in the triangle at left top corner of the cell, is counted in both fluxes $\Phi_{i-\frac{1}{2},j}$ and $\Phi_{i,j+\frac{1}{2}}$. This allows us to take into account fluxes between cells that do not share a face. This kind of exchanges is naturally defined in a geometrical remap.

In the 1D case, the first order Lagrange Remap cell centered scheme can be exactly described using fluxes formalism. The iteration process can be written as follows:

$$m_j^{n+1} = m_j^n - \Delta t \left( (\Phi_\rho)_{j+\frac{1}{2}} - (\Phi_\rho)_{j-\frac{1}{2}} \right),$$

$$(mu)_j^{n+1} = (mu)_j^n - \Delta t \left( P_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} \right) - \Delta t \left( (\Phi_{\rho u})_{j+\frac{1}{2}} - (\Phi_{\rho u})_{j-\frac{1}{2}} \right),$$

$$(mE)_j^{n+1} = (mE)_j^n - \Delta t \left( P_{j+\frac{1}{2}} u_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} u_{j-\frac{1}{2}} \right) - \Delta t \left( (\Phi_{\rho E})_{j+\frac{1}{2}} - (\Phi_{\rho E})_{j-\frac{1}{2}} \right),$$

where $(\Phi_q)_{j+\frac{1}{2}} = u_{j+\frac{1}{2}} \frac{q_j^L + q_{j+1}^L}{2} + sgn(u_{j+\frac{1}{2}}) \frac{q_j^L - q_{j+1}^L}{2}$ is the upwind flux associated to the Lagrangian quantity $q^L$. Lagrangian values are obtained thanks to the explicit discretization of Eqs. (1–2):

$$\mathscr{V}_j^L = \mathscr{V}_j^n + \Delta t(u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}}), \quad \rho_j^L = \rho_j^n \frac{\mathscr{V}_j^n}{\mathscr{V}_j^L},$$

$$u_j^L = u_j^n - \frac{\Delta t}{m_j^n}(P_{j+\frac{1}{2}} - P_{j-\frac{1}{2}}), \quad E_j^L = E_j^n - \frac{\Delta t}{m_j^n}(P_{j+\frac{1}{2}} u_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} u_{j-\frac{1}{2}}),$$

where $P_{j+\frac{1}{2}}$ and $u_{j+\frac{1}{2}}$ are estimated as solution of approximate Riemman solver.

## 3 Splitting Approach

We introduce a two steps algorithm by using a splitting strategy [3]. The first step describes the pressure waves propagation while the second describes pure advection.

### 3.1 One-Dimensional Formulation

We introduce the Euler equations in the formal nonconservative form: denoting by $W = (\rho, u, P)^T$, we have

$$\partial_t W + (uI + B)\partial_x W = 0, \quad \text{where } B = \begin{pmatrix} 0 & \rho & 0 \\ 0 & 0 & \tau \\ 0 & \rho c^2 & 0 \end{pmatrix}.$$

This hyperbolic system has 3 eigenvalues $(u - c, \ u, \ u + c)$. We split the system in the two following ones:

$$\partial_t W + B\partial_x W = 0, \tag{4}$$

$$\partial_t W + uI\partial_x W = 0. \tag{5}$$

The system (4) is hyperbolic with eigenvalues $(-c, 0, c)$. He describes the pressure waves propagation. The system (5) is hyperbolic and has $u$ as a triple eigenvalue. Using the conservative variables, this splitting approach leads to these two systems:

$$\partial_t U + \partial_x \Pi + U \partial_x u = 0, \tag{6}$$

$$\partial_t U + \partial_x (uU) - U \partial_x u = 0, \tag{7}$$

where $U = (\rho, \rho u, \rho E)^T$ are the conservatives variables and $\Pi = (0, P, Pu)^T$ is the pressure contribution.

We propose a 1D discretization of the presented system. The common intermediary state is denoted $(.)^\star$. From $t^n$ to $t^\star$, a implicit discretization of (6) is used and from $t^\star$ to $t^{n+1}$, a explicit discretization of (7) is used. Discretization of the pressure system (6) gives:

$$\rho_j^\star = \rho_j^n - \frac{\Delta t}{h} \rho_j^\star (u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}}), \tag{8}$$

$$\rho_j^\star u_j^\star = \rho_j^n u_j^n - \frac{\Delta t}{h}(P_{j+\frac{1}{2}} - P_{j-\frac{1}{2}}) - \frac{\Delta t}{h}(\rho_j^\star u_j^\star)(u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}}),$$

$$\rho_j^\star E_j^\star = (\rho E)_j^n - \frac{\Delta t}{h}(P_{j+\frac{1}{2}} u_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} u_{j-\frac{1}{2}}) - \frac{\Delta t}{h}(\rho_j^\star E_j^\star)(u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}}).$$

And the discretization of the advection system (7) gives:

$$U_j^{n+1} = U_j^\star - \frac{\Delta t}{h}\left((uU)_{j+\frac{1}{2}}^\star - (uU)_{j-\frac{1}{2}}^\star\right) + \frac{\Delta t}{h}U_j^\star(u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}}).$$

Please note the following important aspects of the discretization:

- Simplification of discrete non conservative terms $U\Delta u$ is needed to build a global conservative scheme from 2 non conservative steps,
- Intercells values $u_{j+\frac{1}{2}}$ and $P_{j+\frac{1}{2}}$ are defined as values at time $t^n + \Delta t/2$ of approximate solutions of the Riemann problem $\tilde{W}(U_j^n, U_{j+1}^n)$, and can be estimated at time $t^n$ explicitly,
- Thanks to Eq. (9), the implicit discretization of the pressure system can be expressed explicitly in variables $\rho$, $u$ and $E$: $\rho_j^\star = \frac{\rho_j^n}{1 + \frac{\Delta t}{h}(u_{j+\frac{1}{2}} - u_{j-\frac{1}{2}})}$, $u_j^\star = u_j^n - \frac{\Delta t}{h\rho_j^n}(P_{j+\frac{1}{2}} - P_{j-\frac{1}{2}})$, $E_j^\star = E_j^n - \frac{\Delta t}{h\rho_j^n}(P_{j+\frac{1}{2}} u_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} u_{j-\frac{1}{2}})$. This expressions match with the 1D discretization of the Lagrangian step of the Lagrange-Remap algorithm but the associated grid is fixed and does not follow the material as previously.
- $(uU)_{j+\frac{1}{2}}^\star$ can be seen as an intercell flux. In order to build a scheme that matches perfectly (in 1D) with the Lagrange-Remap scheme, we define these fluxes as follows: $(uU)_{j+\frac{1}{2}}^\star = u_{j+\frac{1}{2}} U_{j+\frac{1}{2}}^{\star,upwind}$, where $U_{j+\frac{1}{2}}^{\star,upwind} = \frac{U_j^\star + U_{j+1}^\star}{2} + sgn(u_{j+\frac{1}{2}})\frac{U_j^\star - U_{j+1}^\star}{2}$.

Thanks to the previous choices, the presented scheme can be written in the following conservative form:

$$U_j^{n+1} = U_j^n - \frac{\Delta t}{h}\left((uU)_{j+\frac{1}{2}}^\star - (uU)_{j-\frac{1}{2}}^\star\right) - \frac{\Delta t}{h}\begin{pmatrix} 0 \\ P_{j+\frac{1}{2}} - P_{j-\frac{1}{2}} \\ P_{j+\frac{1}{2}}u_{j+\frac{1}{2}} - P_{j-\frac{1}{2}}u_{j-\frac{1}{2}} \end{pmatrix}.$$

(9)

### 3.2 Multidimensional Extension

The previous 1D scheme can be easily extended to 2D and 3D discretizations. Since no deformation of the grid is used here, the extension is easier than Lagrange-Remap for multidimensional problem.

For example, using the same splitting strategy on the 2D Euler equation leads to the two following system:

$$\partial_t U + \nabla \cdot \Pi + U\nabla \cdot \mathbf{u} = 0,$$
$$\partial_t U + \nabla \cdot (\mathbf{u}U) - U\nabla \cdot \mathbf{u} = 0,$$

where $U = (\rho, \rho u, \rho v, \rho E)^T$, $\Pi = \begin{pmatrix} 0 & P & 0 & Pu \\ 0 & 0 & P & pv \end{pmatrix}^T$, and $\mathbf{u} = (u, v)^T$. Discretization of this systems leads to the following scheme:

$$\rho_{i,j}^\star = \frac{\rho_{i,j}^n}{1 + \frac{\Delta t}{h}(u_{i+\frac{1}{2},j} - u_{i-\frac{1}{2},j} + v_{i,j+\frac{1}{2}} - v_{i,j-\frac{1}{2},j})},$$

$$u_{i,j}^\star = u_{i,j}^n - \frac{\Delta t}{\rho_{i,j}^n h}(P_{i+\frac{1}{2},j} - P_{i-\frac{1}{2},j}), \qquad v_{i,j}^\star = v_{i,j}^n - \frac{\Delta t}{\rho_{i,j}^n h}(P_{i,j+\frac{1}{2}} - P_{i,j-\frac{1}{2}}),$$

$$E_{i,j}^\star = E_{i,j}^n - \frac{\Delta t}{\rho_{i,j}^n h}(P_{i+\frac{1}{2},j}u_{i+\frac{1}{2},j} - P_{i-\frac{1}{2},j}u_{i-\frac{1}{2},j})$$

$$- \frac{\Delta t}{\rho_{i,j}^n h}(P_{i,j+\frac{1}{2}}v_{i,j+\frac{1}{2}} - P_{i,j-\frac{1}{2}}v_{i,j-\frac{1}{2}}),$$

$$U_{i,j}^{n+1} = U_{i,j}^\star - \frac{\Delta t}{h}\left((uU)_{i+\frac{1}{2},j}^\star - (uU)_{i-\frac{1}{2},j}^\star + (vU)_{i,j+\frac{1}{2}}^\star - (vU)_{i,j-\frac{1}{2}}^\star\right)$$

$$+ \frac{\Delta t}{h}U_{i,j}^\star(u_{i+\frac{1}{2},j} - u_{i+\frac{1}{2},j} + v_{i,j+\frac{1}{2}} - v_{i,j-\frac{1}{2}}).$$

The fluxes and values at the intercell $(.)_{i+\frac{1}{2},j}$ or $(.)_{i,j+\frac{1}{2}}$ are estimated by an upwinding strategy or a 1D approximate Riemann solver in the corresponding direction. It may be not directly clear that this numerical scheme is conservative, but, it can be shown that it is, using multidimensional version of (9).

An extended version of the scheme has also been implemented: to improve accuracy, gradients of conservative variables are reconstructed before the computationof

the fluxes, in the MUSCL spirit [4]. A *minmod* limiter is used for the slopes. Note that this improvement does not leads to a complete second order in space scheme but notably improves the resolution of contact waves.

## 4  A Practical GPU Implementation

The latest algorithm has been implemented using NVIDIA's SDK language CUDA for GPU in two versions, with and without the gradients reconstruction. The implementation has been tested on a NVIDIA K20M device (2496 cores, 5 GB memory). Since host-device memory transfers are slow, they should be limited to the minimum needs of the application or overlap with kernels executions. Here, data are all allocated in the GPU memory and only data needed for the output files are mirrored in the CPU memory and updated via device to host transfers when required. When running performance analyses, output writing and data transfers are disabled.

The algorithm is divided into 5/6 kernels, each performing one elementary operation: maximum velocity global reduce for CFL condition, Riemann solver, pressure wave propagation, (gradient reconstruction,) flux computation, and final update including advection wave propagation. By doing this decomposition, we build rather small kernels which can fit with the small number of registers per block available on the GPU. These kernels have similar size and the time spent in the different kernels are quite similar, going from 15 to 30 % of the total execution time.

During the optimization process, we try to merge or split some kernels (for example splitting a 2D kernel into a $X$ kernel and a $Y$ one), but the initial decomposition seems to be the most efficient. We provide kernel per kernel analysis and and describe some changes performed during the optimization process. We hope this will help the reader to understand both the implementation and how we try to use as efficiently as possible the GPU. The global reduce kernel is used to compute the maximum velocity from which we compute $\Delta t$ with a given CFL number to ensure the stability of the scheme. The kernel was optimized using techniques presented in the example of a global reduce sum given in the NVIDIA's SDK. The execution time dropped from 15 % to less than 4 % of the total time execution. The other kernels (except the Riemann solver) are limited by the memory bandwidth using 75–85 % of the theoretical memory peak and 15–35 % of the theoretical arithmetic operations peak. The Riemann solver kernel uses about 50 % of the memory peak, and its performance is probably limited by latency. Several reformulations of the kernels where tried (including usage of shared memory, spitting into 2 $(x + y)$ kernels) but we did not reach a better performance.

Since the implementation is mainly memory bounded, we check that we perform efficient read and write operations. Using a Cartesian grid allows us to perform easily coalesced memory access. However, loading data from the neighborhood introduces some misaligned coalesced reads. To avoid misaligned reads and to take advantage of some data reuse, we tried to use shared memory. Unfortunately, kernels using shared memory were not more efficient, mainly because the data reuse is rather

**Fig. 2** Density at time $t = 0.65$—first order scheme on $512 \times 512$ grid



small and the cache memory system of Kepler architectures is efficient. At the very end of the optimization process, we forced the device to store some reused data in registers or in the constant memory. The time execution is distributed among the kernels in the following way: global reduce 4 %, Riemann solver 16 %, pressure wave propagation 14 %, gradient reconstruction 13 %, flux computation 29 %, advection wave propagation 24 %.

## 5 Numerical Results

We use a 2D test case. This test can be seen as an axisymmetrical extension of Sod shock tube. At time $t = 0$, a light fluid at low pressure fulfills a bubble at the center of a square domain. The bubble diameter d is defined such as $d^2 = 0.5$. A dense fluid at high pressure fulfills the remaining space. The case is periodic in both $x$ and $y$ directions. The following values are used to define the initial state: $\rho_1 = 1, P_1 = 1$ and $\rho_2 = 0.125, P_2 = 0.1$. We use a perfect gas with the adiabatic coefficient of 1.4. As in the Sod shock tube, rarefaction waves, shock waves and contact discontinuities appear. The chosen geometry leads to multiple wave interactions. Interactions between shock waves and contact discontinuities produce instabilities. Simulations are run on both $512 \times 512$ and $2048 \times 2048$ Cartesian mesh using double-precision floating points number. Figures 2, 3, 4 and 5 show the density at time $t = 0.65$. Using a CFL number of 0.49, about 3550 time iterations are done on the $512 \times 512$ grid (14600 on the $2048 \times 2048$ grid). Since the first order scheme is quite diffusive, it is not able to capture instabilities. The first order version updates $231.10^6$ ($247.10^6$) cells×iterations/sec, and the minmod version $176.10^6$ ($189.10^6$) cells × iterations/s. Given the precision improvement, the loss of about 25 % performance is acceptable. When using a larger grid ($2018^2$ instead of $512^2$), the computation is less sensitive to side effects and run faster, updating more cells each second (but more iterations are needed to reach the final time).

**Fig. 3** Density at time $t = 0.65$—scheme with min-mod flux limiter on $512 \times 512$ grid



**Fig. 4** Density at time $t = 0.65$—first order scheme on $2048 \times 2048$ grid



**Fig. 5** Density at time $t = 0.65$—scheme with minmod flux limiter on $2048 \times 2048$ grid



## 6 Concluding Remarks and Perspectives

In this paper we have proposed a suitable formulation of Lagrange-Remap schemes for GPU or manycore architecture. A variant algorithm based on operator splitting has been introduced. This two-step algorithm perfectly matches with the Lagrange-

Remap in 1D first-order scheme. Extension to the multidimensional case can be done without any geometrical projection which is time consuming. A 2D GPU implementation has been proposed and performance appears to be very satisfactory. In a future work, second order schemes will be constructed and comparisons with standard forms of Lagrange-Remap algorithms such as BBC [7] or others [1] will be done. This splitting approach will be coupled with an interface capturing method in order to compute multimaterial flows.

# References

1. Benson, D.J.: Computational methods in Lagrangian and Eulerian hydrocodes. Comput. Methods Appl. Mech. Engrg. **99**, 235–394 (1992)
2. Herdman, J.A., Gaudin, W.P., et al.: Accelerating hydrocodes with OpenACC, OpenCL, and CUDA. In: 3rd International Workshop on Performance Modeling, Benchmarking and Simulation (PMBS12), held as part of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC12) (2012)
3. Strang, G.: On the construction and comparison of difference schemes. SIAM J. Numer. Anal. **5**(3), 506–517 (1968)
4. Van Leer, B.: Towards the ultimate conservative difference scheme. ii. Monotonicity and conservation combined in a second-order scheme. J. Comput. Phys. **14**(4), 361–370 (1974)
5. Von Neumann, J., Richtmyer, R.D.: A method for the numerical calculation of hydrodynamic shicks. J. Appl. Phys. **21**(3), 232–237 (1950)
6. Waltz, J.: Operator splitting and time accuracy in Lagrange plus remap solution methods. J. Comput. Phys. **253**, 247–258 (2013)
7. Woodward, P., Colella, P.: The numerical simulation of two-dimensional fluid flow with strong shocks. J. Comput. Phys. **54**(1), 115–173 (1984)

# Efficient Parallel Simulation of Atherosclerotic Plaque Formation Using Higher Order Discontinuous Galerkin Schemes

**Stefan Girke, Robert Klöfkorn and Mario Ohlberger**

**Abstract** The compact Discontinuous Galerkin 2 (CDG2) method was successfully tested for elliptic problems, scalar convection-diffusion equations and compressible Navier-Stokes equations. In this paper we use the newly developed DG method to solve a mathematical model for early stages of atherosclerotic plaque formation. Atherosclerotic plaque is mainly formed by accumulation of lipid-laden cells in the arterial walls which leads to a heart attack in case the artery is occluded or a thrombus is built through a rupture of the plaque. After describing a mathematical model and the discretization scheme, we present some benchmark tests comparing the CDG2 method to other commonly used DG methods. Furthermore, we take parallelization and higher order discretization schemes into account.

## 1 Introduction

Atherosclerotic plaque formation is today seen as a chronic inflammation of the arterial wall which grows over decades and may finally lead to a heart attack in case the artery is occluded or a thrombus is built through a rupture of the plaque. To understand the mechanisms of the chronic inflammation it was recently shown in [11] that genetically modified mice with a cuff around their carotid develop atherosclerotic

S. Girke (✉) · M. Ohlberger
Institute for Computational and Applied Mathematics, University of Münster,
Einsteinstraße 62, 48149 Münster, Germany
e-mail: stefan.girke@wwu.de

M. Ohlberger
e-mail: mario.ohlberger@wwu.de

R. Klöfkorn
National Center for Atmospheric Research, 1850 Table Mesa Drive,
Boulder, CO 80305, USA
e-mail: robertk@ucar.edu

plaque formation up- and downstream of the cuff after they were fed with a high cholesterol diet. A low wall shear stress of the blood onto the arterial wall or highly oscillating blood flow was shown to be an important indicator for the development of plaque because it damages the endothelial layer which initiates the inflammation process in the arterial wall.

At this point our mathematical model (cf. [7]) comes into play which we want to present in Sect. 2: A dysfunction of the endothelial allows low-density lipoproteins (LDL) to enter the artery wall. Once inside the arterial wall, the LDL becomes oxidized which leads to a recruitment of immune cells. The immune cells differentiate into active macrophages when inside the arterial wall starting continuously absorbing the oxidized LDL. Finally, the macrophages differentiate into foam cells, die and build a necrotic core. Smooth muscle cells (SMCs) from the outer regions of the arterial wall can migrate into the lesion and either become an apoptotic cell or migrate around the lesion to form a fibromuscular cap overlaying the plaque. The blood flow in the artery, the wall shear stress onto the endothelial layer and other mechanics are neglected in our model because we concentrate on the inflammation part. Section 3 describes the spatial and temporal discretization of the CDG2 method which was successfully tested for elliptic problems, scalar convection-diffusion equations and compressible Navier-Stokes equations in [3, 8, 9]. We summarize our paper with some 2D and 3D benchmark tests[1] in Sect. 4 and a conclusion in Sect. 5.

## 2 Mathematical Model for Atherosclerotic Inflammation

A variety of mathematical models dealing with atherosclerotic plaque formation exist (cf. [4, 7]). Let $\Omega \subset \mathbb{R}^d$, $d = 2, 3$ be the domain of the arterial wall, $\Gamma_1$ the boundary between the arterial wall and the lumen and $\Gamma_2$ the outer boundary of the arterial wall. Moreover, let $U = (u_1, \ldots, u_6)$ be a vector with six (cellular or molecular) species, where $u_1$ denotes immune cells (i.e. macrophages), $u_2$ SMCs, $u_3$ debris (dead or apoptotic cells), $u_4$ a chemoattractant, $u_5$ non oxidized and $u_6$ oxidized LDL. Then our inflammation model is defined by

$$\partial_t U = -\nabla \cdot (\mathscr{F}(U) - \mathscr{A}(U)\nabla U) + S(U) \quad \forall x \in \Omega, t > 0 \tag{1}$$

with

$$\mathscr{F}(U) = \begin{pmatrix} \chi_{14}\nabla u_4 + \chi_{16}\nabla u_6 \\ \chi_{24}\nabla u_4 - \chi_{21}\nabla u_1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad S(U) = \begin{pmatrix} -d_1 u_1 \\ -d_2 u_2 \\ d_1 u_1 + d_2 u_2 + F_0 u_1 u_3 \\ -\alpha_1 u_1 u_4 - \alpha_2 u_2 u_4 + \gamma u_3 \\ -k u_5 \\ k u_5 \end{pmatrix}, \tag{2}$$

---

[1] Detailed benchmark data: `wwwmath.uni-muenster.de/u/stefan.girke/bmark`.

$$\mathscr{A}(U) := \mathrm{diag}(\mu_1, \mu_2, \mu_3, \mu_4, \mu_5, \mu_6), \quad \chi_{ij} := \chi_{ij}^0 \frac{u_i}{u_j + \chi_{ij}^{\mathrm{th}}}. \tag{3}$$

All parameters are chosen constant and positive (except $\chi_{ij}$), although it is possible (and necessary) to choose them more general, see [7]. All cellular and molecular species are more or less motile by diffusion coefficients $\mu_i$.

The parameters $d_1$ and $d_2$ describe the death rates of immune cells and SMCs. Chemoattractant is neutralized by immune cells and SMCs which is described by $\alpha_1$ and $\alpha_2$. The parameter $k$ describes how fast the native LDL becomes oxidized.

The functions $\chi_{ij}$ are called tactic sensitivity functions. We want to mimic a high sensitivity of a species to the relative gradient of another species on the one hand and a small penalization term $\chi_{ij}^{\mathrm{th}}$ to regularize the tactic movement for small concentrations on the other hand. A lot of other tactic sensitivity functions are possible as well. Our tactic sensitivity functions are defined by constants $\chi_{ij}^0$ and $\chi_{ij}^{\mathrm{th}}$.

A positive $F_0$ indicates a diseased state which may lead to an increase of debris. For a healthy immune system debris would be degraded which is indicated by a negative $F_0$. The parameter $\gamma$ is a production term which is debris dependent.

Boundary conditions are given by

$$\partial_n u_1 = -\beta_1 H(u_4 - u_4^*), \quad \forall x \in \Gamma_1, t > 0, \tag{4}$$

$$\partial_n u_2 = -\beta_2 H(u_4 - u_4^{**}), \quad \forall x \in \Gamma_2, t > 0, \tag{5}$$

$$\partial_n u_5 = -\sigma, \qquad \forall x \in \Gamma_{1,\mathrm{in}} \subset \Gamma_1, t > 0 \tag{6}$$

with Heaviside function $H$ and no-flow for all other boundary conditions. Initial data is be given by some function $u_i(x, 0) = u_i^0(x), i = 1, \ldots, 6, x \in \Omega$. We allow LDL and immune cells to enter the arterial wall through the inner boundary and SMCs to enter through the outer arterial wall. The immune cell (SMC) inflow is triggered when a threshold $u_4^*$ ($u_4^{**}$) of chemoattractant is exceeded. Here, $\beta_1$, $\beta_2$ and $\sigma$ denote constant inflow rates for immune cells, SMCs and LDL, respectively.

## 3 Discretization

The considered discretization is based on the Discontinuous Galerkin (DG) approach and implemented in DUNE- FEM [6] a module of the DUNE framework [2]. The current state of development allows for simulation of convection dominated (cf. [5]) as well as viscous flow (cf. [3]). We consider the CDG2 method from [3] for various polynomial orders in space and 2nd (or 3rd) order in time for the numerical investigations carried out in this paper.

### *3.1 Spatial Discretization*

The spatial discretization is derived in the following way. Given a tessellation $\mathcal{T}_h$ of the domain $\Omega$ with $\cup_{K \in \mathcal{T}_h} K = \Omega$ the discrete solution $U_h$ is sought in the piecewise polynomial space

$$V_h = \{v \in L^2(\Omega, \mathbb{R}^{n_{spec}}) : v|_K \in [\mathscr{P}_k(K)]^{n_{spec}}, \ K \in \mathcal{T}_h\} \quad \text{for some } k \in \mathbb{N},$$

where $n_{spec}$ is the number of species and $\mathscr{P}_k(K)$ is a space containing polynomials up to degree $k$.

We denote with $\Gamma_i$ the set of all intersections between two elements of the grid $\mathcal{T}_h$ and accordingly with $\Gamma$ the set of all intersections, also with the boundary of the domain $\Omega$. The following discrete form is not the most general but still covers a wide range of well established DG methods. For all basis functions $\varphi \in V_h$ we define

$$\langle \varphi, \mathscr{L}_h(U_h) \rangle := \langle \varphi, \mathscr{K}_h(U_h) \rangle + \langle \varphi, \mathscr{I}_h(U_h) \rangle \tag{7}$$

with the element integrals

$$\langle \varphi, \mathscr{K}_h(U_h) \rangle := \sum_{K \in \mathcal{T}_h} \int_K \left( (\mathscr{F}(U_h) - \mathscr{A}(U_h)\nabla U_h) : \nabla \varphi + S(U_h) \cdot \varphi \right), \tag{8}$$

and the surface integrals (by introducing appropriate numerical fluxes $\widehat{\mathscr{F}}_e$, $\widehat{\mathscr{A}}_e$ for the convection and diffusion terms, respectively)

$$\langle \varphi, \mathscr{I}_h(U_h) \rangle := \sum_{e \in \Gamma_i} \int_e \left( \{\{\mathscr{A}(U_h)^T \nabla \varphi\}\}_e : [\![ U_h ]\!]_e + \{\{\mathscr{A}(U_h)\nabla U_h\}\}_e : [\![ \varphi ]\!]_e \right)$$

$$- \sum_{e \in \Gamma} \int_e \left( \widehat{\mathscr{F}}_e(U_h) - \widehat{\mathscr{A}}_e(U_h) \right) : [\![ \varphi ]\!]_e, \tag{9}$$

where $\{\{V\}\}_e = \frac{1}{2}(V^+ + V^-)$ denotes the average and $[\![ V ]\!]_e = (n^+ \otimes V^+ + n^- \otimes V^-)$ the jump of the discontinuous function $V \in V_h$ over element boundaries. For matrices $\sigma, \tau \in \mathbb{R}^{m \times n}$ we use standard notation $\sigma : \tau = \sum_{j=1}^m \sum_{l=1}^n \sigma_{jl}\tau_{jl}$. Additionally, for vectors $v \in \mathbb{R}^m$, $w \in \mathbb{R}^n$, we define $v \otimes w \in \mathbb{R}^{m \times n}$ according to $(v \otimes w)_{jl} = v_j w_l$ for $1 \leq j \leq m$, $1 \leq l \leq n$.

The convective numerical flux $\widehat{\mathscr{F}}_e$ can be any appropriate numerical flux known for standard finite volume methods. For the results presented in this paper we choose $\widehat{\mathscr{F}}_e$ to be the widely used local Lax-Friedrichs numerical flux function.

A wide range of diffusion fluxes $\widehat{\mathscr{A}}_e$ can be found in the literature, for a summary see [1]. We choose the CDG2 flux

$$\widehat{\mathscr{A}}_e(V) := 2\chi_e \left( \mathscr{A}(V) r_e([\![ V ]\!]_e) \right)|_{K_e^-} \quad \text{for } V \in V_h, \tag{10}$$

which was shown to be highly efficient for advection-diffusion equations (cf. [3]). Based on stability results, we choose $K_e^-$ to be the element adjacent to the edge $e$ with the smaller volume. $r_e([\![V]\!]_e) \in [V_h]^d$ is the lifting of the jump of $V$ defined by

$$\int_\Omega r_e([\![V]\!]_e) : \tau = -\int_e [\![V]\!]_e : \{\{\tau\}\}_e \quad \text{for all } \tau \in [V_h]^d. \tag{11}$$

For the numerical experiments in this paper we use $\chi_e = \frac{1}{2}\mathscr{N}_{\mathscr{T}_h}$, where $\mathscr{N}_{\mathscr{T}_h}$ is the maximal number of intersections one element in the grid can have (cf. [3]). We use triangular elements where $\chi_e = 1.5$ for all $e \in \Gamma$, and tetrahedral elements where $\chi_e = 2$ for all $e \in \Gamma$.

## 3.2 Temporal Discretization

The discrete solution $U_h(t) \in V_h$ has the form $U_h(t, x) = \sum_i U_i(t)\varphi_i(x)$. We get a system of ODEs for the coefficients of $U(t)$ which reads

$$U'(t) = f(U(t), t) \text{ in } (0, T] \tag{12}$$

with $f(U(t), t) = M^{-1}\mathscr{L}_h(U_h(t), t)$, $M$ being the mass matrix which is in our case block diagonal or even diagonal, depending on the choice of basis functions. $U(0)$ is given by the projection of $U_0$ onto $V_h$.

For the numerical results we have chosen Diagonally Implicit Runge-Kutta (DIRK) solvers of order 2, 3, or 4 depending on the polynomial order of the basis functions. The DIRK solvers are based on a Jacobian-free Newton-Krylov method (see [10]). The Krylov method is chosen to be GMRES without preconditioner. The implicit solver relies on a **matrix-free** implementation of the discrete operator $\mathscr{L}_h$. In a follow-up paper we will compare this approach to a fully assembled approach.

## 4 Numerical Results

In this section we present some benchmark tests for 2D and 3D focusing on parallelization and higher order DG schemes. Due to the lack of an exact solution $U$ we have computed the $L^2$-error between the discrete solution $U_h$ and a very fine, higher order solution $U_{h'}$. The quadrature order to compute $\|U_h - U_{h'}\|_{L^2(\Omega)}$ was chosen to be $2k + 4$, where $k$ denotes the order of the scheme. All computations are done on an unstructured, tetrahedral mesh.

**Table 1** Accuracy of the CDG2 scheme with 32 threads

| Level | Grid size | Linear | | | Quadratic | | | Cubic | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Time[a] | $L^2$-error | EOC[b] | Time[a] | $L^2$-error | EOC[b] | Time[a] | $L^2$-error | EOC[b] |
| 0 | 80 | 5.72E-1 | 2.42E-3 | – | 2.00E0 | 2.18E-3 | – | 6.52E0 | 1.96E-3 | – |
| 1 | 320 | 5.56E0 | 2.10E-3 | 0.20311 | 2.33E1 | 1.82E-3 | 0.26650 | 8.63E1 | 1.50E-3 | 0.38074 |
| 2 | 1280 | 3.98E2 | 1.82E-3 | 0.21263 | 2.09E2 | 1.34E-3 | 0.43315 | 8.22E2 | 9.26E-4 | 0.69823 |
| 3 | 5120 | 3.33E3 | 1.39E-3 | 0.38944 | 2.21E3 | 7.92E-4 | 0.76429 | 9.12E3 | 4.32E-4 | 1.0993 |
| 4 | 20480 | 3.01E4 | 8.28E-4 | 0.74208 | 2.10E4 | 2.94E-4 | 1.4284 | 8.02E4 | 8.77E-5 | 2.3024 |
| 5 | 81920 | 2.67E5 | 3.21E-4 | 1.3659 | 1.93E5 | 7.26E-5 | 2.0193 | 6.96E5 | 2.33E-5 | 1.9122 |

[a] Total CPU time

[b] Experimental order of convergence

## 4.1 A 2D Numerical Experiment with Six Species

$U_{h'}$ was calculated using the 4th order CDG2 scheme on a grid with 81,920 elements (refinement level 5), i.e. 7,372,800 degrees of freedom. For each $h$-refinement of the grid we bisect the time step size. Results for linear, quadratic and cubic DG schemes can be seen in Table 1. In Fig. 2 (left picture) we compare on a log-log scale the total CPU time of all threads with the $L^2$-error. Although the convergence rate is not as high as from the theory for parabolic problems, we see better rates for higher order schemes. We assume that re-entrant corners are responsible for the reduced convergence rates, see re-entrant corners in left picture of Fig. 1.

The right picture of Fig. 2 shows that the CDG2 is as good as the BR2 scheme and outperforms other DG schemes. A solution where SMCs start migrating into the domain can be found in Fig. 1, right picture.

## 4.2 A 3D Numerical Experiment with Three Species

Despite the fact that the benchmark in Sect. 4.1 can be accomplished with ease in 3D for six species, we want to simplify our model to three species $U = (u_1, u_3, u_4)$:

$$\mathscr{F}(U) = \begin{pmatrix} \chi_{14}\nabla u_4 \\ 0 \\ 0 \end{pmatrix}, \quad S(U) = \begin{pmatrix} 0 \\ F_0 u_1 u_3 \\ -\alpha_1 u_1 u_4 + \gamma u_3 \end{pmatrix}, \tag{13}$$

$$\mathscr{A}(U) = \mathrm{diag}(\mu_1, \mu_3, \mu_4). \tag{14}$$

We cannot trigger the inflammation through an inflow of LDL anymore. Thus, we suppose that the inflammation is triggered by a local, high concentration of debris and keep all other boundary and initial data from the Sect. 4.1.

In the 3D benchmark we examine parallelization using MPI and present in Table 2 strong scaling results for a third order CDG2 scheme on a grid with 113,549 elements

**Fig. 1** *Left* The coarsest grid for the EOC calculations containing 80 elements visualising a re-entrant corner. The angle of 171° stays fixed for all refinements. The considered domain is the cross-section of an arterial wall. *Middle* Initial distribution for the immune cells. We suppose that immune cells are more likely near the inner boundary. *Right* Solution for six species from *left to right, up to down*: Immune cells, SMCs, debris, chemoattractant, native and oxidized LDL. Native LDL has entered the domain through the endothelial layer where it has oxidized. The immune cells have accumulated (hard to recognize due to other effects) around the oxidized LDL and started absorbing them. Dying immune cells have built a necrotic core (debris) which is producing chemoattractant. The increase of chemoattractant has triggered a massive inflow of immune cells from the blood through the endothelial layer. A developed plaque has led SMCs to enter the arterial wall from the outer regions of the artery. To see the development of a fibrous cap (mainly formed by SMCs between the inner boundary and the necrotic core) the mathematical model has to be extended. (Data visualisation: Paraview)



**Fig. 2** Plot CPU time versus $L^2$-error (*Left*). 1st, 2nd and 3rd order CDG2 scheme (*Right*). 1st order CDG, CDG2, Baumann-Oden (BO), Bassy-Rebay (BR2), interior penalty (IP) scheme. (Visualisation of graphs: gnuplot)

and 13,625,880 degrees of freedom. Figure 3 shows the distribution of the processors and a discrete solution of the chemoattractant calculated using first order CDG2.

# 5 Conclusion

We have shown that DG schemes are well suited for solving huge coupled reactive diffusion transport systems. Modern techniques, such as parallelization, help to han-

**Table 2** CPU time for parallel runs using the cubic CDG2 method for computation of 10 time steps

| Processors | 8 | 16 | 32 | 64 | 128 | 256 |
|---|---|---|---|---|---|---|
| CPU time in s | 1177 | 528 | 277 | 142 | 75 | 39 |
| Speedup | – | 2.23 | 4.29 | 8.29 | 15.7 | 30.18 |



**Fig. 3** A 3D cuff model (*Left*). Each colour denotes a processor in a parallel run with 32 processors (*Right*). Isolines of the distribution of the chemoattractant after the inflammation has started

dle large systems in an appropriate CPU time. Furthermore, we have shown that it is possible to model the early stages of atherosclerotic plaque formation. A lot of more work needs to be done: In a future paper we will model the wall shear stress and some more species to understand later stages of atherosclerosis.

# References

1. Arnold, D.N., Brezzi, F., Cockburn, B., Marini, L.D.: Unified analysis of discontinuous galerkin methods for elliptic problems. SIAM J. Numer. Anal. **39**(5), 1749–1779 (2002)
2. Bastian, P., Blatt, M., Dedner, A., Engwer, C., Klöfkorn, R., Kornhuber, R., Ohlberger, M., Sander, O.: A generic grid interface for parallel and adaptive scientific computing. part ii: implementation and tests in dune. Computing **82**(2–3), 121–138 (2008)
3. Brdar, S., Dedner, A., Klöfkorn, R.: Compact and stable discontinuous Galerkin methods for convection-diffusion problems. SIAM J. Sci. Comput. **34**(1), 263–282 (2012). doi: http://dx.doi.org/10.1137/100817528
4. Calvez, V., Houot, J., Meunier, N., Raoult, A., Rusnakova, G., et al.: Mathematical and numerical modeling of early atherosclerotic lesions. ESAIM Proc. **30**, 1–14 (2010)
5. Dedner, A., Klöfkorn, R.: A generic stabilization approach for higher order discontinuous Galerkin methods for convection dominated problems. J. Sci. Comput. **47**(3), 365–388 (2011). doi: http://dx.doi.org/10.1007/s10915-010-9448-0

6. Dedner, A., Klöfkorn, R., Nolte, M., Ohlberger, M.: A generic interface for parallel and adaptive discretization schemes: abstraction principles and the dune-fem module. Computing **90**(3–4), 165–196 (2010)
7. Ibragimov, A., McNeal, C., Ritter, L., Walton, J.: A mathematical model of atherogenesis as an inflammatory response. Math. Med. Biol. **22**(4), 305–333 (2005)
8. Klöfkorn, R.: Benchmark 3d: the compact discontinuous Galerkin 2 scheme. In: Jaroslav (ed.) et al., Finite Volumes for Complex Applications VI Problems & Perspectives, pp. 1023–1033. Springer, Heidelberg (2011)
9. Klöfkorn, R.: Efficient matrix-free implementation of discontinuous Galerkin methods for compressible flow problems. In: Handlovičová, A. et al. (eds.) Algoritmy 2012. 19th Conference on Scientific Computing, Vysoké Tatry, Podbanské, Slovakia, 9–14 Sept 2012. Proceedings of Contributed Papers and Posters. Bratislava: Slovak University of Technology, Faculty of Civil Engineering, Department of Mathematics and Descriptive Geometry, 11–21 (2012). ISBN 978-80-227-3742-5/pbk.
10. Knoll, D.A., Keyes, D.E.: Jacobian-free newton-krylov methods: a survey of approaches and applications. J. Comput. Phys. **193**(2), 357–397 (2004). doi: http://dx.doi.org/10.1016/j.jcp.2003.08.010
11. Kuhlmann, M.T., Cuhlmann, S., Hoppe, I., Krams, R., Evans, P.C., Strijkers, G.J., et al.: Implantation of a carotid cuff for triggering shear-stress induced atherosclerosis in mice. J. Visualized Exp. **59**(e3308), 1–6 (2012)

# A DDFV Scheme for Incompressible Navier-Stokes Equations with Variable Density

**Thierry Goudon and Stella Krell**

**Abstract** We consider the application of "Discrete Duality Finite Volume" methods for the simulation of incompressible heterogeneous viscous flows. We pay attention to the numerical coupling between the mass conservation and the momentum balance equations, together with the divergence free constraint.

## 1 Introduction

This work is concerned with the numerical simulation of the Incompressible Navier–Stokes system

$$\begin{cases} \partial_t(\rho\mathbf{u}) + \mathrm{div}(\rho\mathbf{u}\otimes\mathbf{u}) + \mathrm{div}\left(-2\eta D\mathbf{u} + p\mathrm{Id}\right) = \mathbf{f}, \ \text{in } ]0,T[\times\Omega, \\ \qquad\qquad\qquad\qquad\qquad\qquad \mathrm{div}(\mathbf{u}) = 0, \ \text{in } ]0,T[\times\Omega, \qquad (1) \\ \qquad\qquad\qquad\qquad \partial_t\rho + \mathrm{div}(\rho\mathbf{u}) = 0, \ \text{in } ]0,T[\times\Omega, \end{cases}$$

where the unknowns are the velocity $\mathbf{u}$ :$]0,T[\times\Omega \to \mathbb{R}^2$, the density $\rho$ :$]0,T[\times\Omega \to [0,\infty)$ and the pressure $p$ :$]0,T[\times\Omega \to \mathbb{R}$. Here and below, $\Omega$ is a polygonal open bounded connected subset of $\mathbb{R}^2$, and $T > 0$ is fixed once for all. We denote

T. Goudon · S. Krell (✉)
INRIA, Team COFFEE and Labo. J. A. Dieudonné,
University Nice Sophia Antipolis-CNRS UMR 7351, Nice, France
e-mail: krell@unice.fr

T. Goudon
e-mail: thierry.goudon@inria.fr

D$\mathbf{u} = \frac{1}{2}(\nabla \mathbf{u} + {}^t \nabla \mathbf{u})$ and $[\text{div}(\rho \mathbf{u} \otimes \mathbf{u})]_j = \sum_{i=1}^{2} \partial_i (\rho \mathbf{u}_i \mathbf{u}_j)$ for $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2)$. We supplement the system (1) with the following boundary and initial conditions:

$$\mathbf{u} = \mathbf{g} \text{ on } ]0, T[\times \partial \Omega, \qquad \mathbf{u}(0, .) = \mathbf{u}_{\text{init}} \text{ in } \Omega, \qquad \rho(0, \cdot) = \rho_{\text{init}} \text{ in } \Omega$$

while the pressure is subjected to the condition $\int_\Omega p(t, x) \mathrm{d}x = 0$, for all $t \in \,]0, T[$. When the velocity points inward, we also prescribe the incoming density $\rho_{\text{inc}}$. We assume $\rho_{\text{init}} \in L^\infty(\Omega)$, $\mathbf{u}_{\text{init}} \in (L^\infty(\Omega))^2$, $\mathbf{f} \in (L^2(]0, T[\times \Omega))^2$, $\mathbf{g} \in L^2(]0, T[\times \partial \Omega)$. For the sake of simplicity, we assume that the viscosity $\eta$ is a positive constant.

There is a huge literature on the specific case where $\rho_{\text{init}}(x) = \bar{\rho} > 0$ is a given positive constant: owing to the incompressibility constraint, the density remains constant for ever. For instance, Finite Volume schemes have been recently introduced for the homogeneous Incompressible Navier–Stokes system [5, 7]. However, the situation of heterogeneous flows is much more realistic, and it leads to many difficulties both for the analysis and for numerics, see [3] and the references therein. Dedicated schemes have been introduced and analysed, based either on Finite Difference discretizations, or Finite Element discretizations. However, it is not clear that such methods preserve crucial physical properties like homogeneous solutions, the conservation of the total mass, the positivity of the density. These issues become particularly challenging when the ratio of extreme densities or the Reynolds number increase, or when dealing with highly unstructured meshes. In [3] an original method based on an hybrid Finite Volume/Finite Element approach is developed in order to cope with these difficulties.

In this paper, we address the problem of simulating the system (1) in the framework of the so-called DDFV schemes. These methods have been introduced in [6, 9] to approximate the solution of the Laplace equation on a large class of 2D meshes including non-conformal and distorted meshes. In particular, the scheme does not require "orthogonality" constraints on the mesh, by contrast to classical finite volume methods. Therefore, the method is very appealing in order to handle complex geometries, or to be used in combination to mesh refinements methods in order to follow accurately regions of strong density gradients. The strategy has been extended to a wide class of PDE problems, see e.g. [1, 4, 10], including in higher space dimension [12]. For Navier–Stokes equations, see [10, 11], DDFV schemes provide naturally a staggered discretization: the approximate velocity is stored at the centers and at the vertices of the mesh and the approximate pressure at the edges of the mesh. It turns out that the method can be extended to heterogeneous flows too, with a natural discretization of the density on the edges, the mass conservation being treated by UpWinding techniques. This work is a first attempt in this direction, where we detail how to handle the difficulties induced by the coupling. This approach looks particularly appealing for further extension towards intricate models for mixture flows, with a complex constraint relating the divergence of the velocity field and derivatives of the densities [8].

This paper is organized as follows. In Sect. 2, we detail the construction of the scheme, the numerical issues related to the coupling, and we state a few properties of the scheme. Finally, in Sect. 3, we discuss a few numerical results.

## 2 The DDFV Framework

We refer the reader to [10] for a description of the DDFV scheme for the Stokes problem; with the same notation, we have:

- the DDFV meshes $(\mathfrak{T}, \mathfrak{D})$: $\mathfrak{T}$ combines the primal mesh $\mathfrak{M} \cup \partial\mathfrak{M}$ (whose cells are denoted by $\kappa$), and the dual mesh $\mathfrak{M}^* \cup \partial\mathfrak{M}^*$, (whose cells $\kappa^*$ are built around the vertices $x_{\kappa^*}$ of the primal mesh), see Fig. 1. Next, $\mathfrak{D}$ stands for the diamond mesh, whose cells $D$ are built around the edges $\sigma$ of the primal mesh. For $\phi$ defined on $\mathfrak{D}$ (i.e $\phi_\mathfrak{D} \in \mathbb{R}^\mathfrak{D}$), we denote by $\phi_D$ its value on $D \in \mathfrak{D}$. We use a similar notation for quantities defined on $\mathfrak{T}$, e. g., $\mathbf{u}_\kappa$ (resp. $\mathbf{u}_{\kappa^*}$) for the value on $\kappa$ (resp. $\kappa^*$) of a quantity $\mathbf{u}$ defined on the primal and dual mesh, i.e $\mathbf{u}^\mathfrak{T} \in \left(\mathbb{R}^2\right)^\mathfrak{T}$.
- We define: the discrete gradient $\nabla^\mathfrak{D} : \left(\mathbb{R}^2\right)^\mathfrak{T} \to (\mathrm{M}_2(\mathbb{R}))^\mathfrak{D}$,

$$\nabla^D \mathbf{u}^\mathfrak{T} := \frac{1}{2|D|}\Big[|\sigma|(\mathbf{u}_L - \mathbf{u}_K) \otimes \mathbf{n}_{\sigma K} + |\sigma^*|(\mathbf{u}_{L^*} - \mathbf{u}_{K^*}) \otimes \mathbf{n}_{\sigma^* K^*}\Big], \quad \forall D \in \mathfrak{D},$$

with $\mathbf{n}_{\sigma K}$ the unit vector normal to $\sigma$ oriented from $x_K$ to $x_L$, and $\mathbf{n}_{\sigma^* K^*}$ the unit vector normal to $\sigma^*$ oriented from $x_{K^*}$ to $x_{L^*}$; its discrete dual operator $\mathbf{div}^\mathfrak{T}$ : $(\mathrm{M}_2(\mathbb{R}))^\mathfrak{D} \to \left(\mathbb{R}^2\right)^\mathfrak{T}$, for any $\kappa \in \mathfrak{M}$ and $\kappa^* \in \mathfrak{M}^*$

$$\mathbf{div}^K \xi_\mathfrak{D} := \frac{1}{|K|} \sum_{D \in \mathfrak{D}, D \subset K} |\sigma| \xi_D \mathbf{n}_{\sigma K}, \qquad \mathbf{div}^{K^*} \xi_\mathfrak{D} := \frac{1}{|K^*|} \sum_{D \in \mathfrak{D}, D \subset K^*} |\sigma^*| \xi_D \mathbf{n}_{\sigma^* K^*};$$

its trace $\mathrm{div}^\mathfrak{D} : \left(\mathbb{R}^2\right)^\mathfrak{T} \to \mathbb{R}^\mathfrak{D}$ that is $\mathrm{div}^\mathfrak{D}(\mathbf{u}^\mathfrak{T}) = \mathrm{Tr}(\nabla^\mathfrak{D}\mathbf{u}^\mathfrak{T})$; a discrete strain rate tensor $D^\mathfrak{D} : \left(\mathbb{R}^2\right)^\mathfrak{T} \to (\mathrm{M}_2(\mathbb{R}))^\mathfrak{D}$ that is $D^\mathfrak{D}(\mathbf{u}^\mathfrak{T}) = \frac{1}{2}(\nabla^\mathfrak{D}\mathbf{u}^\mathfrak{T} + {}^t\nabla^\mathfrak{D}\mathbf{u}^\mathfrak{T})$.

As in [3], the system (1) is treated with a time splitting strategy. What is crucial in the construction of the scheme is to verify the compatibility between the discretization of the convection terms in the mass and momentum equations.

## 2.1 Approximation of the Mass Conservation Equation

The mass conservation equation is seen as a transport equation for the density $\rho$ with velocity $\mathbf{u}$. Hence, we set up an approximation based on UpWinding principles.

**Fig. 1** The mesh $\mathfrak{T}$ (*left*). A diamond $D$ with a neighbour diamond $D'$ (*right*)

**Definition 1** *The discrete divergence operator* $\mathrm{divc}^{\mathfrak{D}}$ *is a mapping from* $\mathbb{R}^{\mathfrak{D}} \times \left(\mathbb{R}^2\right)^{\mathfrak{T}}$ *to* $\mathbb{R}^{\mathfrak{D}}$ *defined for all* $\rho_{\mathfrak{D}} \in \mathbb{R}^{\mathfrak{D}}$ *and* $\mathbf{u}^{\mathfrak{T}} \in \left(\mathbb{R}^2\right)^{\mathfrak{T}}$ *by* $\mathrm{divc}^{\mathfrak{D}}(\rho_{\mathfrak{D}}, \mathbf{u}^{\mathfrak{T}}) = \left(\mathrm{divc}^{D}(\rho_{\mathfrak{D}}, \mathbf{u}^{\mathfrak{T}})\right)_{D \in \mathfrak{D}}$, *with*

$$\mathrm{divc}^{D}(\rho_{\mathfrak{D}}, \mathbf{u}^{\mathfrak{T}}) = \frac{1}{|D|} \sum_{\mathfrak{s} = D|D' \in \partial D} |\mathfrak{s}| \left(\left(g_{\mathfrak{s},D}\right)^{+} \rho_D - \left(g_{\mathfrak{s},D}\right)^{-} \rho_{D'}\right), \ \forall D \in \mathfrak{D}$$

*where* $g_{\mathfrak{s},D} = \frac{\mathbf{u}_K + \mathbf{u}_{K^*}}{2} \cdot \mathbf{n}_{\mathfrak{s}D}$ *for* $\mathfrak{s} = [K, x_{K^*}] \in \partial D$, $\mathbf{n}_{\mathfrak{s}D}$ *the unit normal to* $\mathfrak{s}$ *outward of* $D$, $x^{+} = \max(x, 0)$ *and* $x^{-} = -\min(x, 0)$ *for all* $x \in \mathbb{R}$. *When* $\mathfrak{s} \subset \partial \Omega$ *and* $g_{\mathfrak{s},D} < 0$, *we use the prescribed incoming density:* $g_{\mathfrak{s},D}\rho_{\mathrm{inc}}(x_D)$.

Note that $\mathrm{div}^{\mathfrak{D}}(\mathbf{u}^{\mathfrak{T}}) = \mathrm{divc}^{\mathfrak{D}}(\mathbf{1}_{\mathfrak{D}}, \mathbf{u}^{\mathfrak{T}})$, with $\mathbf{1}_{D} = 1$ for all $D \in \mathfrak{D}$: homogeneous solutions are preserved. Let $N \in \mathbb{N}^*$. We note $\delta t = \frac{T}{N}$ and $t_n = n\delta t$ for $n \in \{0, \dots, N\}$. Having at hand the discrete unknowns $\rho_{\mathfrak{D}}^n$ and $\mathbf{u}_{\mathfrak{T}}^n$ at time $t_n$, we first update the density by

$$\frac{\rho_{\mathfrak{D}}^{n+1} - \rho_{\mathfrak{D}}^n}{\delta t} + \mathrm{divc}^{\mathfrak{D}}(\rho_{\mathfrak{D}}^n, \mathbf{u}_{\mathfrak{T}}^n) = 0. \tag{2}$$

It defines $\rho_{\mathfrak{D}}^{n+1} \in \mathbb{R}^{\mathfrak{D}}$. We define a density $\rho_{\mathfrak{T}}^{n+1}$ on the mesh $\mathfrak{T}$ by

$$\rho_K^{n+1} = \frac{1}{|K|} \sum_{D \in \mathfrak{D}, D \subset K} |D \cap K| \rho_D^{n+1}, \quad \rho_{K^*}^{n+1} = \frac{1}{|K^*|} \sum_{D \in \mathfrak{D}, D \subset K^*} |D \cap K^*| \rho_D^{n+1}, \ \forall K \in \mathfrak{M}, K^* \in \mathfrak{M}^*.$$

## 2.2 Approximation of the Non-linear Term $\rho \mathbf{u} \otimes \mathbf{u}$

We remind that $\mathbf{u}^{n+1}$ is defined on the mesh $\mathfrak{T}$; thus the divergence $\rho^n \mathbf{u}^n \otimes \mathbf{u}^{n+1}$ would be naturally defined on $\mathfrak{D}$. This is not compatible with the discretization of the momentum equation which has to be considered on $\mathfrak{T}$. We obtain a meaningful

discretization by going back to the Stokes formula: for $V \in \mathfrak{T}$, $\sum_{\varsigma \in \partial V} \int_{\varsigma} (\rho^n \mathbf{u}^n \cdot \mathbf{n}_{\varsigma,V}) \mathbf{u}^{n+1} ds$ is approached by a formula which looks like $\sum_{\varsigma \in \partial V} F_{\varsigma,V} \mathbf{u}_{\varsigma}^{n+1}$. It remains to explain how to define $F_{\varsigma,V}$ and $\mathbf{u}_{\varsigma}^{n+1}$.

**Definition 2** *We define* $b_m^{\mathfrak{T}} : (F_{\varsigma,V}, \mathbf{u}^{\mathfrak{T}}) \in \mathbb{R}^{\mathfrak{D}} \times (\mathbb{R}^2)^{\mathfrak{T}} \mapsto b_m^{\mathfrak{T}}(F_{\varsigma,V}, \mathbf{u}^{\mathfrak{T}}) \in (\mathbb{R}^2)^{\mathfrak{T}}$, *as follows: for any* $\kappa \in \mathfrak{M}$, $\kappa^* \in \mathfrak{M}^*$,

$$b_{\kappa}^m (F_{\varsigma,V}, \mathbf{u}^{\mathfrak{T}}) = \frac{1}{|\kappa|} \sum_{\sigma \in \partial \kappa} (F_{\sigma,\kappa})^+ \mathbf{u}_{\kappa} - (F_{\sigma,\kappa})^- \mathbf{u}_{\mathsf{L}}$$

$$b_{\kappa^*}^m (F_{\varsigma,V}, \mathbf{u}^{\mathfrak{T}}) = \frac{1}{|\kappa^*|} \sum_{\sigma^* \in \partial \kappa^*} (F_{\sigma^*,\kappa^*})^+ \mathbf{u}_{\kappa^*} - (F_{\sigma^*,\kappa^*})^- \mathbf{u}_{\mathsf{L}^*}.$$

Again, this definition relies on some UpWinding principle. We turn to the definition of $F_{\sigma,\kappa}$. We set $F_{\mathfrak{s},\mathsf{D}}^n = |\mathfrak{s}| \left( (g_{\mathfrak{s},\mathsf{D}}^n)^+ \rho_{\mathsf{D}}^n - (g_{\mathfrak{s},\mathsf{D}}^n)^- \rho_{\mathsf{D}'}^n \right)$. As a matter of fact, (2) recasts as

$$|\mathsf{D}| \frac{\rho_{\mathfrak{D}}^{n+1} - \rho_{\mathfrak{D}}^n}{\delta t} + \sum_{\mathfrak{s} \in \partial \mathsf{D}} F_{\mathfrak{s},\mathsf{D}}^n = 0.$$

We wish to establish a similar conservation relation on the primal cells $\kappa \in \mathfrak{M}$

$$|\kappa| \frac{\rho_{\kappa}^{n+1} - \rho_{\kappa}^n}{\delta t} + \sum_{\sigma \in \partial \kappa} F_{\sigma,\kappa}^n = 0. \tag{3}$$

This requirement guides the construction of $F_{\sigma,\kappa}^n$. To this end, we seek four vectors $W_{\kappa}, W_{\kappa^*}, W_{\mathsf{L}}, W_{\mathsf{L}^*}$ and a function $w \in H^1(\Omega)$ that fulfill

(i) for all $\mathfrak{s} = [\kappa, x_{\kappa^*}] \in \partial \mathsf{D}$, we have $\dfrac{W_{\kappa} + W_{\kappa^*}}{2} \cdot \mathbf{n}_{\mathfrak{s}\mathsf{D}} = \dfrac{F_{\mathfrak{s},\mathsf{D}}^n}{|\mathfrak{s}|}$.

(ii) $w_{|\mathsf{D}}$ is piecewise $\mathbb{P}^1$ on each quarter diamond $\mathscr{Q}$ so that $\mathrm{div}(w_{|\mathsf{D}})$ is constant on $\mathsf{D}$; This constant is imposed to be $\mathrm{div}(w_{|\mathsf{D}}) = \frac{1}{|\mathsf{D}|} \sum_{\mathfrak{s} \in \partial \mathsf{D}} F_{\mathfrak{s},\mathsf{D}}^n$.

Up to a suitable labelling of the vertices of a diamond $\mathsf{D}$, see Fig. 2, we define a function $\Phi : W \in \mathbb{R}^8 \mapsto \Phi(W) \in \mathbb{R}^4$, where the components of $\Phi(W)$ correspond to the inner product $\mathbf{n}_{ij} \cdot (W_i + W_j)/2$. This function is surjective, which proves the existence of $W_{\kappa}, W_{\kappa^*}, W_{\mathsf{L}}, W_{\mathsf{L}^*}$ satistying condition (i). The construction of $w$ on each diamond $\mathsf{D}$ relies on the Nagtegaal device for finite elements methods [13]. On each quarter diamond $\mathscr{Q}_{ij}$ (see Fig. 2), we obtain a $\mathbb{P}^1$ function required to satisfy $w_{|\mathsf{D}} \left( \frac{P_i + P_j}{2} \right) = \frac{W_i + W_j}{2}$ and $\mathrm{div}(w_{|\mathsf{D}})$ is a constant. We find this constant owing to Stokes' formula

**Fig. 2** A suitable labelling of the vertices of a diamond $\mathrm{D}$ and of the quarter diamond $\mathscr{Q}$



$$\int_D \mathrm{div}(w_{|\mathrm{D}}) = \sum_{\mathfrak{s}\in\partial\mathrm{D}} \int_{\mathfrak{s}} w_{|\mathrm{D}} \cdot \mathbf{n}_{\mathfrak{s}\mathrm{D}} = \sum_{\mathfrak{s}\in\partial\mathrm{D}} |\mathfrak{s}| w_{|\mathrm{D}} \left( \frac{P_i + P_j}{2} \right) \cdot \mathbf{n}_{ij}$$

$$= \sum_{\mathfrak{s}\in\partial\mathrm{D}} |\mathfrak{s}| \frac{W_i + W_j}{2} \cdot \mathbf{n}_{ij} = \sum_{\mathfrak{s}\in\partial\mathrm{D}} F^n_{\mathfrak{s},\mathrm{D}}.$$

We conclude that $\mathrm{div}(w_{|\mathrm{D}}) = \dfrac{1}{|\mathrm{D}|} \displaystyle\sum_{\mathfrak{s}\in\partial\mathrm{D}} F^n_{\mathfrak{s},\mathrm{D}}$ holds, which proves (ii). Finally, we set

$F^n_{\sigma,\mathrm{K}} = \displaystyle\int_\sigma w_{|\mathrm{D}} \cdot \mathbf{n}_{\sigma\mathrm{K}}$, which actually leads to the following explicit formula:

$$F^n_{\sigma,\mathrm{K}} = -\frac{|\mathrm{D}\cap\mathrm{L}|}{|\mathrm{D}|} \sum_{\mathfrak{s}\in\partial\mathrm{D},\,\mathfrak{s}\subset\mathrm{K}} F^n_{\mathfrak{s},\mathrm{D}} + \frac{|\mathrm{D}\cap\mathrm{K}|}{|\mathrm{D}|} \sum_{\mathfrak{s}\in\partial\mathrm{D},\,\mathfrak{s}\subset\mathrm{L}} F^n_{\mathfrak{s},\mathrm{D}}.$$

We end up with the mass balance Eq. (3) for $\rho_\mathrm{K}^{n+1}$. The fluxes are conservative since $F^n_{\sigma,\mathrm{K}} = -F^n_{\sigma,\mathrm{L}}$ for $\sigma = \mathrm{K}|\mathrm{L}$. A similar construction can be made for $F^n_{\sigma^*,\mathrm{K}^*}$.

## *2.3 DDFV Schemes for the Navier-Stokes Equation*

The scheme for the momentum equation reads (up to the boundary condition):

$$\begin{cases} \text{Find } \mathbf{u}_\mathfrak{T}^{n+1} \in \left(\mathbb{R}^2\right)^\mathfrak{T} \text{ and } p_\mathfrak{D}^{n+1} \in \mathbb{R}^\mathfrak{D} \text{ such that,} \\[4pt] \frac{\rho_\mathfrak{T}^{n+1}}{\delta t}\mathbf{u}_\mathfrak{T}^{n+1} + \mathbf{div}^\mathfrak{T}(-2\eta\mathrm{D}^\mathfrak{D}\mathbf{u}_\mathfrak{T}^{n+1} + p_\mathfrak{D}^{n+1}\mathrm{Id}) + b_m^\mathfrak{T}(F^n_{\varsigma,V},\mathbf{u}_\mathfrak{T}^{n+1}) = \frac{\rho_\mathfrak{T}^n}{\delta t}\mathbf{u}_\mathfrak{T}^n + \mathbf{f}_\mathfrak{T}^{n+1}, \\[4pt] \mathrm{div}^\mathfrak{D}(\mathbf{u}_\mathfrak{T}^{n+1}) = 0, \qquad \sum_{\mathrm{D}\in\mathfrak{D}} |\mathrm{D}| p_\mathrm{D}^{n+1} = 0. \end{cases}$$

(4)

We refer the reader to [10, 11] for the treatment of the boundary condition. Here, we considered meshes such that we do not need further stabilization terms [2].

**Proposition 1** *The finite volume scheme* (4) *admits a unique solution* $(\mathbf{u}_\mathfrak{T}^{n+1}, p_\mathfrak{D}^{n+1})$. *Going back to the complete problem, if* $\rho_\mathfrak{D}^n = 1$, *then* $\rho_\mathfrak{D}^{n+1} = 1$.

**Fig. 3** Family of meshes. On the *left*: non conformal square mesh; on the *right*: Triangle mesh



**Table 1** Test case 1 on the non conformal square mesh Fig. 3

| NbCell | Ervel | Ratio | Ergradvel | Ratio | Erpre | Ratio |
|--------|-------|-------|-----------|-------|-------|-------|
| 64 | 1.534E-01 | – | 1.662E-01 | – | 60.029 | – |
| 208 | 2.723E-02 | 2.49 | 8.391E-02 | 0.99 | 19.21 | 1.65 |
| 736 | 6.577E-03 | 2.05 | 4.240E-02 | 0.99 | 7.862 | 1.29 |
| 2752 | 1.789E-03 | 1.88 | 2.123E-02 | 1.00 | 3.797 | 1.05 |
| 10624 | 6.434E-04 | 1.48 | 1.061E-02 | 1.00 | 1.900 | 1.00 |

## 3 Numerical Results

We validate the scheme by showing a few numerical experiments, inspired from [3]. The computational domain is $\Omega =]0, 1[^2$. We set $\eta = 1$. We wish to capture explicit solutions of (1) with convenient source term **f** and boundary data **g**. In order to discuss error estimates, a family of meshes is obtained by successive global refinement of the original mesh, see Fig. 3. We compare the relative $L^2(\Omega \times ]0, T[)$-norm of the error obtained with the DDFV scheme, for the pressure (denoted Erpre), the velocity gradient (denoted Ergradvel) and the velocity (denoted Ervel) respectively. On the two tables, we give the number of primal cells (denoted NbCell) and the convergence rates (denoted Ratio). The linear system associated to (4) is solved by a direct method, and $\mathrm{div}(u)$ vanishes at the machine-error order. We can check that the total mass is conserved.

**Test case 1.** The Green-Taylor vortex:

$$\mathbf{u} = \begin{pmatrix} -\cos(2\pi x)\sin(2\pi y)e^{-2t} \\ \sin(2\pi x)\cos(2\pi y)e^{-2t} \end{pmatrix}, \quad p = -\frac{1}{4}(\cos(4\pi x)+\cos(4\pi y))e^{-4t}, \quad \rho = 1.$$

The final time is $T = 1$ and we set $\delta t = 5 \times 10^{-3}$ that ensures the stability of the transport part of the algorithm for all meshes.

The homogeneous case $\rho = 1$ is perfectly preserved by the scheme (Table 1), we observe a first order accuracy on the velocity gradient and the pressure, which seems to be optimal. We obtain a super-convergence for the $L^2$-norm of the velocity.

**Table 2** Test case 2 on the triangle mesh Fig. 3

| NbCell | Ervel | Ratio | Ergradvel | Ratio | ErPre | Ratio | ErDen | Ratio |
|--------|-------|-------|-----------|-------|-------|-------|-------|-------|
| 72 | 1.41E-03 | – | 1.04E-02 | – | 2.82E-02 | – | 2.21E-03 | – |
| 256 | 5.38E-04 | 1.4 | 6.47E-03 | 0.7 | 7.61E-03 | 1.9 | 1.21E-03 | 0.9 |
| 960 | 1.95E-04 | 1.5 | 3.31E-03 | 1.0 | 2.58E-03 | 1.6 | 7.37E-04 | 0.7 |
| 3712 | 8.08E-05 | 1.3 | 1.60E-03 | 1.0 | 1.24E-03 | 1.1 | 4.38E-04 | 0.8 |
| 14592 | 3.75E-05 | 1.1 | 7.63E-04 | 1.0 | 8.35E-04 | 0.6 | 2.41E-04 | 0.9 |

Furthermore, we point out that the convergence rate is not sensitive to the presence of non conformal control volumes.

**Test case 2.** An example of non homogeneous flow:

$$\mathbf{u} = \begin{pmatrix} -y\cos(t) \\ x\cos(t) \end{pmatrix}, \quad p = \sin(x)\sin(y)\sin(t), \quad \rho(r,\theta,t) = 2 + r\cos(\theta - \sin(t)).$$

The final time is $T = 3{,}125 \times 10^{-2}$ and we set $\delta t = 7{,}8125 \times 10^{-5}$ that ensures the stability of the transport part of the algorithm for all meshes. In Table 2, we also give the error on the density (ErDen). Results are coherent with [3], which justifies the validity of the DDFV approach that will be extended to more ambitious situations elsewhere.

# References

1. Andreianov, B., Boyer, F., Hubert, F.: Discrete duality finite volume schemes for leray-lions type elliptic problems on general 2d-meshes. Numer. Methods PDE **23**(1), 145–195 (2007)
2. Boyer, F., Krell, S., Nabet, F.: Inf-Sup Stability of the Discrete Duality Finite Volume Method for the Stokes Problem, Preprint, Inria-CNRS-Univ. Nice (2014)
3. Calgaro, C., Creusé, E., Goudon, T.: An hybrid finite volume-finite element method for variable density incompressible flows. J. Comput. Phys. **227**(9), 4671–4696 (2008). http://math.univ-lille1.fr/simpaf/SITE-NS2DDV/home.html
4. Coudière, Y., Manzini, G.: The discrete duality finite volume method for convection-diffusion problems. SIAM J. Numer. Anal. **47**(6), 4163–4192 (2010)
5. Droniou, J., Eymard, R.: Study of the mixed finite volume method for stokes and navier-stokes equations. Num. Meth. PDEs **25**(1), 137–171 (2009)
6. Domelevo, K., Omnès, P.: A finite volume method for the laplace equation on almost arbitrary two-dimensional grids. Math. Model. Numer. Anal. **39**(6), 1203–1249 (2005)
7. Eymard, R., Herbin, R., Latché, J.-C.: Convergence analysis of a colocated finite volume scheme for the incompressible navier-stokes equations on general 2d or 3d meshes. SIAM J. Numer. Anal. **45**(1), 1–36 (2007)
8. Goudon, T., Vasseur, A.: On a Model for Mixture Flows: Derivation, Dissipation and Stability Properties. Preprint. Inria-CNRS-Univ. Nice (2014)
9. Hermeline, F.: A finite volume method for the approximation of diffusion operators on distorted meshes. J. Comput. Phys. **160**(2), 481–499 (2000)

10. Krell, S.: Stabilized DDFV schemes for Stokes problem with variable viscosity on general 2D meshes. Num. Meth. PDEs, (2011). http://dx.doi.org/10.1002/num.20603
11. Krell, S.: Stabilized DDFV schemes for the incompressible Navier-Stokes equations. In: Proceedings of FVCA6 (Praha), Springer Proceedings in Math, vol. 4, pp. 605–612. (2011)
12. Krell, S., Manzini, G.: The discrete duality finite volume method for the stokes equations on 3d polyhedral meshes. SIAM J. Numer. Anal. **50**(2), 808–837 (2012)
13. Nagtegaal, J.C., Parks, D.M., Rice, J..R.: On numerically accurate finite element solution in the fully plastic range. Comput. Meth. Appl. Mech Eng. **4**,153–177 (197**4**)

# An Efficient Implementation of a 3D CeVeFE DDFV Scheme on Cartesian Grids and an Application in Image Processing

**Niklas Hartung and Florence Hubert**

**Abstract** In this work we describe the implementation of a 3D Center-Vertex-Face/Edge Discrete Duality Finite Volume (CeVeFE DDFV) scheme using only the degrees of freedom (DOF) disposed on a Cartesian grid. These DOF are organised in a three-mesh structure proper to the CeVeFE DDFV setting. Reposing on a diamond structure, the approach presented here greatly simplifies the implementation, also in the case of grids topologically equivalent to the uniform Cartesian one. The numerical scheme is then applied to a problem in image processing, where uniform Cartesian structure of the DOF is naturally imposed by the pixel/voxel structure. A semi-implicit DDFV scheme is used for solving a nonlinear advection-diffusion equation, the subjective surfaces equation, in order to reconstruct the volume of a tumour from noisy 3D SPECT images with signal intensity on the tumour boundary. The matrix of the linear system has a band structure and the method is fast and able to successfully reconstruct the tumour volume.

## 1 Introduction

Discrete Duality Finite Volume (DDFV) schemes, introduced in 2D for the Laplace problem by Hermeline [10], are a possible discretisation strategy applying to very general meshes and a large variety of PDE [3, 6]. A dual or "node" mesh is used in this framework and gradients are defined on a structure called the diamond mesh. A main feature of the DDFV approach is that discrete gradients and divergence operators are defined in a way that a discrete Green formula holds, called "discrete duality".

N. Hartung (✉) · F. Hubert
CNRS, Centrale Marseille, I2M, UMR 7373, Aix Marseille Université, CMI 39 rue Frédéric Joliot-Curie, 13453 Marseille, Cedex 13, France
e-mail: niklas.hartung@univ-amu.fr

F. Hubert
e-mail: florence.hubert@univ-amu.fr

In 3D, several methods have been inspired by the 2D DDFV methodology. The Center-Vertex (CeVe) DDFV schemes have a dual mesh with centers at the vertices of the primary mesh [1, 2, 11]. A different method, called Center-Vertex-Face-Edge (CeVeFE) DDFV scheme, features a third mesh with unknowns at the faces and edges of the primary mesh [5]; this will be our framework.

The DDFV framework can also be used for the discretisation of PDE appearing in image processing, such as level set methods, which are used for a broad spectrum of applications [15]. A curvature-driven level set equation called the subjective surfaces equation has been introduced by Sarti et al. [14] as a tool for the completion of missing boundaries. Along with subsequent extensions, it has been successfully applied to image processing problems [12, 13].

The nonlinearity and the non-divergent form of the curvature-driven level set equation makes particular space discretisation techniques necessary. Several Finite Volume methods have been proposed along with numerical analysis [7, 13]. Recently, stability and convergence of a semi-implicit 2D DDFV scheme was proven [9], but additional vertex unknowns were introduced.

In this work we will detail the efficient implementation of the 3D CeVeFE DDFV method for Cartesian grids. The method will then be used for discretising the subjective surfaces equation on a uniform Cartesian grid. As an application, tumour volume is reconstructed from a 3D SPECT image visualising proliferating cells, which are located at the tumour boundary.

## 2 The 3D CeVeFE DDFV Scheme with Degrees of Freedom on Cartesian Grids

### 2.1 Construction of the Meshes

In the 3D CeVeFE DDFV scheme, three different decompositions of the computational domain $\Omega$ are used, called the primary mesh $\mathscr{M}$, the dual or node mesh $\mathscr{N}$ associated with the vertices of the primary mesh and the tertiary or "face-edge" mesh $\mathscr{F}\mathscr{E}$ associated with the faces and edges of the primary mesh. For the detailed construction of $\mathscr{N}$ and $\mathscr{F}\mathscr{E}$ from a general primary mesh $\mathscr{M}$, we refer to [5].

There is a canonical way to construct these three meshes if we want each gridpoint of a Cartesian structure $\mathscr{T}$ to be associated to exactly one cell of either $\mathscr{M}$, $\mathscr{N}$ or $\mathscr{F}\mathscr{E}$. Referring to each point by its three-dimensional index $(i, j, k)$, $1 \leqslant i \leqslant N_x \in \mathbb{N}$, $1 \leqslant j \leqslant N_y \in \mathbb{N}$, $1 \leqslant k \leqslant N_z \in \mathbb{N}$, we have the following bijections (see Fig. 1):

$$\{(i, j, k) \text{ with } i, j, k \text{ even}\} \mapsto \mathscr{M}, \qquad \{(i, j, k) \text{ with } i, j, k \text{ odd}\} \mapsto \mathscr{N},$$
$$\{(i, j, k) \text{ with } ijk \equiv 2 \mod 4\} \mapsto \mathscr{E}, \quad \{(i, j, k) \text{ with } ijk \equiv 4 \mod 8\} \mapsto \mathscr{F}.$$

We denote the control volumes of the primary mesh $\mathscr{M}$ by K or L (with centers $x_K$ or $x_L$), vertices by $x_A$ or $x_B$, edges by E and faces by F. Control volumes of the dual

**(a)**



**(b)**



**(c)**



**Fig. 1** 3D mesh views for uniform Cartesian grids. Primary mesh centers are marked by *circles*, nodes by *squares*, faces by *upright* and edges by *sideward triangles*. **a** The 3 meshes. **b** Two primary mesh cells together with other cells appearing in the mesh construction. *Extreme left—node cell, middle left—face cell, middle right—edge cell, extreme right—diamond cell*. **c** Mapping of a uniform Cartesian grid onto $\mathcal{M}$, $\mathcal{N}$ and $\mathcal{EF}$

mesh $\mathcal{N}$ will be called A or B. To simplify notations, control volumes of the tertiary mesh $\mathcal{FE}$ will also be called E and F as it will be clear from the context whether the face/edge or the control volume is meant. We also define the center of gravity $x_F$ of a face F and the midpoint $x_E$ of an edge E.

Discrete gradients are defined on a fourth decomposition of $\Omega$ called the diamond mesh $\mathcal{D}$. Each diamond cell $D \in \mathcal{D}$ corresponds to a face-edge couple (F, E) with $E \in \partial F$. In our decomposition of the Cartesian grid $\mathcal{T}$, a diamond cell D will be defined by listing the indices of six points of $\mathcal{T}$: a face center $x_F$, the midpoint $x_E$ of an edge $E \in \partial F$, the two vertices $x_A, x_B \in \partial E$ and the two adjacent primary cell centers $x_K$ and $x_L$. The diamond D is then given by $D := \text{hull}(x_A, x_B, x_L, x_K)$. The

diamond mesh can be subdivided into twelve classes of diamonds, which are listed in Table 1. This classification permits an efficient construction of the diamond mesh. The volume of D is given by $|\mathrm{D}| = \frac{\det(x_{\mathrm{B}} - x_{\mathrm{A}}, x_{\mathrm{F}} - x_{\mathrm{E}}, x_{\mathrm{L}} - x_{\mathrm{K}})}{6} > 0$. For $\mathrm{F} \in \partial\Omega$, the indices that exceed the Cartesian grid are projected onto $\mathscr{T}$, creating degenerate diamonds.

Noting $x_{\mathrm{D}} = \frac{1}{2}(x_{\mathrm{E}} + x_{\mathrm{F}})$, a diamond D can be decomposed into eight tetrahedra $\mathrm{D}_{\mathrm{AKE}}, \mathrm{D}_{\mathrm{ALE}}, \mathrm{D}_{\mathrm{BKE}}, \mathrm{D}_{\mathrm{BLE}}, \mathrm{D}_{\mathrm{AKF}}, \mathrm{D}_{\mathrm{ALF}}, \mathrm{D}_{\mathrm{BKF}}, \mathrm{D}_{\mathrm{BLF}}$ defined by

$$\text{hull}\left(\begin{pmatrix} x_{\mathrm{K}} \\ x_{\mathrm{L}} \end{pmatrix}, \begin{pmatrix} x_{\mathrm{A}} \\ x_{\mathrm{B}} \end{pmatrix}, \begin{pmatrix} x_{\mathrm{E}} \\ x_{\mathrm{F}} \end{pmatrix}, x_{\mathrm{D}}\right). \tag{1}$$

This decomposition permits to define the control volumes C of any of the three meshes as the union of all tetrahedra containing the vertex $x_{\mathrm{C}}$, e.g.

$$\mathrm{K} = \bigcup_{\mathrm{D} \in \mathscr{D}: x_{\mathrm{K}} \in \mathrm{D}} (\mathrm{D}_{\mathrm{AKE}} \cup \mathrm{D}_{\mathrm{BKE}} \cup \mathrm{D}_{\mathrm{AKF}} \cup \mathrm{D}_{\mathrm{BKF}}).$$

With this definition, some boundary volumes, depending on the parity of $N_x$, $N_y$, $N_z$, degenerate automatically. The DDFV meshes are coarser than the canonical mesh associated to the Cartesian grid.

## 2.2 Discrete Gradient and Discrete Divergence Operators

For $u \in \mathbb{R}^{|\mathscr{T}|}$ and a diamond $\mathrm{D} \in \mathscr{D}$, set $u_{\mathrm{C}}$ as a notation for $u(x_{\mathrm{C}})$ where $x_{\mathrm{C}}$ is one of the six points defining D. The discrete gradient $\nabla^d : \mathbb{R}^{|\mathscr{T}|} \mapsto \mathbb{R}^{|\mathscr{D}|}$ is given by

$$\left(\nabla^d(u_{\mathscr{T}})\right)_{\mathrm{D}} = \frac{1}{3|\mathrm{D}|}\left((u_{\mathrm{L}} - u_{\mathrm{K}})\overrightarrow{N_{\mathrm{KL}}} + (u_{\mathrm{F}} - u_{\mathrm{E}})\overrightarrow{N_{\mathrm{EF}}} + (u_{\mathrm{B}} - u_{\mathrm{A}})\overrightarrow{N_{\mathrm{AB}}}\right)$$

for any $\mathrm{D} \in \mathscr{D}$ and with the vectors

$$\overrightarrow{N_{\mathrm{KL}}} = \frac{(x_{\mathrm{B}} - x_{\mathrm{A}}) \times (x_{\mathrm{F}} - x_{\mathrm{E}})}{2}, \ \overrightarrow{N_{\mathrm{AB}}} = \frac{(x_{\mathrm{F}} - x_{\mathrm{E}}) \times (x_{\mathrm{L}} - x_{\mathrm{K}})}{2}, \ \overrightarrow{N_{\mathrm{EF}}} = \frac{(x_{\mathrm{L}} - x_{\mathrm{K}}) \times (x_{\mathrm{B}} - x_{\mathrm{A}})}{2}.$$

These definitions and the structure of Table 1 ensure that $\overrightarrow{N_{\mathrm{XY}}}$ points from X to Y $((\mathrm{X}, \mathrm{Y}) \in \{(\mathrm{K}, \mathrm{L}), (\mathrm{A}, \mathrm{B}), (\mathrm{E}, \mathrm{F})\})$. Note that although there are more edges than faces, $u_{\mathrm{E}}$ and F contribute to the gradient similarly. The discrete divergence $\mathrm{div}^d : \mathbb{R}^{|\mathscr{D}|} \mapsto \mathbb{R}^{|\mathscr{T}|}$ is defined by

$$\left(\mathrm{div}^d(\xi)\right)_{\mathrm{C}} = \frac{1}{|\mathrm{C}|} \sum_{\mathrm{D}: \mathrm{D} \cap \mathrm{C} \neq \emptyset} \xi_{\mathrm{D}} \cdot \overrightarrow{N_{\mathrm{C}}}, \tag{2}$$

with $\overrightarrow{N_{\mathrm{C}}} = \overrightarrow{N_{\mathrm{KL}}}$ if $\mathrm{C} = \mathrm{K} \in \mathscr{M}$, $\overrightarrow{N_{\mathrm{C}}} = \overrightarrow{N_{\mathrm{AB}}}$ if $\mathrm{C} = \mathrm{A} \in \mathscr{N}$, $\overrightarrow{N_{\mathrm{C}}} = \overrightarrow{N_{\mathrm{EF}}}$ if $\mathrm{C} = \mathrm{E} \in \mathscr{E}$ and $\overrightarrow{N_{\mathrm{C}}} = -\overrightarrow{N_{\mathrm{EF}}}$ if $\mathrm{C} = \mathrm{F} \in \mathscr{F}$.

**Table 1** Construction of the diamond mesh. Diamond types are defined via their face/edge representation, noting the orientation of faces (orthogonal to $x$, $y$ or $z$ axis) and enumerating the four edges of the face

| Type | $x_F$ | $x_E$ | $x_K$ | $x_L$ | $x_A$ | $x_B$ |
|---|---|---|---|---|---|---|
| (x, 1) with $i$ odd | $(i, j, k)$ | $(i, j+1, k)$ | $(i-1, j, k)$ | $(i+1, j, k)$ | $(i, j+1, k-1)$ | $(i, j+1, k+1)$ |
| (x, 2) with $i$ odd | $(i, j, k)$ | $(i, j, k-1)$ | $(i-1, j, k)$ | $(i+1, j, k)$ | $(i, j-1, k-1)$ | $(i, j+1, k-1)$ |
| (x, 3) with $i$ odd | $(i, j, k)$ | $(i, j-1, k)$ | $(i-1, j, k)$ | $(i+1, j, k)$ | $(i, j-1, k+1)$ | $(i, j-1, k-1)$ |
| (x, 4) with $i$ odd | $(i, j, k)$ | $(i, j, k+1)$ | $(i-1, j, k)$ | $(i+1, j, k)$ | $(i, j+1, k+1)$ | $(i, j-1, k+1)$ |
| (y, 1) with $j$ odd | $(i, j, k)$ | $(i+1, j, k)$ | $(i, j-1, k)$ | $(i, j+1, k)$ | $(i+1, j, k+1)$ | $(i+1, j, k-1)$ |
| (y, 2) with $j$ odd | $(i, j, k)$ | $(i, j, k+1)$ | $(i, j-1, k)$ | $(i, j+1, k)$ | $(i-1, j, k+1)$ | $(i+1, j, k+1)$ |
| (y, 3) with $j$ odd | $(i, j, k)$ | $(i-1, j, k)$ | $(i, j-1, k)$ | $(i, j+1, k)$ | $(i-1, j, k-1)$ | $(i-1, j, k+1)$ |
| (y, 4) with $j$ odd | $(i, j, k)$ | $(i, j, k-1)$ | $(i, j-1, k)$ | $(i, j+1, k)$ | $(i+1, j, k-1)$ | $(i-1, j, k-1)$ |
| (z, 1) with $k$ odd | $(i, j, k)$ | $(i, j+1, k)$ | $(i, j, k-1)$ | $(i, j, k+1)$ | $(i+1, j+1, k)$ | $(i-1, j+1, k)$ |
| (z, 2) with $k$ odd | $(i, j, k)$ | $(i+1, j, k)$ | $(i, j, k-1)$ | $(i, j, k+1)$ | $(i+1, j-1, k)$ | $(i+1, j+1, k)$ |
| (z, 3) with $k$ odd | $(i, j, k)$ | $(i, j-1, k)$ | $(i, j, k-1)$ | $(i, j, k+1)$ | $(i-1, j-1, k)$ | $(i+1, j-1, k)$ |
| (z, 4) with $k$ odd | $(i, j, k)$ | $(i-1, j, k)$ | $(i, j, k-1)$ | $(i, j, k+1)$ | $(i-1, j+1, k)$ | $(i-1, j-1, k)$ |

These expressions simplify on uniform Cartesian grids (e.g. $|D| = \frac{2}{3}h^3$ for interior diamonds, with voxel length $h$), but no acceleration is obtained by implementing these simplifications, which is why we only present the general case.

## 3 An Application to the Subjective Surfaces Equation

In image processing, uniform Cartesian grids arise naturally because image information is given on pixels or voxels. We will illustrate the performance of the numerical scheme taking an application from this field. The subjective surfaces equation reads

$$\partial_t u + |\nabla u| \text{div}\left(g(|\nabla I|)\frac{\nabla u}{|\nabla u|}\right) = 0 \tag{3}$$

with $g(x) = \frac{1}{1+kx^2}$, $k > 0$, $I$ the (given) image intensity and Dirichlet boundary conditions. Numerically, the solution $u$ of Eq. (3) evolves to a piecewise constant function delimited by regions where $|\nabla I|$ is large. The support of the initial condition $u_0$ is chosen in the region of which the boundary should be determined.

### 3.1 Discretisation of the Subjective Surfaces Equation with CeVeFE DDFV

The meshing described in the previous section has the advantage that the unknowns correspond to the image voxels; we stress that no additional degrees of freedom, nor interpolated values, are used. Following [4, 9], we choose a semi-implicit time discretisation of a regularised form of Eq. (3), which yields a linear scheme:

$$\frac{u^{n+1} - u^n}{\Delta t} + (|\nabla^d u^n| + \varepsilon)\text{div}^d\left(g(|\nabla^d I|)\frac{\nabla^d u^{n+1}}{|\nabla^d u^n| + \varepsilon}\right) = 0, \tag{4}$$

with $\varepsilon > 0$. A symmetric scheme is obtained by multiplying Eq. (4) by the diagonal matrix $\Lambda_n$ with entries $((|\nabla^d u^n|_C + \varepsilon)/|C|)^{-1}$:

$$\Lambda_n u^{n+1} + \Delta t |C| \text{div}^d\left(g(|\nabla^d I|)\frac{\nabla^d u^{n+1}}{|\nabla^d u^n| + \varepsilon}\right) = \Lambda_n u^n. \tag{5}$$

Observe that the matrix $M$ with $Mu = |C|\text{div}^d\left(g(|\nabla^d I|)\frac{\nabla^d u}{|\nabla^d u^n| + \varepsilon}\right)$ is computed in the following way. Let $\lambda_D := \frac{g(|\nabla^d I|_D)}{|\nabla^d u^n|_D + \varepsilon}$, which is known from the previous iteration, and note that in order to calculate this quantity, an approximation of the norm of the full gradient is needed (which basic Finite Difference schemes and some Finite Volume schemes do not yield). Due to the uniform Cartesian grid structure,

$$\overrightarrow{N_{\text{KL}}} \cdot \overrightarrow{N_{\text{EF}}} = \overrightarrow{N_{\text{KL}}} \cdot \overrightarrow{N_{\text{AB}}} = \overrightarrow{N_{\text{EF}}} \cdot \overrightarrow{N_{\text{AB}}} = 0, \tag{6}$$

yielding

$$
\begin{aligned}
(Mu)_{\text{K}} &= \sum_{\text{D} \in \text{D}_{\text{K}}} \frac{\lambda_{\text{D}}}{3|\text{D}|} (u_{\text{L}} - u_{\text{K}}) ||\overrightarrow{N_{\text{KL}}}||^2, & (Mu)_{\text{A}} &= \sum_{\text{D} \in \text{D}_{\text{A}}} \frac{\lambda_{\text{D}}}{3|\text{D}|} (u_{\text{B}} - u_{\text{A}}) ||\overrightarrow{N_{\text{AB}}}||^2, \\
(Mu)_{\text{E}} &= \sum_{\text{D} \in \text{D}_{\text{E}}} \frac{\lambda_{\text{D}}}{3|\text{D}|} (u_{\text{F}} - u_{\text{E}}) ||\overrightarrow{N_{\text{EF}}}||^2, & (Mu)_{\text{F}} &= \sum_{\text{D} \in \text{D}_{\text{F}}} \frac{\lambda_{\text{D}}}{3|\text{D}|} (u_{\text{E}} - u_{\text{F}}) ||\overrightarrow{N_{\text{EF}}}||^2.
\end{aligned}
\tag{7}
$$

Therefore, the meshes $\mathcal{M}$, $\mathcal{N}$ and $\mathcal{FE}$ are not coupled in the resolution of (5), only to the previous time step by $|\nabla u^n|$, accelerating the numerical resolution.

## 3.2 Iterating Over Diamonds

We stress that the quantities needed for the resolution of Eq. (5) can be computed only using the diamond structure. The following information has to be stored for each diamond: the point references explained in Table 1, the volumes of the diamond of its eight constituting tetrahedra (see (1)), and the vectors $\overrightarrow{N_{\text{KL}}}$, $\overrightarrow{N_{\text{AB}}}$ and $\overrightarrow{N_{\text{EF}}}$.

The matrix $M$ can then be assembled efficiently by iterating over $\text{D} \in \mathcal{D}$ and by computing for each diamond the contributions at the indices corresponding to $x_{\text{K}}$, $x_{\text{L}}$, $x_{\text{A}}$, $x_{\text{B}}$, $x_{\text{E}}$ and $x_{\text{F}}$ via the formulas (7). Similarly, the measure of the control volumes can be assembled from the eight tetrahedra constituting the diamonds. These procedures, including the construction of the diamond mesh, are easily vectorised.

## 3.3 DDFV Solution

DDFV solutions, defined on overlapping meshes, naturally give rise to averaged discrete solutions [2, 3]. In our case, based on the solution $u$ of (5) at the final time $T$, on each mesh ($\mathcal{M}$, $\mathcal{N}$ and $\mathcal{FE}$) a cell-wise piecewise constant function $(u_{\mathcal{M}}, u_{\mathcal{N}}, u_{\mathcal{FE}})$ is defined. The DDFV solution is

$$u_{DDFV} = \frac{1}{3} \left( u_{\mathcal{M}} + u_{\mathcal{N}} + u_{\mathcal{FE}} \right),$$

which is constant on each tetrahedron constituting the diamond cells. In order to visualize the DDFV solution on the Cartesian grid, it is projected on the cells of $\mathcal{T}$:

$$u_{DDFV}^{cart} = \left( \frac{1}{|\text{C}|} \int_{\text{C}} u_{DDFV} \right)_{\text{C} \in \mathcal{T}}.$$

**Fig. 2** 2D cuts of a 3D SPECT image (through then *center* of the tumour and parallel to the $x$, $y$ and $z$ axes, respectively) showing the density of proliferating tumour cells, which are localised at the boundary. The tumour reconstruction is marked by the *black line*

This averaging is the price we pay for avoiding additional unknowns (as compared to [9], in 2D); indeed, in 3D it is crucial to reduce the number of degrees of freedom. It is important to note that the use of $u_{DDFV}^{cart}$ is generally necessary and cannot be replaced by the evaluation of $u$. Indeed, the weak coupling of the three meshes due to the semi-implicit time discretisation allows $u$ to contain local checkerboard structures caused by noise whereas $u_{DDFV}$ is smooth.

### 3.4 Numerical Results

The numerical scheme is illustrated on 3D SPECT images visualising proliferating tumour cells. These cells are mainly localised on the tumour boundary but do not cover the entire surface, notably due to physical constraints such as bones. We want to obtain the volume and shape of the tumour based on these images. In practice, tumour diameters are often measured manually and volume is approximated with an ellipsoid formula. The numerical method described above permits to obtain a less heuristical estimation of the volume, also indicating the shape. Voxels of the $N_x \times N_y \times N_z$-image are numbered in a classical way by $N(i, j, k) = i + N_x \cdot (j - 1) + N_y \cdot N_x \cdot (k - 1)$, such that $\Lambda_n + \Delta t M$ is a band matrix. Figure 2 shows the three different 2D cuts of the original image and the reconstructed tumour volume.

## 4 Conclusion

We have presented here the implementation of a 3D CeVeFE DDFV scheme using a Cartesian structure without introducing artificial unknowns, which is an important property in view of the high computational complexity in 3D. It comes at the cost of a mild smoothing by projecting the discrete solution on underlying voxels.

The implementation presented here finds an application in image processing, where uniform Cartesian grids naturally arise. The fact that DDFV schemes can be used on degenerate meshes makes the implementation relevant for non-uniform

Cartesian grids, for example the highly deformed Kershaw meshes appearing in porous media [8]. Similar band matrix profiles can be obtained in these cases.

We have successfully used a 3D DDFV discretisation of the subjective surfaces equation for the reconstruction of the tumour shape on an exemplary SPECT image. A subsequent step would be to test the performance of an automatised version on a large number of images and to compare it to the ellipsoid formula.

It should also be stressed that the DDFV framework is one out of many possible discretisation strategies, each with their advantages and shortcomings. It is hoped that this work permits an easy access to the DDFV approach.

# References

1. Andreianov, B., Bendahmane, M., Hubert, F.: On 3d ddfv discretization of gradient and divergence operators. part ii. Comput. Method Appl. M. **13**(4), 369–410 (2013)
2. Andreianov, B., Bendahmane, M., Hubert, F., Krell, S.: On 3d ddfv discretization of gradient and divergence operators. Part I. IMA J. Numer. Anal. **32**(4), 1574–1603 (2012)
3. Andreianov, B., Boyer, F., Hubert, F.: Discrete duality finite volume schemes for leray-lions type elliptic problems on general 2d meshes. Numer. Meth. Part D E **23**(1), 145–195 (2007)
4. Corsaro, S., Mikula, K., Sarti, A., Sgallari, F.: Semi-implicit covolume method in 3d image segmentation. J. Sci. Comput. **28**(6), 2248–2265 (2006)
5. Coudière, Y., Hubert, F.: A 3d discrete duality finite volume method for nonlinear elliptic equations. SIAM J. Sci. Comput. **33**(4), 1739–1764 (2011)
6. Domelevo, K., Omnes, P.: A finite volume method for the laplace equation on almost arbitrary two-dimensional grids. Math. Model Numer. Anal. **39**(6), 1203–1249 (2005)
7. Eymard, R., Handlovičová, A., Mikula, K.: Study of a finite volume scheme for the regularized mean curvature flow level set equation. IMA J. Numer. Anal. **31**(3), 813–846 (2011)
8. Eymard, R., Henry, G., Herbin, R., Hubert, F., Klöfkorn, R., Manzini, G.: 3D benchmark on discretization schemes for anisotropic diffusion problems on general grids in: finite volumes for complex applications VI, 895–930 (2011)
9. Handlovičová, A., Kotorová, D.: Numerical analysis of a semi-implicit ddfv scheme for the regularized curvature driven level set equation in 2d. Kybernetika **49**, 829–854 (2013)
10. Hermeline, F.: A finite volume method for the approximation of diffusion operators on distorted meshes. J. Comput. Phys. **160**(2), 481–499 (2000)
11. Hermeline, F.: A finite volume method for approximating 3d diffusion operators on general meshes. J. Comput. Phys. **288**(16), 5763–5786 (2009)
12. Mikula, K., Peyriéras, N., Remešíková, M., Stasova, S.: Segmentation of 3D cell membrane images by PDE methods and its applications. Comput. Biol. Med. **41**, 326–339 (2011)
13. Mikula, K., Remešíková, M.: Finite volume schemes for the generalized subjective surface equation in image segmentation. Kybernetika **45**(4), 646–656 (2009)
14. Sarti, A., Malladi, R., Sethian, J.: Subjective surfaces: a method for completing missing boundaries. Proc. Nat. Acad. Sci. **12**(97), 6258–6263 (2000)
15. Sethian, J.: Level Set Methods and Fast Marching Methods: Evolving Interfaces in Computational Geometry, Fluid Mechanics, Computer Vision, and Material Science. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, New York (1999)

# MPFA Algorithm for Solving Stokes-Brinkman Equations on Quadrilateral Grids

Oleg Iliev, Ralf Kirsch, Zahra Lakdawala and Galina Printsypar

**Abstract**  This work is concerned with the development of a robust and accurate numerical method for solving the Stokes-Brinkman system of equations, which describes a free fluid flow coupled with a flow in porous media. Quadrilateral boundary fitted grid with a sophisticated finite volume method, namely MPFA O-method, is used to discretize the system of equations. Numerical results for two examples are presented, namely, channel flow and flow in a ring with a rolled porous medium.

## 1 Introduction

There is much work invested in the numerical simulation of coupled free fluid and porous media flow, and one still aims towards faster, more robust, and more accurate simulators. This problem arises in many applications, such as filtration, membranes, hydrology, and so on. To account for the free fluid flow coupled with the porous media flow, different approaches exist (see e.g. the detailed discussion in [5] and references therein). We are concerned with a mathematical model using Stokes-Brinkman system of equations (e.g. [3, 4] and references therein), and our primary target is application of the developed algorithms to solving filtration problems.

---

O. Iliev · R. Kirsch · Z. Lakdawala
Department of Flow and Material Simulations, Fraunhofer ITWM, Kaiserslautern, Germany
e-mail: oleg.iliev@itwm.fraunhofer.de

R. Kirsch
e-mail: ralf.kirsch@itwm.fraunhofer.de

Z. Lakdawala
e-mail: zahra.lakdawala@itwm.fraunhofer.de

O. Iliev · G. Printsypar (✉)
Numerical Porous Media Center (NumPor), King Abdullah University of Science
and Technology (KAUST), Thuwal, Saudi Arabia
e-mail: galina.printsypar@kaust.edu.sa

One of the challenges in solving filtration problems is the complicated shape of the computational domain. The voxel grid based methods from [3, 4] were successfully applied in numerous industrial and academic problems, but in the case of very complicated domains, a high number of voxels is needed. In 2D case, boundary adapted quadrilateral grids due to their adaptation to the boundary and efficient use of computational resources, in certain cases may be a better choice, compared to the Cartesian grids. However, more sophisticated discretization and numerical algorithms need to be developed in this case. Here we will focus on a numerical algorithm which adapts MPFA O-method (see e.g. [1, 2]). MPFA is widely used, e.g., in solving scalar elliptic equations, but up to the authors knowledge, there was no investigation of MPFA discretization for Stokes-Brinkman problems, and this paper aims to fill this gap. Furthermore, the discontinuity of the coefficients in the Stokes-Brinkman model requires special interpolation at the porous-fluid interface [3]. However, due to lack of space, the latter will not be described here and will be a subject of another paper.

## 2 Modeling of Coupled Free and Porous Media Flow

There exist different models for simulating a free fluid flow coupled with a flow in porous media. One of the popular approaches is to use Stokes and Darcy flow problems with interface coupling conditions. Another approach is to use the Navier-Stokes-Brinkman model that had been presented for the fluid flow through a filter element in our earlier work (see [3] and references therein). In this study we are concerned with a reduced model using the Stokes-Brinkman equations, which read

$$\rho \frac{\partial \mathbf{u}}{\partial t} - \nabla \cdot (\mu \nabla \mathbf{u}) + \mu \mathbf{K}^{-1} \mathbf{u} = -\nabla p \,, \quad \mathbf{x} \in \Omega,$$
$$\nabla \cdot \mathbf{u} = 0 \,, \quad \mathbf{x} \in \Omega. \tag{1}$$

Here $\mathbf{u}$ and $p$ denote the fluid velocity vector and the fluid pressure, respectively. Moreover, $\mu$ is the fluid dynamic viscosity. $\Omega$ is the computational domain, which consists of two nonintersecting subsets, namely the fluid domain $\Omega_f$ and the porous domain $\Omega_p$. $\mathbf{K}$ is the intrinsic permeability of the porous medium in $\Omega_p$ and $\mathbf{K}^{-1} = 0$ in $\Omega_f$.

A typical set of boundary conditions looks as follows

$$\mathbf{u}(\mathbf{x}) = \mathbf{u}_{\text{in}}(\mathbf{x}), \ \mathbf{x} \in \Gamma_i; \quad \sigma \cdot \mathbf{n} = 0, \ \mathbf{x} \in \Gamma_o; \quad \mathbf{u}(\mathbf{x}) = 0, \ \mathbf{x} \in \Gamma_s; \tag{2}$$

where $\Gamma_i$, $\Gamma_o$ and $\Gamma_s$ denote the inlet, outlet, and solid wall boundaries, $\sigma$ is the stress tensor, $\mathbf{n}$ is the outward unit normal to $\Gamma_o$. More details can be found, for example, in [3, 5].

# 3 Chorin Type Method for Stokes-Brinkman Equations

*Chorin type algorithm and its discretization.* At first, let us introduce for simplicity operators $\mathscr{D} = \nabla \cdot \mu \nabla$ and $\mathscr{B} = \mu \mathbf{K}^{-1}$. Then, the system of equations (1) to be solved reads

$$\rho \frac{\partial \mathbf{u}}{\partial t} - \mathscr{D}\,\mathbf{u} + \mathscr{B}\,\mathbf{u} = -\nabla p, \quad \mathbf{x} \in \Omega, \tag{3}$$

$$\nabla \cdot \mathbf{u} = 0, \quad \mathbf{x} \in \Omega. \tag{4}$$

For a given discretization time step $\tau > 0$ and an initial time moment $t_0 \geq 0$, we define $t_n = t_0 + n\,\tau$, $n = 0, 1, \ldots$. Let $\mathbf{u}^n$ and $p^n$ denote the approximation of $\mathbf{u}$ and $p$ at time $t_n$. The fractional time step discretization can be written as

$$\frac{\rho}{\tau}\left(\mathbf{u}^* - \mathbf{u}^n\right) - \mathscr{D}\mathbf{u}^* + \mathscr{B}\mathbf{u}^* = -\nabla p^n, \tag{5}$$

$$\frac{\rho}{\tau}\left(\mathbf{u}^{n+1} - \mathbf{u}^*\right) + \mathscr{B}\mathbf{u}^{n+1} - \mathscr{B}\mathbf{u}^* = -\left(\nabla p^{n+1} - \nabla p^n\right), \tag{6}$$

$$\nabla \cdot \mathbf{u}^{n+1} = 0, \tag{7}$$

where $\mathbf{u}^*$ is a prediction to the fluid velocity.

*Finite volume integral formulation.* The computational domain $\Omega$ is subdivided into a set $\mathscr{V}$ of quadrilateral finite volumes $v$. After integrating the system of equations (5)–(7) over each $v$ and performing some transformations we obtain the following Chorin type algorithm (for more details see [3])

$$\int_v \frac{\rho}{\tau}\left(\mathbf{u}^* - \mathbf{u}^n\right)d\mathbf{x} - \int_v \mathscr{D}\mathbf{u}^* d\mathbf{x} + \int_v \mathscr{B}\mathbf{u}^* d\mathbf{x} = -\int_v \nabla p^n d\mathbf{x}, \quad v \in \mathscr{V}; \tag{8}$$

$$-\int_v \nabla \cdot \left(V\left(\frac{\rho}{\tau} + \mathscr{B}\right)\right)^{-1} \nabla p' d\mathbf{x} = -\int_v \nabla \cdot \mathbf{u}^* d\mathbf{x}, \quad v \in \mathscr{V}; \tag{9}$$

$$\mathbf{u}_v^{n+1} = \mathbf{u}_v^* - \int_v \left(V\left(\frac{\rho}{\tau} + \mathscr{B}\right)\right)^{-1} \nabla p' d\mathbf{x}, \quad v \in \mathscr{V}; \tag{10}$$

$$p_v^{n+1} = p_v^n + p_v', \quad v \in \mathscr{V}; \tag{11}$$

where index '$v$' denotes volume averaged variables, $p'$ is the pressure correction, $V$ is the measure of the finite volume $v$, $V = mes(v)$.

# 4 Space Discretization Using MPFA

The equations are discretized using the cell-centered collocated finite volume approach on quadrilaterals. Due to the complex quadrilateral grid arrangement and varying (discontinuous) coefficients in the Stokes-Brinkman case, special attention

**Fig. 1** Primary and dual grids



is paid to the spatial discretization of the terms. The Multipoint Flux Approximation method (MPFA) (see [1, 2]) is employed to approximate the following type of terms $\int_v \nabla \cdot \alpha \nabla \phi d\mathbf{x}$ and $\int_v \nabla \phi d\mathbf{x}$, where $\phi$ is a scalar function representing one of the velocity components or the pressure, $\alpha$ is a constant coefficient in our case, but this discretization technique can account for a full tensor coefficient $\alpha$.

At first, let us describe some necessary details. We employ two grids, name primary and dual grids. As shown in Fig. 1, the primary grid consists of quadrilaterals marked in red. The dual grid marked in blue is formed by connecting the centers of primary quadrilaterals and the midpoints of their edges. The duals are further divided into four parts, each part belongs to a different quadrilateral.

*Approximation of $- \int_v \nabla \cdot \alpha \nabla \phi d\mathbf{x}$*
Employing the Gauss' divergence theorem, we get the flux $f_v$ through the boundary $\partial v$ of the control volume $v$

$$f_v = - \int_v \nabla \cdot \alpha \nabla \phi d\mathbf{x} = - \oint_{\partial v} (\alpha \nabla \phi) \cdot \mathbf{n} \, dl. \tag{12}$$

The evaluation of the flux is reduced to calculation of the integral in Eq. (12), which is standard for the MPFA method (see [1]). We can write a flux expression for the quadrilateral $v$ as follows

$$f_v = \sum_{i=0}^{n_e-1} f_i \approx \sum_{i=0}^{n_e-1} \left( \sum_{j=0}^{n_p-1} t_{ij} \hat{\phi}_j \right). \tag{13}$$

As shown in Fig. 1, $f_i$ is the flux through half edges $e_i$ with $n_e = 8$ denoting the number of half edges of the primary cell and $n_p = 4$ denoting the number of primary quadrilaterals in the dual cell, $t_{ij}$ is called the transmissibility coefficient, computed via the MPFA O-method (see Eq. (26) in [1]). $\hat{\phi}_j$ is the value of function $\phi$ at the center of the $j$th primary quadrilateral in the dual cell (see Fig. 1), which also represent unknowns in the discretized system.

*Approximation of $\int_v \nabla \phi d\mathbf{x}$*

Here the MPFA method is extended in order to approximate the gradient operators in the last terms of Eqs. (8) and (10). Using the Green's theorem and fundamental relations between ordinary and line integrals, for the first gradient component we obtain

$$\int_v \frac{\partial \phi}{\partial x} d\mathbf{x} = \oint_{\partial v} 0 dx + \phi dy = \oint_{\partial v} (\phi, 0) \cdot (dy, -dx) = \oint_{\partial v} (\phi, 0) \cdot \mathbf{n} dl$$

$$= \oint_{\partial v} \phi n_x dl \approx \sum_{i=0}^{n_e - 1} \phi_i s_i n_{x,i}. \tag{14}$$

where $\phi_i$ are the values of $\phi$ on $e_i$ (see Fig. 1), $s_i$ is the measure of $e_i$, $s_i = mes(e_i)$, $\mathbf{n}_i = (n_{x,i}, n_{y,i})$ is the outward unit normal vector of $e_i$. Values $\phi_i$ can also be estimated using MPFA method like fluxes $f_i$ in (13) using values $\hat{\phi}_j$ but with different transmissibility coefficients (see Eq. (24) in [1]). Similar procedure can be carried out to approximate the second component of the gradient $\int_v \partial \phi / \partial y d\mathbf{x}$.

## 5 Numerical Results

MPFA is one of the discretization techniques which was developed to approximate fluxes with full tensor coefficient on irregular grids. Standard discretization techniques such as two-point flux approximation can lead to unphysical or inaccurate results for irregular grids. In this section we present a channel example and compare its solution on different regular and irregular grids along with the analytical solution. Example for a radial Stokes-Brinkman flow in a ring is also presented on irregular grid. Obtained values are compared to their analytic counterparts.

*Stationary channel flow between two parallel plates.* The domain is a 2D channel of width $d$ (along $x_1$) between two infinitely large plates. The equations are considered in Cartesian coordinates. We assume that the flow is stationary $\partial \mathbf{u} / \partial t \equiv 0$, and there is no flow in $x_1$ direction, $u_1 \equiv 0$. Moreover, the pressure is linear in $x_2$ and does not depend on $x_1$, i.e. $\partial p / \partial x_2 \equiv -C_p$. The boundary conditions are "no-slip", where $\mathbf{u}(x_1 = 0) = \mathbf{u}(x_1 = d) = 0$. Then, in infinitely long channel the analytical solution yields

$$u_2(x_1) = \frac{C_p}{2\mu} \left( dx_1 - x_1^2 \right) = \frac{p_{\text{drop}}}{2\mu L} \left( dx_1 - x_1^2 \right),$$

where $p_{\text{drop}}$ is the pressure drop along the channel at distance $L$, $\mu$ is the viscosity. Then, the maximum centerline velocity reads

$$u_{2,max} \left( \frac{d}{2} \right) = \frac{P_{drop} d^2}{8 \mu L} = 1.5 u_{in}, \tag{15}$$

**Fig. 2** Different quadrilateral grids: stretched Cartesian grid (*left*), fish-bone grid (*middle*), and mosaic grid (*right*)



**Fig. 3** Velocity distributions on different grids, i.e. stretched Cartesian grid (*left*), fish-bone grid (*middle*), and mosaic shaped quadrilaterals (*right*). Resolution is $80 \times 40$ quadrilaterals

where $u_{in}$ is the inflow or average velocity.

Here, we present the results for the channel problem with a simple rectangular computational domain. The numerical algorithm was tested on different complex grids (see Fig. 2), such as stretched Cartesian grid, fish-bone grid, and mosaic grid, with different resolutions. The inflow parameters used for this problem are the viscosity $\mu = 0.2$ kg/ms, the inlet velocity $u_{in} = 0.015$ m/s, the fluid density $\rho = 800$ kg/m$^3$. Dimensions of the channel are $d = 1$ m and $L = 5$ m.

The velocity profiles on the different grids are illustrated in Fig. 3. Table 1 summarizes the results on different computational domains for different grid resolutions. According to the analytical solution (15), the maximum centerline velocity is $u_{2,max} = 0.0225$ m/s and the pressure drop is $p_{drop} = 0.18$ Pa. It can be seen that the complex grids result in a comparable pressure and velocity accuracy. Moreover the results on different resolutions also compare well. It can be concluded that the

**Table 1** Summary of results on different computational domains for the flow between two parallel plates

| Geometry | Resolution | $p_{drop}$ (Pa) | $u_{2,max}$ (m/s) |
|---|---|---|---|
| Cartesian | $80 \times 40$ | 0.1749 | 0.0228 |
| | $40 \times 20$ | 0.19 | 0.02269 |
| Fish-bone | $80 \times 40$ | 0.1745 | 0.0228 |
| | $80 \times 20$ | 0.1859 | 0.02214 |
| | $40 \times 20$ | 0.185 | 0.02268 |
| Mosaic | $80 \times 40$ | 0.1735 | 0.02268 |
| | $80 \times 20$ | 0.1829 | 0.02192 |
| | $40 \times 20$ | 0.1882 | 0.02247 |

Convergence and grid study

**Table 2** Input parameters for the ring problem

| | | | | | | |
|---|---|---|---|---|---|---|
| Inner radius, $R_1$ | (m) | 1 | Inflow velocity, $u_r^{in}$ | (m/s) | $6.36692 \times 10^{-6}$ |
| Outer radius, $R_2$ | (m) | 5 | Permeability, $K$ | (m$^2$) | $1 \times 10^{-12}$ |
| Inner porous radius, $r_1$ | (m) | 2.516375 | Fluid density, $\rho$ | (kg/m$^3$) | 800 |
| Outer porous radius, $r_2$ | (m) | 3.48156 | Viscosity, $\mu$ | (kg/ms) | 0.2 |

algorithm works well on an arbitrary complex composition of quadrilateral grids for the Stokes system of equations.

*Radial Stokes-Brinkman flow in a ring.* The domain is a 2D ring (annulus) which contains an additional ring-shaped porous medium. The inlet is located on the outer circle boundary and the outlet is on the inner circle boundary. They are separated by the porous medium. On the inlet and the outlet we use Dirichlet boundary conditions for the fluid velocity $u_r(R_1) = u_r^{out}$ and $u_r(R_2) = u_r^{in}$. Using the continuity equation, the solution for the fluid velocity is given by

$$u_r(r, \varphi) = u_r(r) = \frac{R_2\, u_r^{in}}{r}, \quad u_\varphi(r, \varphi) \equiv 0, \text{ for } r \in [R_1, R_2] \text{ and } \varphi \in [0, 2\pi],$$

where $\mathbf{u} = (u_r, u_\varphi)$ is the fluid velocity in polar coordinates $r$ and $\varphi$. All other parameters are introduced in Table 2. The momentum equation in the porous media in polar coordinates reads

$$\frac{\mu}{K} \frac{R_2\, u_r^{in}}{r} = -\frac{\partial p}{\partial r}, \text{ for } r \in [r_1, r_2].$$

Then, an analytical expression for the pressure drop across the porous medium yields

$$p(r_2) - p(r_1) = -\frac{\mu}{K} R_2\, u_r^{in} In\left(\frac{r_2}{r_1}\right).$$

**Fig. 4** Numerical results for the ring with a rolled porous medium inside

Figure 4 (left) illustrates the computational domain discretized using the adaptive grid with resolution $120 \times 29$. The geometries in red and blue denote the porous and fluid regions, respectively. In the Fig. 4 (middle, right), the pressure and velocity profiles are shown. The computed pressure difference is 2061.2 KPa. This compares well to the analytically computed pressure difference, which is equal to 2067.1 KPa. Note that the analytical pressure difference in the fluid domain is neglected as it is very small compared to the pressure drop across the porous media.

## 6 Summary

In this paper, we have discussed the numerical algorithm and its associated discretization for solving the Stokes and Stokes–Brinkman system of equations on quadrilateral grids. The discretization method employs MPFA O-method to approximate first and second order terms in the Eqs. (8)–(11) as described in Sect. 3. The results for the two examples compare well to the analytical solutions.

## References

1. Aavatsmark, I.: An introduction to multipoint flux approximations for quadrilateral grids. Comput. Geosci. **6**, 405–432 (2002)
2. Aavatsmark, I.: Multipoint flux approximation methods for quadrilateral grids. In: 9th International Forum on Reservoir Simulation, Abu Dhabi (2007)
3. Ciegis, R., Iliev, O., Lakdawala, Z.: On parallel numerical algorithms for simulating industrial filtration problems. Berichte des Fraunhofer ITWM. **114**, 24 (2007)
4. Iliev, O., Laptev, V.: On numerical simulation of flow through oil filters. Comput. Visual. Sci **6**, 139–146 (2004)
5. Laptev, V.: Numerical solution of coupled flow in plain and porous media. Ph.D. thesis, Technical University Kaiserslautern (2004)

# Nonlinear Monotone FV Schemes for Radionuclide Geomigration and Multiphase Flow Models

Ivan Kapyrin, Kirill Nikitin, Kirill Terekhov and Yuri Vassilevski

**Abstract** We present applications of the nonlinear monotone finite volume method to radionuclide transport and multiphase flow in geological media models. The scheme is applicable for full anisotropic discontinuous permeability or diffusion tensors and arbitrary conformal polyhedral cells. We consider two versions of the nonlinear scheme: two-point flux approximation preserving positivity of the solution and compact multi-point flux approximation that provides discrete maximum principle. We compare the new nonlinear schemes with the conventional linear two-point and multi-point (O-scheme) flux approximations. Both new nonlinear schemes have compact stencils and a number of important advantages over the traditional linear discretizations. Two industrial applications are discussed briefly: radionuclides transport modeling within the radioactive waste safety assessment and multiphase flow modeling of oil recovery process.

## 1 Introduction

A simple and accurate conservative method applicable to general conformal meshes and full anisotropic tensor permeability coefficients, is much-in-demand among engineers. The maximum principle is one of the important properties of solutions of

I. Kapyrin · K. Nikitin (✉) · K. Terekhov · Y. Vassilevski
Institute of Numerical Mathematics, Gubkina 8, Moscow, Russia
e-mail: ivan.kapyrin@gmail.com

K. Terekhov
e-mail: kirill.terehov@gmail.com

I. Kapyrin · K. Nikitin
Institute of Nuclear Safety, B. Tulskaya 52, Moscow, Russia
e-mail: nikitin.kira@gmail.com

Y. Vassilevski
Moscow Institute of Physics and Technology, Institutski 9, Dolgoprudny, M.R., Russia
e-mail: yuri.vassilevski@gmail.com

partial differential equations (PDEs) such as the diffusion or heat equation. Its discrete counterpart is a very desirable property to have in a numerical scheme. Unfortunately, the schemes satisfying the discrete maximum principle (DMP) impose severe limitations on mesh regularity [6] and problem coefficients. Violation of the DMP leads to various numerical artifacts, such as heat flow from a cold material to a hot one, that can be amplified by physics non-linearity.

The classical two-point finite volume (FV) scheme for diffusion problems defines a two-point flux approximation (TPFA) across a mesh face as a difference of two concentrations at neighboring cells times a transmissibility coefficient. It results in a system of algebraic equation with an M-matrix with diagonal dominance in rows, which implies immediately the DMP [15]. However, accuracy of this scheme depends on mesh geometry and mutual orientation of mesh faces and principle directions of the diffusion tensor. More precisely, the co-normal vector for a face must be collinear to the vector connecting neighboring collocation points, which is clearly the impossible requirement for arbitrary tensors and/or arbitrary polyhedral cells. The multi-point flux approximation (MPFA) scheme solves accuracy problem by using more than two points in the flux stencil [1] and a matrix of transmissibility coefficients. The MPFA scheme provides a second-order accurate approximation of concentrations but is only conditionally stable and conditionally monotone [14].

A new research direction pioneered by Le Potier [7] uses a two-point flux stencil with two coefficients that depend on the concentrations in neighboring cells. Nonlinear FV schemes with TPFA proposed in [3, 5, 7, 9, 10, 13, 18] guarantee solution positivity on general meshes and for general tensor coefficients.

For general meshes and coefficients the DMP requires a nonlinear multi-point flux approximation. For diffusion problems, such schemes were proposed in [8, 19] using auxiliary unknowns at mesh vertices. Later an interpolation-free multi-point nonlinear approximation of diffusive fluxes was proposed for two-dimensional [11] and three-dimensional cases [2, 4]. The resulting scheme has the minimal stencil and reduces to the classical two-point FV scheme on Voronoi or rectangular meshes and for scalar (and, in a few cases, diagonal tensor) coefficients.

In this article, we present two our FV schemes for the steady-state diffusion equation with anisotropic coefficients: both schemes work on general polyhedral meshes and have a compact stencil, the first preserves non-negativity of the discrete solution and the second satisfies the DMP. We also briefly consider two applications of the nonlinear schemes to subsurface flows: simulation of radionuclides geomigration from a nuclear waste disposal and multiphase flow modeling of oil recovery process.

The paper outline is as follows. In Sect. 2 we introduce our nonlinear FV schemes for the steady-state diffusion equation. In Sect. 3 we present a new parallel toolkit and two industrial applications of the presented FV schemes.

## 2 Nonlinear Finite Volume Methods

Let $\Omega$ be a three-dimensional polyhedral domain with boundary $\Gamma$. The mixed form of the diffusion equation for unknown concentration $c$ with the Dirichlet boundary condition is as follows:

$$\mathbf{q} = -\mathbb{K}\nabla c, \quad \operatorname{div} \mathbf{q} = f \ \text{ in } \ \Omega,$$
$$c = g \ \text{ on } \ \Gamma. \tag{1}$$

Here $\mathbb{K}(\mathbf{x})$ is a symmetric positive definite discontinuous (possibly anisotropic) diffusion tensor, $f(\mathbf{x})$ is a source term, and $g(\mathbf{x})$ is a boundary data.

A discretization scheme can have two additional properties: discrete maximum (or minimum) principle and non-negativity of the discrete solution. The minimum principle states that for $f \geq 0$ the concentration $c(\mathbf{x})$ satisfies:

$$\min_{\mathbf{x}\in\bar{\Omega}} c(\mathbf{x}) \geq \min\{0, \ \min_{\mathbf{x}\in\Gamma} g(\mathbf{x})\}.$$

The maximum principle is formulated similarly. In the following we shall refer to both principles as the maximum principle. Non-negativity is a weaker property which stems from the minimum principle: for non-negative $f$ and $g$ one has non-negative $c(\mathbf{x})$. A numerical scheme can provide non-negativity of $c$ but violate the discrete maximum principle (DMP) and thus can produce oscillations.

The cell-centered FV scheme uses one degree of freedom, $C_T$, per cell $T$ collocated at cell barycenter $\mathbf{x}_T$. Integrating the mass balance Eq. (1) over $T$ and using the divergence theorem, we obtain:

$$\sum_{f\in\partial T} \sigma_{T,f}\, \mathbf{q}_f \cdot \mathbf{n}_f = \int_T f \, d\mathbf{x}, \quad \mathbf{q}_f = \frac{1}{|f|}\int_f \mathbf{q}\, ds, \tag{2}$$

where $\mathbf{q}_f \cdot \mathbf{n}_f$ is the total flux across face $f$, and $\sigma_{T,f}$ is either 1 or $-1$ depending on the mutual orientation of normal vector to face $\mathbf{n}_f$ and the outer normal to cell boundary $\mathbf{n}_T$.

Both nonlinear flux approximation schemes exploit the same idea of vector expansion. First we need to find a triplet of three vectors $\mathbf{t}_{1*}$ connecting $\mathbf{x}_{T_1}$ with other collocation points such that the co-normal vector $\ell_f = \mathbb{K} \cdot \mathbf{n}_f$ can be expanded

$$\ell_f = \alpha_{1a}\, \mathbf{t}_{1a} + \beta_{1b}\, \mathbf{t}_{1b} + \gamma_{1c}\, \mathbf{t}_{1c}, \quad \alpha_{1a} \geq 0, \ \beta_{1b} \geq 0, \ \gamma_{1c} \geq 0, \tag{3}$$

where $a, b, c$ are indexes of neighboring cells.

Since the flux normal component is the directional derivative along the co-normal vector $\ell_f$, it can be represented as the sum of three directional derivatives along $\mathbf{t}_{1*}$ which are approximated by central differences:

$$(\mathbf{q}_f \cdot \mathbf{n}_f)_h^{(1)} = \alpha_{1a}\, (C_a - C_1) + \beta_{1b}\, (C_b - C_1) + \gamma_{1c}\, (C_c - C_1). \tag{4}$$

**Fig. 1** Two representations of co-normal vector $\ell_1 = -\ell_2 = \mathbb{K} \cdot \mathbf{n}_e$ (2D example)



For the opposite co-normal vector $-\ell_e$ we have similar representation with another triplet and central differences, see Fig. 1 for the 2D example:

$$(-\mathbf{q}_f \cdot \mathbf{n}_f)_h^{(2)} = \alpha_{2k} \, (C_k - C_2) + \beta_{2l} \, (C_l - C_2) + \gamma_{2m} \, (C_m - C_2). \qquad (5)$$

Our flux discretization is a linear combination of approximations (4) and (5) with coefficients $\mu_+$ and $\mu_-$. For the sake of approximation the linear combination should be convex:

$$\mu_+ + \mu_- = 1.$$

The second equation for $\mu_\pm$ is dictated by the goal of the method:

- To obtain the two-point discretization, we get rid of unwanted concentrations in the flux stencil:

$$\mu_+(\alpha_{1a} \, C_a + \beta_{1b} \, C_b + \gamma_{1c} \, C_c) - \mu_-(\alpha_{2k} \, C_k + \beta_{2l} \, C_l + \gamma_{2m} \, C_m) = 0.$$

- To provide the DMP, we balance the contributions of one-sided fluxes:

$$\mu_+(\mathbf{q}_f \cdot \mathbf{n}_f)_h^{(1)} = \mu_-(-\mathbf{q}_f \cdot \mathbf{n}_f)_h^{(2)},$$

so that either (4) or (5) can be used in assembling the discrete fluxes in (2). This helps us to preserve compactness of the stencil for both cells $T_1$ and $T_2$ even with the multi-point fluxes (4), (5).

FV method with the nonlinear TPFA provides non-negativity of the discrete solution [3, 9], whereas FV method with the nonlinear MPFA provides the DMP [2, 11]. In the case of $\mathbb{K}$-orthogonal mesh vectors $\mathbb{K} \, \mathbf{n}_f$ and $\mathbf{t}_{12}$ are collinear, both nonlinear flux approximations reduce by construction to the conventional linear TPFA which provides at least first order accuracy. In general case, the linear TPFA may not pro-

**Fig. 2** Statement of the test case: two Dirichlet boundary conditions and full anisotropic diffusion tensor



vide approximation at all, whereas the linear MPFA may not provide the DMP or positivity. This statement is illustrated by an extended test case from [12], we consider all four schemes and two cases of Dirichlet boundary conditions: $G_0 = 0$, $G_1 = 2$ and $G_0 = 10$, $G_1 = 12$. Figure 2 presents the set up of the problem and Table 1 shows monotonicity and DMP violation by the schemes.

## 3 Applications

Means for the development of parallel numerical models of complex phenomena on general polyhedral meshes are provided by data structures and algorithms from the open source package Integrated Numerical Modelling Object-oriented Supercomputing Technologies (INMOST) [17]. FV discretization assumes that the processor possessing a mesh cell has access to data in neighboring cells. If a cell adjoins to the boundary of the local submesh associated with a processor, some of its neighbors belong to other processors. For each local submesh we generate additional layers of ghost cells composed of these neighbors. The ghost cells contain exact copy of data of the associated normal cells. The main difference between the ghost cell and the normal cell is that the ghost cell data should be actualized after any update of the normal cell data. Actualization involves inter-processor communications that move the data from normal cells to their ghost copies. Mesh data structure implemented in INMOST allows simple design of a numerical scheme on each mesh cell and is very convenient even for single processor implementations. Both applications presented in this paper are built using INMOST toolkit.

First we consider application of the nonlinear FV schemes for the black-oil model [12, 16]. The black oil model describes the three-phase flow of water, oil and gas components in the underground reservoir. If the reservoir pressure drops below certain threshold, then oil is split into a liquid phase and gaseous phase at thermodynamic equilibrium. In this case the water phase does not exchange mass with the other phases, while the liquid and the gaseous phases exchange mass. The model consists of mass conservation equations for each of the components and Darcy's velocity equations for each phase:

**Table 1** Minimum and maximum concentration values for the problem with the Dirichlet boundary conditions. Orthogonal grid with $h = 1/40$

| Scheme | $G_0 = 0,\ G_1 = 2$ | | $G_0 = 10,\ G_1 = 12$ | |
|---|---|---|---|---|
| | $C_{min}$ | $C_{max}$ | $C_{min}$ | $C_{max}$ |
| lin.TPFA | $1.3 \times 10^{-5}$ | 1.889 | 10.00 | 11.889 |
| nonl. TPFA | $1.6 \times 10^{-10}$ | 1.948 | **9.972** | 11.940 |
| MPFA | $\mathbf{-5.5 \times 10^{-2}}$ | **2.087** | 9.945 | **12.087** |
| nonl. MPFA | $1.2 \times 10^{-9}$ | 1.993 | 10.00 | 11.993 |

$$\mathbf{u}_\alpha = -\frac{k_{r\alpha}}{\mu_\alpha} \mathbb{K}\Big(\nabla p_\alpha - \rho_\alpha(p)\mathbf{g}\nabla z\Big), \quad \alpha = w, o, g, \tag{6}$$

where $\mathbb{K}$ is the absolute permeability tensor, $z$ is the depth, $\mathbf{g}$ is the gravity term, $p_\alpha$, $S_\alpha$ are *unknown* pressure and saturation, $\mu_\alpha$ and $k_{r\alpha}$ are the formation viscosity and relative phase permeability, $\rho_\alpha$ are the densities at current conditions for the phase $\alpha = w, o, g$.

We use the fully implicit scheme in time and Newton method to solve the nonlinear system at each time step. Construction of the Jacobian matrix is based on partial derivatives with respect to primary variables (oil pressure $p$, water and gas saturations $S_w$, $S_g$) of discrete Darcy fluxes. The latter are obtained either by the conventional linear TPFA or MPFA or by the nonlinear TPFA or MPFA presented above (the diffusion tensor should be replaced with absolute permeability tensor).

Dependence of the method coefficients on primary variables leads to the extension of the Jacobian stencil [12, 16]. For instance, in case of the nonlinear TPFA one has

$$-(\mathbb{K}\nabla p)_f^h \cdot \mathbf{n}_f = D_f^+(p)p_+ - D_f^-(p)p_-. \tag{7}$$

Coefficients $D_f^\pm$ must be differentiated as dependent on primary variables in neighboring cells: $\Delta D_p^\pm = \sum_{T_i \in \Sigma_{T_*}} L_{p,i}^\pm \Delta p_{T_i}$, where $\Sigma_{T*} = \Sigma_{T_+} \cup \Sigma_{T_-}$, $\Sigma_{T_\pm}$ is the set of cells forming the stencil for cell $T_\pm$, $L_{p,i}^\pm$ are the coefficients of differentiation. Wider stencil $\Sigma_{T*}$ for Jacobian results in more dense Jacobian matrix and more expensive Jacobian-vector multiplication and Jacobian preconditioning compared to the conventional linear TPFA. On the other hand, the linear TPFA is often inconsistent.

An example for three-phase water-flooding with several wells in heterogeneous media using nonlinear TPFA scheme is shown in Fig. 3.

The second application of the nonlinear FV schemes is related to validation of safe subsurface disposal of radioactive wastes (RW). In this application two main tasks must be solved, the groundwater (GW) flow problem and the transport in porous media problem, which may be strongly coupled in some cases. The novel FV schemes are implemented within the code Geomigration of Radionuclides (GeRa). This code is developed to model the major significant processes for radwaste disposal safety: saturated and unsaturated flow, density-driven flow, reactive transport with decay, heat transport. The basis for all these numerical models are the discretizations of the

**Fig. 3** Example of three-phase flow in heterogeneous media using nonlinear TPFA scheme. *Left* computational grid and geological layers. *Right* water saturation field



**Fig. 4** Example problem: groundwater flow in a realistic heterogeneous media. *Left* computational grid and geological layers. *Right* pressure head and flow streamlines

diffusion and advection operators. The computational meshes are assumed arbitrary polyhedral. The code involves the triangular prismatic and the octree-hexahedral mesh generators. In the first generator the resulting meshes may contain triangular prisms, tetrahedra and pyramids. The octree hexahedral generator cuts and adapts the cells to the domain boundary and interfaces between geological layers leading to complicated polyhedral cells.

The GW flow problem may be solved by FV scheme with either the linear TPFA (may be inconsistent) and MPFA (may be non-monotone) or the nonlinear TPFA and MPFA (both consistent and monotone). For the temporal discretization the operator-splitting scheme or the implicit scheme may be used. The first one treats the advection operator explicitly and the diffusion operator implicitly. Advection may be modeled using the conventional first-order accurate FV scheme with piecewise-constant solution or the second-order accurate TVD-scheme with linear reconstruction of discrete solution on the cells. For the diffusion operator any of the four flux approximation schemes (linear/nonlinear TPFA/MPFA) may be applied. The implicit scheme solves the coupled advection-diffusion problem using the nonlinear FV method for diffusion and local linear solution reconstruction for advection.

Numerical experiments with GeRa show robustness of the nonlinear schemes: the resulting matrices are reasonably well conditioned and the solutions remain non-negative or satisfy the DMP. In case of large complicated grids and heterogeneous tensor coefficients the schemes provide the best solution, as they allow to solve efficiently the generated grid equations and they are consistent.

Figure 4 (left) presents a filtration model with three geological layers, single well and outflow boundary with a prescribed water head. Water head solution and flow streamlines obtained using the FV scheme with the nonlinear TPFA is shown on Fig. 4 (right).

# References

1. Aavatsmark, I., Eigestad, G., Mallison, B., Nordbotten, J.: A compact multipoint flux approximation method with improved robustness. Num. Meth. Part. Diff. Equ. **24**(5), 1329–1360 (2008)
2. Chernyshenko, A.: Generation of adaptive polyhedral meshes and numerical solution of 2nd order elliptic equations in 3D domains and on surfaces. Ph.D. thesis, INM RAS, Moscow (2013)
3. Danilov, A., Vassilevski, Y.: A monotone nonlinear finite volume method for diffusion equations on conformal polyhedral meshes. Russ. J. Numer. Anal. Math. Model. **24**(3), 207–227 (2009)
4. Gao, Z.M., Wu, J.M.: A small stencil and extremum-preserving scheme for anisotropic diffusion problems on arbitrary 2D and 3D meshes. J. Comp. Phys. **250**, 308–331 (2013)
5. Kapyrin, I.V.: A family of monotone methods for the numerical solution of three-dimensional diffusion problems on unstructured tetrahedral meshes. Dokl. Math. **76**(2), 734–738 (2007)
6. Korotov, S., Křížek, M., Neittaanmäki, P.: Weakened acute type condition for tetrahedral triangulations and the discrete maximum principle. Math. Comp. **70**(233), 107–119 (2001)
7. Le Potier, C.: Schema volumes finis monotone pour des operateurs de diffusion fortement anisotropes sur des maillages de triangle non structures. C. R. Acad. Sci. Paris Ser. I **341**, 787–792 (2005)
8. Le Potier, C.: Finite volume scheme satisfying maximum and minimum principles for anisotropic diffusion operators. In: Eymard, R., Hérard, J.M. (eds.) Finite Volumes for Complex Applications, pp. 103–118. Wiley-ISTE, London (2008)
9. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: Interpolation-free monotone finite volume method for diffusion equations on polygonal meshes. J. Comp. Phys. **228**(3), 703–716 (2009)
10. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: A monotone finite volume method for advection-diffusion equations on unstructured polygonal meshes. J. Comp. Phys. **229**, 4017–4032 (2010)
11. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: Minimal stencil finite volume scheme with the discrete maximum principle. Russ. J. Numer. Anal. Math. Model. **27**(4), 369–385 (2012)
12. Nikitin, K., Terekhov, K., Vassilevski, Y.: A monotone nonlinear finite volume method for diffusion equations and multiphase flows. Comp. Geosci. (2013). doi:10.1007/s10596-013-9387-6
13. Nikitin, K., Vassilevski, Y.: A monotone nonlinear finite volume method for advectiondiffusion equations on unstructured polyhedral meshes in 3d. Russ. J. Numer. Anal. Math. Model. **25**(4), 335–358 (2010)
14. Nordbotten, J.M., Aavatsmark, I., Eigestad, G.T.: Monotonicity of control volume methods. Numer. Math. **106**(2), 255–288 (2007)
15. Stoyan, G.: On maximum principles for monotone matrices. Linear Algebra Appl. **78**, 147–161 (1986)
16. Terekhov, K., Vassilevski, Y.: Two-phase water flooding simulations on dynamic adaptive octree grids with two-point nonlinear fluxes. Russ. J. Numer. Anal. Math. Model. **28**(3), 267–288 (2013)
17. Vassilevski, Y., Konshin, I., Kopytov, G., Terekhov, K.: INMOST—program platform and graphic environment for development of parallel numerical models on general meshes (in Russian). Moscow university publishing, Moscow (2013)

18. Yuan, A., Sheng, Z.: Monotone finite volume schemes for diffusion equations on polygonal meshes. J. Comp. Phys. **227**(12), 6288–6312 (2008)
19. Yuan, G., Sheng, Z.: The finite volume scheme preserving extremum principle for diffusion equations on polygonal meshes. J. Comp. Phys. **230**(7), 2588–2604 (2011)

# Numerical Modelling of Viscous and Viscoelastic Fluids Flow in the Channel with T-Junction

**Radka Keslerová, Karel Kozel and David Trdlička**

**Abstract** In this work the numerical solution of the viscous and viscoelastic fluids flow for generalized Newtonian and Oldroyd-B fluids are considered. The governing system of equations is the system of generalized Navier-Stokes equations for incompressible laminar fluids flow. For the stress tensor on the right hand side of this system two different mathematical models for viscous and viscoelastic fluids flow are used, Newtonian model and Oldroyd-B model. For the numerical simulation of generalized Newtonian and Oldroyd-B fluids flow in the tested domain a cross model for viscosity function $\mu(\dot{\gamma})$ is considered. The finite volume method combined with the artificial compressibility method is used for the spatial discretization. For the time discretization the explicit multistage Runge-Kutta scheme is used. Computational domain is formed by the branched channel with one inlet and two outlet parts. The crosssection is square and the branch is perpendicular to the main pipe. The numerical results of generalized Newtonian and generalized Oldroyd-B fluids flow obtained by this method are presented.

## 1 Introduction

Branching of pipes occurs very often in many technical or biological applications. It is to be in human body in the complex branching system of blood vessels. Thus, this work is motivated by medical area of research. The blood can be characterized by shear-thinning property. In this work also the viscoelastic character is considered.

R. Keslerová (✉) · K. Kozel · D. Trdlička
CTU in Prague, Karlovo nám. 13, Praha, Czech Republic
e-mail: keslerov@marian.fsik.cvut.cz

K. Kozel
e-mail: Karel.Kozel@fs.cvut.cz

D. Trdlička
e-mail: David.Trdlicka@fs.cvut.cz

This leads to study of generalized Newtonian and Oldroyd-B fluids flow in the branched channel with T-junction.

## 2 Mathematical Model

The fundamental system of equations is the system of generalized Navier–Stokes equations for incompressible fluids. This system is based on the system of balance laws of mass and momentum for incompressible fluids.

$$\operatorname{div} \boldsymbol{u} = 0 \tag{1}$$

$$\rho \frac{\partial \boldsymbol{u}}{\partial t} + \rho (\boldsymbol{u}.\nabla)\boldsymbol{u} = -\nabla P + \operatorname{div} \mathsf{T} \tag{2}$$

where $P$ is the pressure, $\rho$ is the constant density, $\boldsymbol{u}$ is the velocity vector. The symbol $\mathsf{T}$ represents the stress tensor.

### 2.1 Stress Tensor

For the different choice of fluids model the different model of the stress tensor on the right hand side of the system of Navier-Stokes equations is used. For viscous flows with the representative of Newtonian fluids the simple model called *Newtonian model* is considered (see e.g. [1, 2])

$$\mathsf{T} = 2\mu \mathsf{D} \tag{3}$$

where $\mu$ is the dynamic viscosity and tensor $\mathsf{D}$ is the symmetric part of the velocity gradient.

In the case of viscoelastic fluids, the simplest viscoelastic model can be used. This model is denoted as *Maxwell model*

$$\mathsf{T} + \lambda_1 \frac{\delta \mathsf{T}}{\delta t} = 2\mu \mathsf{D} \tag{4}$$

where $\lambda_1$ is the relaxation time. The symbol $\frac{\delta}{\delta t}$ represents upper convected time derivative which is defined for general tensor by the relation (8).

By combination of these two presented models (Newtonian and Maxwell) the behaviour of mixture of viscous and viscoelastic fluids can be described. This model is called Oldroyd-B model and it has the form

$$\mathsf{T} + \lambda_1 \frac{\delta \mathsf{T}}{\delta t} = 2\mu \left( \mathsf{D} + \lambda_2 \frac{\delta \mathsf{D}}{\delta t} \right). \tag{5}$$

where symbols $\lambda_1$ and $\lambda_2$ are the relaxation and retardation time (with dimension of time).

The stress tensor $\mathsf{T}$ is decomposed to the Newtonian (viscous) part $\mathsf{T}_s$ and viscoelastic part $\mathsf{T}_e$, ($\mathsf{T} = \mathsf{T}_s + \mathsf{T}_e$). The tensor $\mathsf{T}_s$ is defined by Newtonian model (3) and the viscoelastic tensor $\mathsf{T}_e$ is defined by Maxwell model (4)

$$\mathsf{T}_s = 2\mu_s \mathsf{D}, \qquad \mathsf{T}_e + \lambda_1 \frac{\delta \mathsf{T}_e}{\delta t} = 2\mu_e \mathsf{D}, \tag{6}$$

where

$$\frac{\lambda_2}{\lambda_1} = \frac{\mu_s}{\mu_s + \mu_e}, \qquad \mu = \mu_s + \mu_e. \tag{7}$$

The upper convected derivative $\frac{\delta}{\delta t}$ used in the viscoelastic stress tensor is defined for the general tensor $\mathsf{M}$ by the relation, for more details see [2]

$$\frac{\delta \mathsf{M}}{\delta t} = \frac{\partial \mathsf{M}}{\partial t} + (\boldsymbol{u}.\nabla)\mathsf{M} - (\mathsf{WM} - \mathsf{MW}) - (\mathsf{DM} + \mathsf{MD}) \tag{8}$$

where $\mathsf{D}$ is symmetric part and $\mathsf{W}$ is antisymmetric part of the velocity gradient

$$\mathsf{D} = \frac{1}{2}(\nabla \boldsymbol{u} + \nabla \boldsymbol{u}^T) = \frac{1}{2}\begin{pmatrix} 2u_x & u_y + v_x & u_z + w_x \\ u_y + v_x & 2v_y & v_z + w_y \\ w_x + u_z & w_y + v_z & 2w_z \end{pmatrix} \tag{9}$$

and

$$\mathsf{W} = \frac{1}{2}(\nabla \boldsymbol{u} - \nabla \boldsymbol{u}^T) = \frac{1}{2}\begin{pmatrix} 0 & u_y - v_x & u_z - w_x \\ v_x - u_y & 0 & v_z - w_y \\ w_x - u_z & w_y - v_z & 0 \end{pmatrix}. \tag{10}$$

## 2.2 Generalizing: Cross Model

Both considered mathematical models (Newtonian and Oldroyd-B) could be generalized for numerical simulation of the blood flow. In this case the viscosity $\mu$ is no more constant but is defined by viscosity function according to the shear-thinning cross model (for more details see [8])

$$\mu(\dot{\gamma}) = \mu_\infty + \frac{\mu_0 - \mu_\infty}{(1 + (\lambda\dot{\gamma})^b)^a}, \quad \dot{\gamma} = 2\sqrt{\frac{1}{2}\text{tr}\,\mathsf{D}^2} \tag{11}$$

the following parameters have been used for the blood flow simulations presented in this paper: $\mu_0 = 1.6 \times 10^{-1}$ Pa s, $\mu_\infty = 3.6 \times 10^{-3}$ Pa s, $a = 1.23$, $b = 0.64$, $\lambda = 8.2$ s. The governing system of Eqs. (1), (2) is completed by the equation for the viscoelastic part of the stress tensor, therefore this system can be rewritten as follows

$$\text{div } \boldsymbol{u} = 0 \tag{12}$$

$$\rho \frac{\partial \boldsymbol{u}}{\partial t} + \rho(\boldsymbol{u}.\nabla)\boldsymbol{u} = -\nabla P + \text{div } \mathsf{T} \tag{13}$$

$$\mathsf{T} = \mathsf{T}_s + \mathsf{T}_e, \qquad \mathsf{T}_s = 2\mu(\dot{\gamma})\mathsf{D} \tag{14}$$

$$\frac{\partial \mathsf{T}_e}{\partial t} + (\boldsymbol{u}.\nabla)\mathsf{T}_e = \frac{2\mu_e}{\lambda_1}\mathsf{D} - \frac{1}{\lambda_1}\mathsf{T}_e + (\mathsf{W}\mathsf{T}_e - \mathsf{T}_e\mathsf{W}) + (\mathsf{D}\mathsf{T}_e + \mathsf{T}_e\mathsf{D}). \tag{15}$$

## 3 Numerical Solution

Numerical solution of the presented mathematical models is based on cell-centered finite volume method using explicit Runge–Kutta time integration. Steady state solution is achieved for $t \to \infty$. In this case the artificial compressibility method can be applied. It means that the continuity equation is completed by the time derivative of the pressure (for more details see e.g. [3–5, 7]). The system of equations (including the modified continuity equation) could be rewritten in the vector form.

$$\tilde{R}_\beta W_t + F_x^c + G_y^c + H_z^c = F_x^v + G_y^v + H_z^v + S, \quad \tilde{R}_\beta = \text{diag}(\frac{1}{\beta^2}, 1, \dots, 1) \tag{16}$$

where $\beta \in \mathbf{R}^+$, $W$ is vector of unknowns, $F^c, G^c, H^c$ and $F^v, G^v, H^v$ are inviscid and viscous fluxes and $S$ denotes the source term

$$W = \begin{pmatrix} p \\ u \\ v \\ w \\ t_1 \\ \vdots \\ t_6 \end{pmatrix}, \quad F^c = \begin{pmatrix} u \\ u^2 + p \\ uv \\ uw \\ ut_1 \\ \vdots \\ ut_6 \end{pmatrix}, \quad G^c = \begin{pmatrix} v \\ uv \\ v^2 + p \\ vw \\ vt_1 \\ \vdots \\ vt_6 \end{pmatrix}, \quad H^c = \begin{pmatrix} w \\ uw \\ vw \\ w^2 + p \\ wt_1 \\ \vdots \\ wt_6 \end{pmatrix}, \tag{17}$$

$$F^v = \begin{pmatrix} 0 \\ 2\mu(\dot{\gamma})u_x \\ \mu(\dot{\gamma})(u_y + v_x) \\ \mu(\dot{\gamma})(u_z + w_x) \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad G^v = \begin{pmatrix} 0 \\ \mu(\dot{\gamma})(u_y + v_x) \\ 2\mu(\dot{\gamma})v_y \\ \mu(\dot{\gamma})(v_z + w_y) \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad H^v = \begin{pmatrix} 0 \\ \mu(\dot{\gamma})(u_z + w_x) \\ \mu(\dot{\gamma})(v_z + w_y) \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \tag{18}$$

$$S = \begin{pmatrix} 0 \\ t_{1x} + t_{2y} + t_{3z} \\ t_{2x} + t_{4y} + t_{5z} \\ t_{3x} + t_{5y} + t_{6z} \\ 2\frac{\mu_e}{\lambda_1}u_x - \frac{t_1}{\lambda_1} + 2u_x t_1 + (u_y + v_x)t_2 + (u_z + w_x)t_3 \\ \frac{\mu_e}{\lambda_1}(u_y + v_x) - \frac{t_2}{\lambda_1} + u_y t_4 + u_z t_5 + u_y t_1 + w_y t_3 + u_x t_2 + v_y t_2 \\ \frac{\mu_e}{\lambda_1}(u_z + w_x) - \frac{t_3}{\lambda_1} + u_y t_5 + u_z t_6 + u_z t_1 + v_z t_2 + u_x t_3 + w_z t_3 \\ 2\frac{\mu_e}{\lambda_1}v_y - \frac{t_4}{\lambda_1} + (u_y + v_x)t_2 + 2v_y t_4 + (v_z + w_y)t_5 \\ \frac{\mu_e}{\lambda_1}(v_z + w_y) - \frac{t_5}{\lambda_1} + v_x t_3 + v_z t_6 + u_z t_2 + v_z t_4 + v_y t_5 + w_z t_5 \\ 2\frac{\mu_e}{\lambda_1}w_z - \frac{t_6}{\lambda_1} + (u_z + w_x)t_3 + (v_z + w_y)t_5 + 2w_z t_6 \end{pmatrix}. \quad (19)$$

Equation (16) is discretized in space by the finite volume method and the arising system of ODEs is integrated in time by the multistage Runge–Kutta scheme ([5, 6])

$$\begin{aligned} W_i^n &= W_i^{(0)} \\ W_i^{(s)} &= W_i^{(0)} - \alpha_{s-1}\Delta t \mathcal{R}(W)_i^{(s-1)} \\ W_i^{n+1} &= W_i^{(M)} \qquad s = 1, \dots, M, \end{aligned} \quad (20)$$

whith $M = 3$, $\alpha_0 = \alpha_1 = 0.5$, $\alpha_2 = 1.0$, the steady residual $\mathcal{R}(W)_i$ is defined by finite volume method as

$$\mathcal{R}(W)_i = \frac{1}{C_i} \sum_{k=1}^{6} \left[ \left( \overline{F}_k^c - \overline{F}_k^v \right) \Delta Sx_k - \left( \overline{G}_k^c - \overline{G}_k^v \right) \Delta Sy_k - \left( \overline{H}_k^c - \overline{H}_k^v \right) \Delta Sz_k \right] + \overline{S}, \quad (21)$$

where $C_i$ is the volume of the primary grid cell. The symbols $\overline{F}_k^c, \overline{G}_k^c, \overline{H}_k^c$ and $\overline{F}_k^v, \overline{G}_k^v, \overline{H}_k^v$ denote the numerical approximation of the inviscid and viscous physical fluxes. The symbol $\overline{S}$ denotes the numerical approximation of the source term. The symbols $\Delta Sx_k$, $\Delta Sy_k$ and $\Delta Sz_k$ respectively represent the volume of the $k$th-surface of the primary cell in the $x$, $y$, $z$ direction.

The inviscid numerical fluxes are computed as an aritmetic average of the inviscid numerical fluxes of two neighbouring finite volume cells

$$\overline{F}_k^c = \frac{1}{2}\left( F_i^c + F_k^c \right), \quad \overline{G}_k^c = \frac{1}{2}\left( G_i^c + G_k^c \right), \quad \overline{H}_k^c = \frac{1}{2}\left( H_i^c + H_k^c \right) \quad (22)$$

where index $i$ denotes the index of the primary cell and the index $k$ is the index of the neighboring cells. The mesh in the computational domain is assumed structured with hexahedral cells. The variables of the source term are computed in the same meaning.

In the definition of the viscous fluxes there are the partial derivatives of velocity with respect to the spatial coordinates $x$, $y$, $z$. The numerical approximation of these derivatives need to be discretized. Integrating these over a dual cell and using Green's theorem results in, e.g. for $\overline{u}_x$

$$\bar{u}_x = \frac{1}{C_k} \sum_{m=1}^{8} u_m \Delta S x_m, \tag{23}$$

where $C_k$ is the volume of the dual cell corresponding to the $k$-th surface of the basic cell, the shape of dual cells is octahedron. The symbol $\Delta S x_m$ is the volume of the $m$th-surface of the dual cell in the $x$ direction. Similarly for other velocity derivatives.

*Steady Boundary Conditions.* The flow is modelled in a bounded computational domain where a boundary is divided into three mutually disjoint parts: a solid wall, an outlet and an inlet. At the inlet Dirichlet boundary condition for velocity vector (parabolic profile) is used and for a pressure and the stress tensor Neumann boundary condition is used. At the outlet parts the pressure value is given and for the velocity vector and the stress tensor Neumann boundary condition is used. The homogeneous Dirichlet boundary condition for the velocity vector is used on the wall. For the pressure and stress tensor Neumann boundary condition is considered.

## 4 Numerical Results

This section deals with the comparison of the numerical results of generalized Newtonian and generalized Oldroyd-B fluids flow. Numerical tests are performed in an idealized branched channel with the square cross-section. Figure 1 (left) shows the shape of the tested domain. The computational domain is discretized using a structured, wall fitted mesh with hexahedral cells. The domain is divided to four blocks,

| | | |
|---|---|---|
| Block #1 | $60 \times 30 \times 30$ cells | black |
| Block #2 | $30 \times 30 \times 30$ cells | red |
| Block #3 | $40 \times 30 \times 30$ cells | blue |
| Block #4 | $30 \times 40 \times 30$ cells | green |

As initial condition the following model parameters are used: $\mu_e = 0.001$ Pa s, $\mu_s = 0.009$ Pa s, $\lambda_1 = 0.06$ s, $U_0 = 0.1\,\text{ms}^{-1}$, $L_0 = 0.01$ m, $\rho = 1000\,\text{kg m}^{-3}$ In the outlet parts the pressure is given by values: 0.0005 Pa (main channel) and 0.00025 Pa (branch). Using these data, fully developed Poiseuille velocity profile (for Newtonian fluid) is prescribed at the inlet (Dirichlet condition). At the outlet homogeneous Neumann conditions for the velocity components and a constant pressure are prescribed. On the vessel walls no-slip homogeneous Dirichlet conditions are prescribed for the velocity field. In the case of the Oldroyd-B and generalized Oldroyd-B models, homogeneous Neumann conditions are imposed for the components of the extra stress tensor at all boundaries. In Fig. 1 (right) the axial velocity profile for fully developed flow close to the branching is shown. The lines for

**Fig. 1** Structure of the computed domain (*left*) and axial velocity profile for steady fully developed flow of tested fluids (*right*). **a** Structure of the domain. **b** Axial velocity profile



**Fig. 2** Velocity isolines of steady flows for generalized Newtonian and Oldroyd-B fluids. **a** Newtonian. **b** Generalized Newtonian. **c** Oldroyd-B. **d** Generalized Oldroyd-B

Newtonian and Oldroyd-B fluids are similar to the parabolic line, as was assumed. From this velocity profile is clear that the shear thinning fluids attain lower maximum velocity in the central part of the channel (close to the axis of symmetry) which is compensated by the increase of local velocity in the boundary layer close to the wall. In Fig. 2 the velocity isolines and the cuts through the main channel and the small branch are shown.

The axial velocity isolines for all tested fluids (Newtonian, generalized Newtonian, Oldroyd-B and generalized Oldroyd-B) are shown in the Fig. 3. It can be

**Fig. 3** Axial velocity isolines in the center-plane area. **a** Newtonian. **b** Generalized Newtonian. **c** Oldroyd-B. **d** Generalized Oldroyd-B

observed from Fig. 3 that the size of separation region for generalized Newtonian and generalized Oldroyd-B fluids is smaller than for Newtonian and Oldroyd-B fluids.

## 5 Conclusion

In this paper a finite volume solver for incompressible laminar viscous and viscoelastic fluids flow in the branching channel with T-junction was described. Used mathematical models (Newtonian and Oldroyd-B) were generalized by the shear-thinning cross model for numerical solution of generalized Newtonian and Oldroyd-B fluids flow. These types of flow were numerical modelled in three dimensional domain with square cross-section.

## References

1. Bodnar, T., Sequeira, A.: Numerical study of the significance of the non-newtonian nature of blood in steady flow through stenosed vessel. Adv. Math. Fluid Mech. 83–104 (2010). http://www.springer.com/mathematics/applications/book/978-3-642-04067-2
2. Bodnar, T., Sequeira, A., Prosi, M.: On the shear-thinning and viscoelastic effects of blood flow under various flow rates. Appl. Math. Comput. **217**, 5055–5067 (2010)
3. Chorin, A.J.: A numerical method for solving incompressible viscous flow problem. J. Computat. Phys. **135**, 118–125 (1967)
4. Dvořák, R., Kozel, K.: Mathematical Modelling in Aerodynamics (in Czech). CTU, Praha (1996)

5. Keslerová, R., Kozel, K.: Numerical modelling of incompressible flows for Newtonian and non-Newtonian fluids. Math. Comput. Simul. **80**, 1783–1794 (2010)
6. LeVeque, R.: Finite-Volume Methods for Hyperbolic Problems. Cambridge University Press, Cambridge (2004)
7. Louda, P., Kozel, K., Příhoda, J., Beneš, L., Kopáček, T.: Numerical solution of incompressible flow through branched channels. Comput. Fluids **46**, 318–324 (2011)
8. Vimmr, J., Jonášová, A.: Non-Newtonian effects of blood flow in complete coronary and femoral bypasses. Math. Comput. Simul. **80**, 1324–1336 (2010)

# Gradient Evaluation on a Quadtree Based Finite Volume Grid

**Zuzana Krivá, Angela Handlovičová and Karol Mikula**

**Abstract**  Many problems described by nonlinear PDEs need good approximations of gradients on finite volumes. Using finite volume methods, this can be difficult task if discretization of a computational domain does not fulfill the classical orthogonality property. Such a situation can occur, e.g., during coarsening in image processing using quadtree grids. We present a construction of an adjusted quadtree grid for which the connection of representative points of two adjacent finite volumes is perpendicular to their common boundary. On the other hand, for such an adjusted grid, the intersection of representative points connection with a finite volume boundary is not a middle point of their common edge. In this paper we present a new method of gradient evaluation for such a situation.

## 1  The Computational Grid

In this section we introduce our finite volume computational grid, its construction and its properties. Our purpose is to build the grid using large elements for regions with homogeneous values of a solution function—in our experiment representing image intensities. To this purpose we first build a graded quadtree, i.e. the quadtree, in which the difference in a level between adjacent cells is constrained, in our case to one. Grids associated with such trees are often used in order to produce procedures that are easier to implement. Moreover, in our case it is an inevitable requirement to

Z. Krivá (✉) · A. Handlovičová · K. Mikula
Department of Mathematics, Faculty of Civil Engineering, Slovak University of Technology in Bratislava, Radlinského 11, 813 68 Bratislava, Slovak Republic
e-mail: kriva@math.sk

A. Handlovičová
e-mail: angela.handlovicova@stuba.sk

K. Mikula
e-mail: mikula@math.sk

**Fig. 1** An example of the original quadtree grid together with the representative points of its elements (*on the left*). This grid is transformed into the consistent one (*on the right*)



be able to adjust the quadtree to the consistent finite volume grid. The consistent grid possesses the important property that the connection of two representative points of two adjacent finite volumes is perpendicular to their common boundary, which is an important fact when we use the classical finite volume discretization [2]. An example of a quadtree and a corresponding consistent grid is displayed in Fig. 1.

**Building the quadtree.** Let us suppose that our data is given on a regular non-adaptive square grid (which corresponds e.g. to the pixel structure of an image). First we build the quadtree by merging the elements with similar values from the smaller cells to the larger cells, i.e. from leaves to the root. The old values are either unchanged, or replaced by averaging the values from the processed area. During this process, the information about successful or unsuccessful merging is stored in a binary field with the size corresponding to the image. Moreover, this information is stored in such a way that it enables us to create a graded quadtree with a prescribed ratio of elements. It can be also used as a stopping criterion during *traversing* the quadtree and to test the configurations of elements—the leaves of the quadtree.

As we have already mentioned, in order to simplify creating the linear system matrix, where access to neighbors is needed, and to enable creating the consistent grid, we require that the ratio of sides of two adjacent squares is 1:1, 1:2 or 2:1. The used technique of building the quadtree adaptive grids is described in [4]. It uses the following *coarsening criterion*: the cells are merged if a difference in their intensities is below a prescribed tolerance $\varepsilon$.

**Adjustment to the quadtree based consistent grid.** The quadtree grid (Fig. 1 left) is *inconsistent* in the sense, that we cannot find the unique representative points of the adjacent grid elements—finite volumes—such that the connection of their representative points is perpendicular to their common boundary. The adaptive grid fulfilling this condition is called *consistent* and it is an *admissible* mesh in the sense of [2]. However, the basic quadtree grid can be adjusted to a consistent one procedurally: we must adjust the shape, if two adjacent finite volumes $p$ and $q$ are of different size. If we denote the length of a common edge in the original quadtree by $h$ and we shift the "hanging node" by $v = \frac{h}{3}$ (e.g. in Fig. 2 we shift $X$ to $X'$), then the connection of representative points is perpendicular to the shifted common boundary. This fact (and also the fact that $\frac{BX'}{PQ} = \frac{2}{3}$) follows from the similarity of triangles $\triangle AQP$ and $\triangle XX'B$ with the ratio of their adjacent sides 1:3. The area of $p$ is also evaluated procedurally—it depends on a configuration of its neighbors.

**Fig. 2** Adjustment to the consistent grid. $|XX'| = v = \frac{1}{3}h$. XB$=\frac{2}{3}$PA, hence $\frac{BX'}{PQ} = \frac{2}{3}$. Examples of the shapes where the intersection of the connection of representative points and a common edge $\sigma$ is not the midpoint of $\sigma$

**Notations.** Let every finite volume $p$ of measure $|p|$ have a representative point $X_p$ lying in its center or in the center of the original square for an adjusted element of the consistent grid. The common interface of $p$ and $q$ is a line segment—an edge $\sigma_{pq}$ with a nonzero measure in $\mathbf{R}$ denoted by $|\sigma_{pq}|$ and $d_{pq} = |X_q - X_p|$ is the distance of representative points. Let us denote by $X_\sigma$ auch a point of $\sigma_{pq}$, which represents the *intersection of the line segment $X_p X_q$ and $\sigma_{pq}$*. In our consistent grid, $X_p X_q$ is perpendicular to $\sigma$, but the intersection $X_\sigma$ is not the midpoint of $\sigma$ in the general case. Let us denote by $X_\sigma^*$ the *midpoint* of the edge $\sigma$. By $\mathscr{E}_p$ we denote the set of all edges $\sigma$ of $p$. When we speak about a unit outer normal vector to $\sigma \in \mathscr{E}_p$, we denote it by $\mathbf{n}_{p\sigma}$.

## 2 Approximation of the Gradient on the Consistent Grid

Our method for evaluation of gradients on finite volumes is based on [3]. Such a method works locally in that sense that we consider also representative points on finite volume edges, but not values at the corners. Then, with a help of these points we only need access to neighbors sharing a common edge, which is important when working on adaptive grids.

When solving PDEs where nonlinearities depend on the solution gradient, the method from [3] works as follows:

1. for edges $\sigma$ of a finite volume $p$ we define representative points $X_\sigma^*$—their midpoints, it must hold $X_\sigma^* = X_\sigma$,
2. with a help of these points, we evaluate the norm of gradient on $p$ locally using the consequence of the Stokes formula, see (3)–(4),
3. discrete equation for the finite volume $p$ is derived locally,
4. values of solution in $X_\sigma^*$ are obtained by using conservation principle.

In the consistent adaptive grid $X_\sigma^* \neq X_\sigma$ in general. Such a situation occurs on edges containing a hanging node in the original quadtree grid. The most critical shape in this sense is the sharp element where $X_\sigma$ is not the midpoint on any of the edges (Fig. 2 right).

Let us suppose the linear approximation of the solution over the finite volume $p$. At $X \in p$ any linear function can be written as

$$u(X) = u(X_p) + \nabla u \cdot (X - X_p) = u_p + \nabla u \cdot (X - X_p). \tag{1}$$

If $X = X_\sigma$ it holds

$$u_\sigma - u_p = \nabla u \cdot (X_\sigma - X_p), \tag{2}$$

where $u_\sigma$, $u_p$ represent values of the solution at points $X_\sigma$ and $X_p$. The gradient of the linear function is a constant vector in $\mathbf{R}^2$, thus also over a control volume $p$. It will be denoted by $\nabla u$. Then it holds

$$\nabla u = \frac{1}{|p|} \int_p \nabla u \, dX = \frac{1}{|p|} \int_{\partial p} u \mathbf{n}_p \, dS = \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} \int_\sigma (u_p + \nabla u \cdot (X - X_p)) \mathbf{n}_{p\sigma} \, dS$$

$$= \frac{1}{|p|} u_p \sum_{\sigma \in \mathscr{E}_p} |\sigma| \mathbf{n}_{p\sigma} + \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma| \nabla u \cdot (X_\sigma^* - X_p) \mathbf{n}_{p\sigma}. \tag{3}$$

The term $\sum_{\sigma \in \mathscr{E}_p} |\sigma| \mathbf{n}_{p\sigma} = \mathbf{0}$ and the expression $|\sigma| \nabla u (X_\sigma^* - X_p) \mathbf{n}_{p\sigma}$ represents the precise integration of a linear function over the edge $\sigma$. Thus we have

$$\nabla u = \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma| \nabla u \cdot (X_\sigma^* - X_p) \mathbf{n}_{p\sigma}. \tag{4}$$

On the edges, where $X_\sigma \neq X_\sigma^*$, we can express

$$X_\sigma^* - X_p = (X_\sigma - X_p) + (X_\sigma^* - X_\sigma). \tag{5}$$

Then $\nabla u$ can be split into two parts

$$\nabla u = \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma| \nabla u \cdot (X_\sigma - X_p) \mathbf{n}_{p\sigma} + \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma| \nabla u \cdot (X_\sigma^* - X_\sigma) \mathbf{n}_{p\sigma}. \tag{6}$$

The part of $\nabla u$ given by the first term of (6) will be denoted as $(\nabla u)^A$ and due to (2) it can be evaluated as

$$(\nabla u)^A = \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|(u_\sigma - u_p)\mathbf{n}_{p\sigma}. \tag{7}$$

The second term of (6) is a *correction* of $(\nabla u)^A$ and it depends on the unknown gradient.

## 2.1 Evaluation of the Gradients with Corrections

In the following text we use subscripts in two ways: if they represent derivatives, we use $x$ or $y$ and if they represent the vector components, we use 1 or 2. Let us denote the correction vector $(X_\sigma^* - X_\sigma)$ by $\mathbf{c}_\sigma = ((c_\sigma)_1, (c_\sigma)_2)$. We will work with $(\nabla u)^A = ((u_x)^A, (u_y)^A)$, $\mathbf{n}_{p\sigma} = ((n_{p\sigma})_1, (n_{p\sigma})_2)$ and the unknown vector $\nabla u = (u_x, u_y)$. Now (6) can be rewritten into the form

$$(u_x, u_y) = ((u_x)^A, (u_y)^A) + \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|((c_\sigma)_1 u_x + (c_\sigma)_2 u_y)((n_{p\sigma})_1, (n_{p\sigma})_2).$$
$$\tag{8}$$

We see that (8) represents the linear system of two equations with two unknowns $u_x$ and $u_y$ which can be adjusted to the following form:

$$u_x \left(1 - \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|(c_\sigma)_1(n_{p\sigma})_1\right) + u_y \left(-\frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|(c_\sigma)_2(n_{p\sigma})_1\right) = (u_x)^A,$$

$$u_x \left(-\frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|(c_\sigma)_1(n_{p\sigma})_2\right) + u_y \left(1 - \frac{1}{|p|} \sum_{\sigma \in \mathscr{E}_p} |\sigma|(c_\sigma)_2(n_{p\sigma})_2\right) = (u_y)^A.$$

We rewrite the system into such a form that we can see that the coefficient matrix denoted by $B$ depends only on the shape of a grid element, but not on its size (level). Let us denote: $\mathbf{N}_{p\sigma} = \frac{|\sigma|\mathbf{n}_{p\sigma}}{l}$ and $\mathbf{C}_\sigma = \frac{\mathbf{c}_\sigma}{l}$, where $l$ is the edge length of the square in the non adjusted quadtree. We have:

$$u_x \left(1 - \frac{l^2}{|p|} \sum_{\sigma \in \mathscr{E}_p} (C_\sigma)_1(N_{p\sigma})_1\right) \quad + \quad u_y \left(-\frac{l^2}{|p|} \sum_{\sigma \in \mathscr{E}_p} (C_\sigma)_2(N_{p\sigma})_1\right) \quad = (u_x)^A,$$
$$\tag{9}$$

$$u_x \left(-\frac{l^2}{|p|} \sum_{\sigma \in \mathscr{E}_p} (C_\sigma)_1(N_{p\sigma})_2\right) \quad + \quad u_y \left(1 - \frac{l^2}{|p|} \sum_{\sigma \in \mathscr{E}_p} (C_\sigma)_2(N_{p\sigma})_2\right) = (u_y)^A.$$

**Fig. 3** *Example 1 Left* the consistent quadtree grid with the inspected element. For $u_2(X) = \frac{1}{3}(x^3 + y^3)$ we compare $(\nabla u)_p$ and $(\nabla u)_p^A$ with the values of the gradient evaluated analytically in $(x_p, y_p)$. The values are scaled with darker values representing larger differences. *Middle* $||(\nabla u)^A| - |\nabla u_{exact}||$. *Right* $||\nabla u| - |\nabla u_{exact}||$

The elements of the coefficient matrix in (9) can be evaluated procedurally traversing the quadtree, or we can construct $B$ using its properties mentioned later. $B$ can be also precalculated in advance for every shape (there is only limited number of shapes in the consistent quadtree grid)—we can store $B^{-1}$ and evaluate $\nabla u = B^{-1}(\nabla u)^A$.

*Example 1* Let us take the consistent quadtree grid built over a uniform grid with $32 \times 32$ elements (Fig. 3 left). We inspect specific functions defined on $[-1.25, 1.25] \times [-1.25, 1.25]$: we consider the norm of the gradient evaluated analytically, the norm of $(\nabla u)^A$ and $\nabla u$ obtained by solving (9). First let us take the function $u_1(X) = \frac{1}{2}(x^2 + y^2)$. We take the sharp marked element (Fig. 3 left) with the representative point $(x_p, y_p) = (0.742, -0.89)$. First $u_\sigma$ is set to the exact value evaluated using $u(X)$. The approximated gradient—the vector $(\nabla u)^A$ evaluated without correction is equal to $(-1, 711, -1.801)$. After correction using (9) it is equal to $(0.860, -1.03)$, while the analytical gradient at this point has the value $(x_p, y_p)$ given above. In practical tasks, $u_\sigma$ is obtained by an interpolation. Thus we consider also that $u_\sigma$ is obtained by a linear interpolation between $u_p$ and $u_q$, its neighbor. It is interesting that in such case the approximated gradient of the quadratic function $u_1(X)$ obtained by (9) is equal to the analytical one. Secondly, let us take the function $u_2(X) = \frac{1}{3}(x^3 + y^3)$, the selected volume as in the previous case and $u_\sigma$ obtained by a linear interpolation. The analytical value of the gradient is $(0.551, 0.807)$, using (7) we get $(\nabla u)^A = (0.813, 0.987)$ and using (9) $\nabla u = (0.5572, 0.813)$. Figure 3 depicts differences of norms of $(\nabla u)_p^A$ and analytical gradient evaluated in the representative points of grid elements $(x_p, y_p)$ (middle) and the norms of $(\nabla u)_p$ obtained by (9) and the analytical gradient (right) for the function $u_2(X)$ in $(x_p, y_p)$. At the end we explored $L_2$ norms of errors $|\nabla u| - |\nabla u_{exact}|$ evaluated on four consistent adaptive grids obtained by consequent refinement of the grid from Fig. 3: every finite volume of a corresponding quadtree grid was divided into four subvolumes and afterwards the grid was adjusted to the consistent one. We have obtained following results: 0.0619, 0.0173, 0. 0051 and 0.00158.

*Properties of the coefficient matrix $B$*. The nonzero corrections occur only if one of the edgepoints of $\sigma$ is the shifted node. Let the edge vector $\boldsymbol{\sigma}$ be oriented from the shifted node to the quadtree corner. It can be shown that:

1. on the aligned edge $\sigma$, the correction $\mathbf{c}_\sigma$ can be expressed like $\mathbf{c}_\sigma = \frac{\sigma}{10}$, on the vertical or horizontal edge $\mathbf{c}_\sigma = \frac{\sigma}{4}$,

2. $\sum\limits_{\sigma \in \mathscr{E}_p} (C_\sigma)_1 (N_{p\sigma})_1 = - \sum\limits_{\sigma \in \mathscr{E}_p} (C_\sigma)_2 (N_{p\sigma})_2$,

3. $\sum\limits_{\sigma \in \mathscr{E}_p} (C_\sigma)_1 (N_{p\sigma})_2 = \sum\limits_{\sigma \in \mathscr{E}_p} (C_\sigma)_2 (N_{p\sigma})_1$,

4. It holds that the matrix $B$ is regular ($det(B) > 0$) and the system (9) has always a unique solution. It can be proved using properties 1, 2 and 3.

# 3 Numerical Solution of the Regularized Perona-Malik Equation on the Consistent Adaptive Grid

In this section we present one experiment—solution of the regularized Perona-Malik equation [1] on a rectangular domain $\Omega \subset \mathbf{R}^2$ discretized with help of a consistent adaptive grid. The scaling interval $I = [0, T]$ is discretized into scale steps with $t^n = t^n + \tau$, $\tau$ is the scale step size, on the boundaries we keep the zero Neumann boundary conditions. So we solve the problem

$$\partial_t u - \nabla \cdot (g(|\nabla G_s * u|)\nabla u) = 0, \quad \text{in } Q_T \equiv I \times \Omega, \tag{10}$$

where $g(s) = \frac{1}{1+Ks^2}$, $K > 0$ is the Perona-Malik function slowing down the diffusion in the vicinity of edges and $G_s(x)$ is the smoothing kernel. In our algorithm we realize the convolution $\nabla(G_s * u) = G_s * \nabla u$ by solving the linear heat equation. We apply one or several steps of the adaptive scheme for time $T_s$ corresponding to $s$ to both $x$ and $y$ coordinates of the gradient, then we evaluate the norm of the gradients and apply the Perona-Malik function $g$ to get the diffusion coefficient denoted by $g_p^{s,n-1}$.

Let us denote by $u_\sigma^n$ the value of the solution in $X_\sigma$ at the time step $t^n$. The derivative in the direction $\mathbf{n}_{p\sigma}$ is approximated by $\nabla u^n \cdot \mathbf{n}_{p\sigma} \approx \frac{\left(u_\sigma^n - u_p^n\right)}{d_{p\sigma}}$. The diffusion coefficient $g_p^{s,n-1}$ is constant all over $p$, thus the flux over $\sigma$ can be approximated by

$$F_{p\sigma}^n = g_p^{s,n-1} \frac{|\sigma|}{d_{p\sigma}} \left(u_\sigma^n - u_p^n\right). \tag{11}$$

A good way to evaluate $\frac{|\sigma|}{d_{p\sigma}}$ is to consider the neighbor $q$ sharing $\sigma$ with $p$. Then we can express (11) with a help of the transmissivity coefficient $T_{pq} = \frac{|\sigma|}{d_{pq}}$ and the ratio of $d_{p\sigma}$ and $d_{q\sigma}$, where $d_{p\sigma}$ and $d_{q\sigma}$ are distances of representative points from $X_\sigma$.

**Fig. 4** *Numerical experiment* The artificial noisy image, the filtered image and the fixed adaptive grid

If $\sigma \perp X_p X_q$ in the non adjusted grid, $\frac{d_{p\sigma}}{d_{q\sigma}} = 1$, otherwise, $\frac{d_{p\sigma}}{d_{q\sigma}} = \frac{4}{1}$ or $\frac{1}{4}$. For $T_{pq}$ it holds that if one edgepoint of $\sigma$ is a hanging node in the nonadjusted quadtree, then $T_{pq} = \frac{2}{3}$, otherwise $T_{pq} = 1$. The approximated flux (11) can be expressed as

$$F_{p\sigma}^n = T_{pq} \left( 1 + \frac{d_{q\sigma}}{d_{p\sigma}} \right) g_p^{s,n-1} \left( u_\sigma^n - u_p^n \right). \tag{12}$$

Now we solve the linear system, where the set of equations for all finite volumes $p$

$$(u_p^n - u_p^{n-1}) \, |p| = \tau \sum_{\sigma \in \mathscr{E}_p} F_{p\sigma}^n \tag{13}$$

is accompanied by a set of equations for every $u_\sigma^n$, $\sigma \in \mathscr{E}_p$, obtained from the relationship $F_{p\sigma}^n = -F_{q\sigma}^n$ resulting in

$$u_\sigma^n = \frac{d_{q\sigma} g_p^{s,n-1} u_p^n + d_{p\sigma} g_q^{s,n-1} u_q^n}{d_{q\sigma} g_p^{s,n-1} + d_{p\sigma} g_q^{s,n-1}}.$$

We present here a numerical experiment where we begin with a regular grid and continue to use it until the decrease of elements is sufficient. Then we run the adaptive algorithm on the same adaptive grid. Advantage of this approach is that for the fixed adaptive grid we can store all necessary information, e.g. configurations of neighbors, matrix $B$, etc. We consider the image of the size $128 \times 128$ disturbed by the additive noise. We performed 13 scale steps with $\tau = 1$, with $K = 1000$ in the Perona-Malik function $g$ and the time of presmoothing $T_s = 0.6$. The number of grid elements was reduced to $\frac{1}{3}$ after 5 scale steps, and then we continued on the fixed grid. The parameter $\varepsilon$ used in the coarsening criterion is set to 0.01. Figure 4 shows the data itself, the filtered data and the adaptive grid fixed after 5 scale steps.

# References

1. Catté, F., Lions, P., Morel, J., Coll, T.: Image selective smoothing and edge detection by nonlinear diffusion. SIAM J. Numer. Anal. **129**, 182–193 (1992)
2. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. Handb. Numer. Anal. **7**, 713–1018 (2000)
3. Eymard, R., Handlovičová, A., Mikula, K.: Study of a finite volume scheme for the regularized mean curvature flow level set equation. IMA J. Numer. Anal. **31**(3), 813–846 (2010)
4. Krivá, Z., Mikula, K.: An adaptive finite volume scheme for solving nonlinear diffusion equations in image processing. J. Visual Commun. Image Represent. **13**, 22–35 (2002)

# 3D Lagrangian Segmentation with Simultaneous Mesh Adjustment

**Karol Mikula and Mariana Remešíková**

**Abstract**  We present a method for 3D image segmentation based on the Lagrangian approach. The segmentation model is a 3D analogue of the geodesic active contour model [1] and it contains an additional tangential movement term that allows us to control the quality of the mesh during the evolution process. The model is discretized by the finite volume approach. Segmentation of zebrafish cell images is shown to illustrate the performance of the method.

## 1 Introduction

A large number of existing 3D image segmentation techniques are based on PDE models representing evolution of 2D surfaces in 3D. Most of them use the level set approach due to its favorable properties with respect to possible topological changes. The other alternative is the Lagrangian approach that directly evolves a 2D surface without viewing it as an isosurface of a three-dimensional function. Because of its two-dimensional character, this technique offers a possibility to obtain faster algorithms. However, even if we do not have to deal with any topological changes in the course of the computation, a Lagrangian method can face the problem of mesh deterioration as a discretized surface evolves. Therefore, in order to successfully apply such methods, we need to have at disposal a mechanism for controlling the quality of the surface discretization during the computation.

K. Mikula · M. Remešíková (✉)
Faculty of Civil Engineering, Department of Mathematics and Descriptive Geometry, Slovak University of Technology, Radlinského 11, 81368 Bratislava, Slovakia
e-mail: remesikova@math.sk

K. Mikula
e-mail: mikula@math.sk

Our paper presents a Lagrangian method for 3D image segmentation that allows to adjust the mesh quality along with the surface evolution. The segmentation model contains two normal movement components—one is given by the gradient of an image edge detector function and the other one depends on the edge detector itself and the mean curvature of the evolving surface. An additional tangential velocity term is added in order to be able to redistribute the mesh points during the evolution. The corresponding PDE is discretized by a finite volume technique and the redistribution is designed so that all control volumes have the same area for $t \to \infty$. The performance of the method is illustrated by examples using microscope images of zebrafish cells.

## 2 The Segmentation Model

Let $I : \mathbb{R}^3 \supset \Omega \to \mathbb{R}$ be an image intensity function. There are several possibilities how to detect the edges in the image; one of them is to use the edge detector function $e : \Omega \to \mathbb{R}$ of the form

$$e(x, y, z) = \frac{1}{1 + K \|\nabla I(x, y, z)\|^2} \tag{1}$$

where $K$ is a positive real constant.

Now let $X$ be a two-dimensional Riemannian sphere with metric $g_X$ and $F : X \to \Omega \times \langle 0, t_s \rangle$ its time-dependent embedding in $\Omega$. The image of $F^t = F(\cdot, t)$ will be denoted by $S^t$. The surface $S^0$ will represent the initial estimate of the surface of the segmented object and $S^{t_s}$ will be the result of the segmentation procedure that should be as close to the actual surface of the segmented object as possible. We let $F$ evolve by the 3D analogue of the geodesic active contour model [1],

$$\partial_t F = a \left( \nabla e \cdot N \right) N + b e \Delta_{g_F} F \tag{2}$$

where $N$ is a unit normal to $S$ and $\Delta_{g_F} F$ denotes the Laplace-Beltrami operator with respect to the metric $g_F$ induced on $X$ by $F$. It is known that $\Delta_{g_F} F$ is equal to the mean curvature vector of $F$. As we can see from (1), the curvature term is dominant in regions with low intensity changes where $e$ is close to 1 and its gradient is close to 0. On the contrary, the gradient of $e$ becomes significant near the edges where $e$ decreases and approaches 0 for large values of $K$ and $\|\nabla I(x, y, z)\|$. The parameters $a \in \mathbb{R}_+$, $b \in \mathbb{R}_+$ are added to control the influence of the two terms on the segmentation process.

In order to be able to redistribute the mesh points along the surface during the evolution, we enrich (2) with a tangential velocity term. The new model reads

$$\partial_t F = a \left( \nabla e \cdot N \right) N + b e \Delta_{g_F} F + v_T = v_N + v_T \tag{3}$$

where $v_T$ is a tangential vector field on $S$ and $v_N$ denotes the normal component of the evolution, $v_N = a\,(\nabla e \cdot N)\,N + b e \Delta_{g_F} F$.

In our case, we use an area-oriented tangential redistribution [6] derived from the evolution of the induced metric $g_F$. Both metrics $g_X$ and $g_F$ induce measures on $X$; let us denote them by $\mu_X$ and $\mu_F$. The Radon-Nikodým derivative $G = \frac{\partial \mu_F}{\partial \mu_X}$ is called the *area density* of $F$. It evolves along with $F$ as [3]

$$\partial_t G = \left(v_N \cdot h + \operatorname{div}_{g_F} w_T\right) G \tag{4}$$

where $h$ is the mean curvature vector of $F$, $w_T$ is a vector field on $X$ obtained as the pull-back of $v_T$ along $F$ and $\operatorname{div}_{g_F}$ denotes the divergence with respect to the metric $g_F$. From this follows the evolution of the area of $S$,

$$\partial_t A = \int_X \left(v_N \cdot h + \operatorname{div}_{g_F} w_T\right)\,\mathrm{d}\mu_F = \int_X v_N \cdot h\,\mathrm{d}\mu_F. \tag{5}$$

The embedding $F^t$ is called *area uniform with respect to* $g_X$ if its area density $G^t$ is constant. Our redistribution method is based on the requirement $G^t \longrightarrow_{t\to\infty} C$ that is equivalent to the practically more convenient dimensionless condition

$$\frac{G^t}{A^t} \xrightarrow[t\to\infty]{} C.$$

This can be achieved, for example, if $\frac{G}{A}$ satisfies

$$\partial_t \left(\frac{G}{A}\right) = \omega \left(C - \frac{G}{A}\right) \tag{6}$$

where $\omega \in \mathbb{R}_+ \times \langle 0, t_s \rangle$ represents the redistribution speed. Since we know how both $G$ and $A$ evolve, the combination of (4) and (5) with (6) implies that $w_T$ has to satisfy

$$\operatorname{div}_{g_F} w_T = v_N \cdot h - \frac{1}{A}\int_X v_N \cdot h\,\mathrm{d}\mu_F + \omega \left(C\frac{A}{G} - 1\right). \tag{7}$$

Since this condition does not uniquely determine $w_T$, we suppose, in addition, that $w_T$ is a gradient field, that means $w_T = \nabla_{g_F} \psi$, $\psi : X \times \langle 0, t_s \rangle \to \mathbb{R}$. Thus we obtain

$$\Delta_{g_F} \psi = v_N \cdot h - \frac{1}{A}\int_X v_N \cdot h\,\mathrm{d}\mu_F + \omega \left(C\frac{A}{G} - 1\right) \tag{8}$$

that yields a unique solution if we prescribe the value of $\psi$ in one point of $X$.

**Fig. 1** The surface discretization mesh. *Left* the triangulation of the topological sphere $X$. *Right* the corresponding approximation of the embedded surface $F^n(X)$

## 3 Numerical Approximation of the Segmentation Model

The time discretization of our segmentation model (3) is semi-implicit,

$$\frac{F^n - F^{n-1}}{\tau} = a \left( \nabla e \cdot N^{n-1} \right) N^{n-1} + be\Delta_{g_{F^{n-1}}} F^n + v_T^{n-1}. \tag{9}$$

The space discretization is based on the finite volume approach and it includes two meshes—the mesh discretizing the surface $S^n$ and the voxel grid of the image used to approximate $e$ and $\nabla e$. First, let us consider a triangulation of $X$ which is a simplicial complex homeomorphic to $X$. The corresponding homeomorphism induces a triangular structure on $X$ consisting of vertices $X_i$, $i = 1 \ldots n_v$, edges $e_j$, $j = 1 \ldots n_e$, and triangles $\mathcal{T}_k$, $k = 1 \ldots n_t$.

Now we construct the control volume mesh (Fig. 1). The point $X_i$ is the common vertex of $m$ mesh triangles $\mathcal{T}_1, \ldots, \mathcal{T}_m$ and $m$ edges $e_1, \ldots, e_m$, where $e_p$ connects $X_i$ with its neighbor $X_{i_p}$ (we use local indexing for simplicity). The triangle $\mathcal{T}_p$ admits a barycentric coordinate system—each point of the triangle can be expressed as $P = \lambda_1 X_i + \lambda_2 X_{i_p} + \lambda_3 X_{i_{p+1}}$ where $\lambda_1 + \lambda_2 + \lambda_3 = 1$. Let $B_p$ be the barycenter of $\mathcal{T}_p$ and $C_p$ the center of $e_p$, $p = 1 \ldots m$, and let the barycentric subdivision of $\mathcal{T}_p$ be constructed using these points. The control volume $V_i$ corresponding to $X_i$ is constructed as the union of the triangles $\mathcal{V}_{p,1} = M_i C_p B_p$ and $\mathcal{V}_{p,2} = M_i B_p C_{p+1}$ for $p = 1 \ldots m$ where we set $C_{m+1} = C_1$. Each triangle contains two control volume edges $\sigma_{p,1} = C_p B_p$, $\sigma_{p,2} = B_p C_{p+1}$.

The manifold $X$ can be embedded in $\mathbb{R}^3$ by $\bar{F}^n$, a piecewise linear approximation of $F^n$. First, we set $\bar{F}^n(X_i) = F^n(X_i)$. Then, for any triangle $\mathcal{T}_p$ with vertices $X_i$, $X_{i_p}$, $X_{i_{p+1}}$, we set $\bar{F}^n(\lambda_1 X_i + \lambda_2 X_{i_p} + \lambda_3 X_{i_{p+1}}) = \lambda_1 F^n(X_i) + \lambda_2 F^n(X_{i_p}) + \lambda_3 F^n(X_{i_{p+1}})$. The embedding $\bar{F}^n$ induces a metric $g^n$ on $X$ which induces a measure $\mu^n$ on $X$.

The surface $\bar{S}^n = \bar{F}^n(X)$ is a polyhedron with vertices $\bar{F}^n(X_i) = F^n(X_i) = F_i^n$, edges $\bar{e}_j^n = \bar{F}^n(e_j)$ and triangular faces $\bar{\mathcal{T}}_p^n = \bar{F}^n(\mathcal{T}_p)$. The approximation of the unit normal to $S^n$ at $F_i^n$ is denoted by $N_i^n$. We will use the notation $v_{p,1}^n$, $v_{p,2}^n$ for

the outward unit normals to $\bar{F}^n(\sigma_{p,1})$ and $\bar{F}^n(\sigma_{p,2})$ in the plane of $\bar{\mathcal{T}}_p^n$. Further, $\theta_{p,1}^n$ and $\theta_{p,2}^n$ will represent the angles of $\mathcal{T}_p$ adjacent to $X_{i_p}$ and $X_{i_{p+1}}$, respectively, measured in the metric $g^n$.

Integrating (9) over $V_i$, we obtain

$$
\int_{V_i} \frac{F^n - F^{n-1}}{\tau} \mathrm{d}\mu_{F^{n-1}} = \int_{V_i} a \left( \nabla e \cdot N^{n-1} \right) N^{n-1} \mathrm{d}\mu_{F^{n-1}}
$$
$$
+ \int_{V_i} be \Delta_{g_{F^{n-1}}} F^n \mathrm{d}\mu_{F^{n-1}} + \int_{V_i} v_T^{n-1} \mathrm{d}\mu_{F^{n-1}}.
\tag{10}
$$

The term on the left hand side can be approximated simply by

$$
\int_{V_i} \frac{F^n - F^{n-1}}{\tau} \mathrm{d}\mu_{F^{n-1}} \approx \mu^n(V_i) \frac{F_i^n - F_i^{n-1}}{\tau}.
\tag{11}
$$

In order to approximate $\|\nabla I\|$, $e$ and $\nabla e$, we use the voxel structure of the image $I$. Let us suppose that the voxels are cubes with side length $h$. The voxel with coordinates $x \in \mathbb{N}$, $y \in \mathbb{N}$, $z \in \mathbb{N}$ will be denoted by $P_j$, $j = (x, y, z)$. Since $X$ is embedded in the image domain $\Omega$, the voxel coordinates corresponding to $F_i^n$ are obtained simply by rounding its coordinates. The representative value of $I$ and $e$ in $P_j$ will be denoted by $I_j$ and $e_j$. Further, $v_1$, $v_2$ and $v_3$ are the standard basis vectors in $\mathbb{R}^3$. The 6 voxel faces will be represented by $F_j^{\pm p}$, $p = 1, 2, 3$.

First, let us construct the approximation of $\nabla I$ in the barycenter $c_j^{\pm p}$ of $F_j^{\pm p}$. The derivative in the direction of $v_p$ is discretized by

$$
D^{\pm p} I_j = \pm \left( I_{j \pm v_p} - I_j \right) / h.
$$

For the other two directions $v_q$, $q \neq p$, we will use the values of $I$ in the centers of the voxel edges $F_j^{\pm p, \pm q}$; we denote them by $I_{j \pm \frac{1}{2} v_p \pm \frac{1}{2} v_q}$. Then we use

$$
D^{\pm p, q} I_j = \frac{I_{j \pm \frac{1}{2} v_p + \frac{1}{2} v_q} - I_{j \pm \frac{1}{2} v_p - \frac{1}{2} v_q}}{h}, \quad I_{j \pm \frac{1}{2} v_p \pm \frac{1}{2} v_q} = \frac{I_j + I_{j \pm v_p} + I_{j \pm v_q} + I_{j \pm v_p \pm v_q}}{4}.
$$

Finally, we take

$$
Q_j^{\pm p} = \left( (D^{\pm p} I_j)^2 + \sum_{p \neq q} (D^{\pm p, q} I_j)^2 \right), \quad \|\nabla I(x, y, z)\|^2 \approx \left( \sum_{p=1}^{3} (Q_j^{+p} + Q_j^{-p}) \right) / 6.
\tag{12}
$$

The gradient of $e$ is computed analogously.

Now, the surface normal at $F_i^n$ is approximated by the arithmetic mean of the normals to all triangles containing $F_i^n$. This completes the approximation of the first term on the right hand side of (10). As for the second term, we use

$$\int_{V_i} be\Delta_{g_{F^{n-1}}} F^n \, d\mu_{F^{n-1}} \approx b_i e_i \frac{1}{2} \sum_{p=1}^{m} \left( \cot \theta_{i,p-1,1}^{n-1} + \cot \theta_{i,p,2}^{n-1} \right) (F_i^n - F_{i_p}^n) \quad (13)$$

where we used the cotangent scheme [4] to discretize the Laplace-Beltrami operator. The value $e_i$ is the value of $e$ in the voxel containing $F_i^n$. We consider $\theta_{i,0,1}^{n-1} = \theta_{i,m,1}^{n-1}$.

The last term to discretize is the integral of the tangential velocity. Since $w_T^n$ is a gradient field, the following version of the Stokes theorem applies [2, 6]

$$\int_{V_i} v_T^{n-1} \, d\mu_{F^{n-1}} = \int_{\partial V_i} \psi^{n-1} v_i^{n-1} \, dH_{\mu_{F^{n-1}}} - \int_{V_i} \psi^{n-1} h^{n-1} \, d\mu_{F^{n-1}}.$$

This yields the approximation

$$\int_{V_i} v_T^{n-1} \, d\mu_{F^{n-1}} \approx \sum_{p=1}^{m} \left( \|\sigma_{i,p,1}\|_{n-1} \psi_{i,p,1}^{n-1} v_{i,p,1}^{n-1} + \|\sigma_{i,p,2}\|_{n-1} \psi_{i,p,2}^{n-1} v_{i,p,2}^{n-1} \right) \\ - \mu^{n-1}(V_i) \psi_i^{n-1} h_i^{n-1} \quad (14)$$

where $\|\cdot\|_{n-1}$ denotes the length computed by the metric $g^{n-1}$ and $\psi_{i,p,1}^{n-1}$, $\psi_{i,p,2}^{n-1}$ are the values of $\psi^{n-1}$ in the midpoints of $\sigma_{i,p,1}$ and $\sigma_{i,p,2}$. They are obtained from the values of $\psi^{n-1}$ in the vertices $X_i$ by linear interpolation.

The function $\psi$ is computed from (8) where, again, we use the cotangent scheme to discretize the Laplace-Beltrami operator of $\psi^{n-1}$. This scheme is also used to approximate the mean curvature vector $h$, namely

$$h_i^{n-1} = \frac{1}{\mu^{n-1}(V_i)} \sum_{p=1}^{m} \left( \cot \theta_{i,p-1,1}^{n-1} + \cot \theta_{i,p,2}^{n-1} \right) (F_i^n - F_{i_p}^n). \quad (15)$$

The area of $S^{n-1}$ is approximated by

$$A^{n-1} = \sum_{i=1}^{n_v} \mu^{n-1}(V_i). \quad (16)$$

Alternatively, $A(t^{n-1})$ could be approximated as

$$A(t^{n-1}) = \int_X G(x, t^{n-1}) \, d\mu_X \approx \sum_{i=1}^{n_v} G_i^{n-1} \mu_X(V_i).$$

This leads to an approximation of the volume density $G^{n-1}$. Since we did not particularly specify $\mu_X$, we can assume that $\mu_X(X) = 1/C$ and $\mu_X(V_i) = \mu_X(X)/n_v$ for all $i = 1 \ldots n_v$. Then we can set

**Fig. 2** Cell nucleus segmentation—the data, the initial surface and the segmented surface shown in two different 2D slices

$$G_i^{n-1} = \mu^{n-1}(V_i)\frac{n_v}{\mu_X(X)} = Cn_v\mu^{n-1}(V_i), \qquad C\frac{A}{G} \approx \frac{A^{n-1}}{n_v\mu^{n-1}(V_i)}. \tag{17}$$

## 4 Experiments

Finally, we present two examples of segmentation of biological images. The images display cell nuclei and cell membranes of a zebrafish embryo. Segmentation of cell nuclei and cells has a large number of applications [5]. Particularly, segmentation in form of a triangulated surface can be easily used to compute the area of the surface of a cell or to evaluate the shape of a cell. Before segmenting, the images were pre-filtered by the geodesic mean curvature flow method [1].

In both experiments, we used a relatively large value of $\omega$. Since the tangential direction is approximated, for such large values, the points tend to deviate from the surface where they should be situated [6]. In order to overcome this difficulty, in each time step we first perform the corresponding normal movement, then the tangential movement and afterwards we project the new vertices on the surface obtained by the normal movement alone.

The first experiment deals with the nucleus image. We show segmentation of a single cell nucleus. We performed 400 time steps and the model parameters were set to $n_v = 258$, $\tau = 0.001$, $h = 1.0$, $\omega = 100.0$, $a = 1.0$, $b = 200.0$ for time steps $1 \ldots 200$ and $b = 1.0$ after. The initial condition was a sphere centered in a manually estimated nucleus center. Figure 2 shows two different 2D slices of the data, the initial surface and the segmentation result. Figure 3 shows the effect of the tangential redistribution of mesh points during the computation. We can see that the tangential movement leads to more evenly distributed mesh points and thus a more correct representation of the surface. Quantitatively evaluated, the ratio of the minimal and maximal control volume area was 0.176 when no redistribution was applied while it reached 0.894 when the redistribution step was included.

In the second experiment, we segmented several cells from the membrane image. Membrane data are usually of a worse quality and more difficult to segment than nucleus data. We performed 400 time steps and we used $n_v = 258$, $\tau = 0.003$, $h = 1.0$, $\omega = 100.0$, $a = 3.0$, $b = 20.0$ for time steps $1 \ldots 200$ and $b = 1.0$ after.

**Fig. 3** Cell nucleus segmentation. *Left* the segmented nucleus surface obtained with no tangential redistribution. *Right* the surface obtained with tangential redistribution of mesh points, $\omega = 100.0$



**Fig. 4** Cell membrane image segmentation—2D slices of the data, of the initial condition and of the segmented surfaces

**Fig. 5** Cell membrane image segmentation—the segmented surfaces



Similarly to the case of nucleus segmentation, the initial surface was a sphere (of the same radius for all cells). Figure 4 shows a 2D slice of the image, the initial surfaces and the segmented cells. Figure 5 shows the whole segmented cells.

# References

1. Caselles, V., Kimmel, R., Sapiro, G.: Geodesic active contours. Int. J. Comput. Vision **22**, 61–79 (1997)
2. Dziuk, G., Elliott, C.: Finite elements on evolving surfaces. IMA J. Numer. Anal. **27**, 262–292 (2007)
3. Mantegazza, C.: Lecture Notes on Mean Curvature Flow. Springer, Berlin (2011)
4. Meyer, M., Desbrun, M., Schroeder, P., Barr, A.: Discrete differential geometry operators for triangulated 2manifolds. Vis. Math. **III**, 35–57 (2003)
5. Mikula, K., Peyrieras, N., Remesikova, M., Stasova, O.: Segmentation of 3d cell membrane images by pde methods and its applications. Comput. Biol. Med. **41**(6), 326–339 (2011)
6. Mikula, K., Remesikova, M., Sarkoci, P., Sevcovic, D.: Surface evolution with tangential redistribution of points (2013) (submitted)

# A Model Reduction Framework for Efficient Simulation of Li-Ion Batteries

**Mario Ohlberger, Stephan Rave, Sebastian Schmidt and Shiquan Zhang**

**Abstract** In order to achieve a better understanding of degradation processes in lithium-ion batteries, the modelling of cell dynamics at the mircometer scale is an important focus of current mathematical research. These models lead to large-dimensional, highly nonlinear finite volume discretizations which, due to their complexity, cannot be solved at cell scale on current hardware. Model order reduction strategies are therefore necessary to reduce the computational complexity while retaining the features of the model. The application of such strategies to specialized high performance solvers asks for new software designs allowing flexible control of the solvers by the reduction algorithms. In this contribution we discuss the reduction of microscale battery models with the reduced basis method and report on our new software approach on integrating the model order reduction software pyMOR with third-party solvers. Finally, we present numerical results for the reduction of a 3D microscale battery model with porous electrode geometry.

M. Ohlberger · S. Rave (✉)
Center for Nonlinear Science and Applied Mathematics Muenster, Einsteinstrasse 62,
48149 Muenster, Germany
e-mail: stephan.rave@uni-muenster.de

M. Ohlberger
e-mail: mario.ohlberger@uni-muenster.de

S. Schmidt
Fraunhofer Institute for Industrial Mathematics ITWM, Fraunhofer-Platz 1,
67663 Kaiserslautern, Germany
e-mail: sebastian.schmidt@itwm.fraunhofer.de

S. Zhang
School of Mathematics, Sichuan University, 610064 Chengdu, China
e-mail: shiquanz3@gmail.com

# 1 Introduction

A major cause for the failure of rechargeable lithium-ion batteries is the deposition of metallic lithium at the negative battery electrode (Li-plating). Once established, this metallic phase can grow in the form of dendrites to the positive electrode, ultimately short-circuiting the cell. As Li-plating is initiated at the interface between active electrode particles and the electrolyte, understanding of this phenomenon is only gained through physical models accounting for effects on the micrometer-scale. This in turn requires highly resolved meshes in the model discretization.

A thermodynamically consistent microscale battery model was developed in [7]. Based on a finite volume discretization [9], this model has been implemented at Fraunhofer ITWM in the battery simulation software BEST [8]. However, since such microscale discretizations lead to very large, highly nonlinear equation systems, simulations can currently only be performed on small portions of the cell and parameter studies testing different charging regimes or operating conditions are very time consuming. It is therefore desirable to combine microscale modeling with model order reduction strategies which are able to reduce the computation time while at the same time keeping the microscopic features of the model.

The reduced basis method is a well-established approach for model order reduction of problems given by parametric partial differential equations and has been successfully adapted to various industrial applications (see references in [5]). In this approach, the original equation is projected onto a low-dimensional discrete function space which has been constructed from the solution trajectories of the high-dimensional problem for selected parameters of a well-chosen training set. The applicability of the method to nonlinear finite volume discretizations has been been shown in [4, 5]. Results for the model order reduction of a pseudo-2D battery model using similar techniques have been presented in [6].

A major challenge for the implementation of reduced basis schemes lies, however, in their integration with (already existing) PDE solvers: in those schemes the solver has to be controlled by the reduction algorithm which, apart from solving the high-dimensional problem, now also has to provide the reduction data needed to perform the low-dimensional simulations. Moreover, the solver is usually unable to perform the reduced computations, which are based on different data structures. This often leads to insertion of model reduction specific algorithms into the solver's code base, while in a separate code base the solution algorithm for the reduced problem is re-implemented [3]. As a result, code is duplicated and the adoption of a different model reduction strategy requires changes in both code bases.

After discussing the application of the reduced basis method to the microscale model from [7], we present the design of our new model reduction software pyMOR [1] which is specifically tailored to address these problems by offering a deep and flexible integration with external PDE solvers. We will conclude with first numerical results for the reduction of the full 3D-model with porous electrode geometries, underlining the potential of the model reduction approach.

## 2 Reduction of the Microscale Model

Our work is based on the microscale battery model introduced in [7]. Under the assumption of a globally constant temperature $T$, this model is given by a system of partial differential equations for the concentration of Li$^+$-ions $c$ and the electrical potential $\phi$ on each part of the domain, i.e. the positive and negative electrodes, the electrolyte and the current collectors. Each of these systems is of the form

$$\frac{\partial c}{\partial t} + \nabla \cdot N = 0, \qquad \nabla \cdot j = 0,$$

where $N = -(\alpha(c, \phi)\nabla c + \beta(c, \phi)\nabla\phi)$, $j = -(\gamma(c, \phi)\nabla c + \delta(c, \phi)\nabla\phi)$ with the coefficients $\alpha, \beta, \gamma, \delta$ depending on the domain for which the system is given. While these coefficients can be considered constant in first approximation, a strong nonlinearity enters the model through the interface conditions between electrolyte and active particles in the electrodes. These conditions are given by prescribing the normal interface fluxes of concentration and potential into the electrolyte via the Butler-Volmer kinetics, i.e.

$$j_s \cdot n = j_e \cdot n = 2k\sqrt{c_e c_s(c_{max} - c_s)}\sinh\left(\frac{\phi_s - \phi_e - U_0(\frac{c_s}{c_{max}})}{2RT} \cdot F\right),$$

and $N_s \cdot n = N_e \cdot n = j_s \cdot n/F$. Here the subscripts $s$ ($e$) denote the value of the respective quantity in the active particle (electrolyte) domain at the interface, and $n$ is the unit normal at the interface pointing into the electrolyte. $U_0$ denotes the open circuit potential, $k$ is a reaction rate, $c_{max}$ the maximum Li-ion concentration in the particle and $T$ the temperature. The constants $F$ and $R$ denote the Faraday and universal gas constants. The system is closed via appropriate boundary conditions as well as interface conditions for the current collectors. E.g. a constant charge rate $I$ corresponds to the Neumann boundary condition $j \cdot n = -I$ at the positive electrode side of the domain.

### 2.1 Discretization

A discretization of the model based on a cell centered finite volume scheme has been introduced in [9]. In this discretization, the interface conditions between electrolyte and active particles are incorporated into the numerical fluxes and the implicit Euler method is used for time discretization. As a result, one obtains nonlinear equation systems of the form

$$\left[\frac{1}{\Delta t}(c_\mu^{(t+1)} - c_\mu^{(t)})\right] + A_\mu\left(\begin{bmatrix} c_\mu^{(t+1)} \\ \phi_\mu^{(t+1)} \end{bmatrix}\right) = 0, \qquad c_\mu^{(t)}, \phi_\mu^{(t)} \in V_h \qquad (1)$$

with $A_\mu$ denoting the finite volume space operator acting on the discrete function space $V_h \oplus V_h$. The subscript indicates the dependence of the solution on a certain set of parameters $\mu$ (we consider the charge rate and temperature in our example below). The discrete equation systems are solved in BEST with a Newton scheme utilizing an algebraic multigrid solver for the linear systems in each Newton step.

## 2.2 Reduced Basis Approximation

The reduced basis method is based on the idea of performing a Galerkin projection of the high-dimensional discrete equations (1) onto low-dimensional subspaces $\tilde{V}_c, \tilde{V}_\phi \subset V_h$ constructed from solutions of (1) for appropriately selected parameters. Under this projection, (1) is transformed into

$$\left[ \begin{matrix} \frac{1}{\Delta t}(\tilde{c}_\mu^{(t+1)} - \tilde{c}_\mu^{(t)}) \\ 0 \end{matrix} \right] + \left\{ P_{\tilde{V}} \circ A_\mu \right\} \left( \left[ \begin{matrix} \tilde{c}_\mu^{(t+1)} \\ \tilde{\phi}_\mu^{(t+1)} \end{matrix} \right] \right) = 0, \qquad \tilde{c}_\mu^{(t)} \in \tilde{V}_c, \ \tilde{\phi}_\mu^{(t)} \in \tilde{V}_\phi, \ (2)$$

where $P_{\tilde{V}}$ denotes the orthogonal projection onto the reduced space $\tilde{V} := \tilde{V}_c \oplus \tilde{V}_\phi$. After this projection has been performed in a preceding "offline-phase", the resulting low-dimensional system can be solved quickly for new parameter values in a following "online-phase".

For the selection of $\tilde{V}_c$ and $\tilde{V}_\phi$ a large variety of algorithms has been considered ([5] and references therein), many of which are based on a greedy search over a prescribed (or adaptively refined) training set of parameters: in each round of the algorithm, an error estimator is used to search the training set for the parameter $\mu^*$ to which the solution of (1) is worst approximated by the solution of the reduced problem (2). The high-dimensional solution trajectory $[c_{\mu^*}^{(t)}, \phi_{\mu^*}^{(t)}]$ is then computed and $\tilde{V}_c$, $\tilde{V}_\phi$ are enlarged by vectors from the linear span of this trajectory via an appropriate extension algorithm. As the reduced spaces are constructed from solutions of the full microscale model, characteristic features, e.g. concentration hotspots in certain electrode regions due to local particle geometry, are still representable within these spaces, despite their low dimensionality.

While posed on low-dimensional spaces, problem (2) still depends on evaluations of the high-dimensional operator $A_\mu$. This dependency can be removed by application of the so-called empirical operator interpolation method [4]. In this approach, the given operator is only evaluated at a small number of degrees of freedom (DOFs) of the discrete space. The evaluation of the full operator is then approximated via linear combination with a pre-computed (collateral) interpolation basis. The interpolated operator can be evaluated quickly, independently of the dimension of $V_h$, due to the locality of finite volume operators: the evaluation of $A_\mu$ at $M$ degrees of freedom only requires the knowledge of its argument at $M' \leq C \cdot M$ DOFs with $C$ being determined by the maximum number of cell neighbours in the given grid. If we

denote by $\tilde{A}_\mu : \mathbb{R}^{M'} \to \mathbb{R}^M$ the restricted operator and by $R_{M'} : V_h^2 \to \mathbb{R}^{M'}$, $I_M : \mathbb{R}^M \to V_h^2$ the operators given by projection onto the interpolation DOFs and linear combination with the collateral basis, we obtain the fully reduced equation systems

$$\left[ \begin{array}{c} \frac{1}{\Delta t}(\tilde{c}_\mu^{(t+1)} - \tilde{c}_\mu^{(t)}) \\ 0 \end{array} \right] + \left\{ (P_{\tilde{V}} \circ I_M) \circ \tilde{A}_\mu \circ R_{M'} \right\} \left( \left[ \begin{array}{c} \tilde{c}_\mu^{(t+1)} \\ \tilde{\phi}_\mu^{(t+1)} \end{array} \right] \right) = 0. \quad (3)$$

The linear operators $P_{\tilde{V}} \circ I_M$ and $R_{M'}$ can be pre-evaluated during the offline-phase for a given basis of $\tilde{V}$, completely eliminating high-dimensional operations from (3). For the determination of the interpolation DOFs and collateral basis, greedy search strategies can again be utilized [4].

## 3 A New Software Framework

The implementation of reduced basis schemes involves several building blocks: solution of the detailed problem (1) for a given parameter, projection of the operators, extension of the reduced spaces (high-dimensional operations), as well as solution of the reduced problem (3), estimation of the reduction error and greedy algorithms (low-dimensional operations). In previous software approaches [3], the implementation of all high-dimensional operations takes place in the solver code, whereas the low-dimensional operations are implemented in a separate model reduction software. As a consequence, both code bases have to be adapted if the reduction strategy shall be modified. This can slow down implementation of new algorithms significantly if the solver is developed by a different team than the model reduction software. Moreover, despite the fact that (1) and (3) are of the same mathematical structure, both software packages need to implement the same algorithm for solving the respective problems. In particular, for empirical operator interpolation the restricted operator $\tilde{A}_\mu$ has to be implemented again for the reduced scheme.

The design of pyMOR mitigates these difficulties by exploiting the observation that all aforementioned building blocks can be implemented in terms of operations on the following types of objects, either provided by implementations in pyMOR itself (usually low-dimensional objects) or by external solvers (usually high-dimensional objects):

- **Vector arrays** store collections of vectors, supporting basic linear algebra operations, e.g. computation of linear combinations of vectors or scalar products. Selected DOFs can be extracted for the implementation of operator interpolation.
- **Operators** represent linear or nonlinear operators, bilinear forms or functionals. Operators can be applied to vector arrays. Linear solvers are exposed through application of the inverse operator, Jacobians and restricted operators can be formed.

**Fig. 1** Detailed simulation of battery model with DUNE on a $48 \times 24 \times 24\ \mu m^3$ computational domain with random electrode geometry. Coloring indicates $Li^+$ concentration in active particles (electrolyte not displayed)



**Fig. 2** Sketch of the interface concept for the integration of pyMOR with external solvers



- **Discretizations** encode as containers for operators the mathematical structure of a given discrete problem and implement algorithms for solving the problem in terms of the operators they contain.

All algorithms in pyMOR are implemented in terms of the interfaces provided by these classes. As an important consequence, there is no distinction between high- and low-dimensional objects in pyMOR except for the different types of vector arrays or operators that represent them. In particular, the same discretization class can be used to solve (1) as well as (3) or (2). The reduction process merely consists in the replacement of operators of a given discretization object by the corresponding projected operators. For empirical interpolation, pyMOR implements a generic interpolated operator which can be used to efficiently interpolate any restrictable operator in pyMOR. The evaluation of the restricted operator $\tilde{A}_\mu$ can still be performed by the same code used to evaluate the full operator $A_\mu$ (Fig. 1).

As a consequence of this design, the model reduction algorithms in pyMOR are completely decoupled from the development of the high-dimensional discretizations (cf. Fig. 2).

**Table 1** Constants used in numerical example

| Domain | $\alpha$ | $\beta$ | $\gamma$ | $\delta$ | $c_0$ | $c_{max}$ | k |
|---|---|---|---|---|---|---|---|
| Electrolyte | $1.622 \cdot 10^{-6}$ | 0 | $-5.171 \cdot 10^{-5} \cdot T$ | 0.02 | $1.200 \cdot 10^{-3}$ | – | – |
| Pos. electrode | $1.0 \cdot 10^{-10}$ | 0 | 0 | 0.38 | $2.057 \cdot 10^{-2}$ | $2.367 \cdot 10^{-2}$ | 0.2 |
| — current coll. | 0 | 0 | 0 | 0.38 | 0 | – | – |
| Neg. electrode | $1.0 \cdot 10^{-10}$ | 0 | 0 | 10 | $2.639 \cdot 10^{-3}$ | $2.468 \cdot 10^{-2}$ | 0.002 |
| — current coll. | 0 | 0 | 0 | 10 | 0 | – | – |

$c_0$ denotes initial concentration, furthermore, $U_0(x) = -0.132 + 1.41 \cdot \exp(-3.52x)$ for the negative and $U_0(x) = 4 + 0.07 \cdot \tanh(-22x + 12) - 0.1 \cdot (1/(1.002 - x)^{0.37} - 1.6) - 0.045 \cdot \exp(-72x^8) + 0.01 \cdot \exp(-200(x - 0.19))$ for the positive electrode, $R = 8.314$, $F = 9.6487 \cdot 10^4$

## 3.1 Implementational Aspects

Following the line of most other model order reduction packages, we chose with Python a scripting language for the implementation of pyMOR. Such languages offer a high amount of interactivity, making it very easy to experiment with various variants of model reduction algorithms.

While there is no underlying assumption of how the communication through the abstract interfaces is handled, we favour, where possible, a tight integration of external solvers with pyMOR. In particular for shared-memory solvers, an attractive option is the compilation of the solver code as a shared library which then can be directly loaded as a Python extension module. Apart from offering the easiest and at the same time most efficient way of integration, an additional benefit is the direct accessibility of solver data structures from Python which can be exploited to quickly augment the high-dimensional code with additional features. This route of development has also been chosen for the ongoing integration of pyMOR with BEST within the publicly founded MULTIBAT project.

## 4 Numerical Results

In order to provide a testbed for our reduction framework, an experimental implementation of the battery model has been developed based on the PDELab discretization module for the DUNE software framework [2] (cf. Fig. 1). As a first experiment, we considered a small 3D test problem with randomly generated electrode geometry, for which we evaluated the approximation quality of the reduced basis projection (2). We chose constant material properties resulting in the coefficients in Table 1. The computational domain was of size $4.8 \cdot 10^{-3} \times 2.4 \cdot 10^{-2} \times 2.4 \cdot 10^{-2}$ (cm$^3$) which was meshed with a regular $40 \times 20 \times 20$ grid. The width of the electrodes (current collectors) was 10 (5) grid cells. The positive (negative) electrode was filled to 61.4 % (74.2 %) with particle cells. 20 time steps of length 30 (s) were made.

**Table 2** Relative $L^\infty - L^2$ errors for the reduced basis approximation (2) of the high-dimensional model (1)

| dim $\tilde{V}_c = \tilde{V}_{phi}$ | 8 | 16 | 24 | 32 |
|---|---|---|---|---|
| Concentration | $8.7 \cdot 10^{-3}$ | $1.9 \cdot 10^{-3}$ | $1.2 \cdot 10^{-3}$ | $4.3 \cdot 10^{-4}$ |
| Potential | $1.3 \cdot 10^{-3}$ | $2.1 \cdot 10^{-4}$ | $7.7 \cdot 10^{-5}$ | $1.5 \cdot 10^{-5}$ |

The parameters, charge rate $I$ and temperature $T$, were allowed to vary in the intervals $[10^{-4}, 10^{-3}]$ (A/cm$^2$) and $[250, 350]$ ($K$). The reduced spaces were constructed with the POD-Greedy algorithm [5] on a training set of $3 \times 3$ equidistant parameters, using the true reduction error for snapshot selection. During each extension step, both reduced spaces were extended separately by orthogonally projecting the selected trajectory onto the respective reduced space and then enlarging the space with the first POD mode of the trajectory of projection errors. In Table 2, the maximum reduction error over the whole parameter space is estimated for different basis sizes by computation of the errors for 20 randomly selected new parameters.

# References

1. pyMOR—Model Order Reduction with Python. http://www.pymor.org
2. Bastian, P., Blatt, M., Dedner, A., Engwer, C., Klöfkorn, R., Ohlberger, M., Sander, O.: A generic grid interface for parallel and adaptive scientific computing. Part I: abstract framework. Computing **82**(2–3), 103–119 (2008)
3. Drohmann, M., Haasdonk, B., Kaulmann, S., Ohlberger, M.: A software framework for reduced basis methods using Dune-RB and RBmatlab. In: Dedner, A., Flemisch, B., Klöfkorn, R. (eds.) Advances in DUNE, pp. 77–88. Springer, Berlin Heidelberg (2012)
4. Drohmann, M., Haasdonk, B., Ohlberger, M.: Reduced basis approximation for nonlinear parametrized evolution equations based on empirical operator interpolation. SIAM J. Sci. Comput. **34**(2), A937–A969 (2012)
5. Haasdonk, B., Ohlberger, M.: Reduced basis method for finite volume approximations of parametrized linear evolution equations. m2an. Math. Model. Numer. Anal. **42**(2), 277–302 (2008)
6. Iliev, O., Latz, A., Zausch, J., Zhang, S.: On some model reduction approaches for simulations of processes in Li-ion battery. In: Proceedings of Algoritmy 2012, Conference on Scientific Computing, Vysoké Tatry, Podbanské, Slovakia, pp. 161–171. Slovak University of Technology in Bratislava (2012)
7. Latz, A., Zausch, J.: Thermodynamic consistent transport theory of li-ion batteries. J. Power Sources **196**(6), 3296–3302 (2011)
8. Less, G.B., Seo, J.H., Han, S., Sastry, A.M., zausch, J., latz, A., schmidt, S., wieser, C., kehrwald, D., fell, S.: Micro-scale modeling of li-ion batteries: parameterization and validation. J. Electrochem. Soc. **159**(6), A697 (2012)
9. Popov, P., Vutov, Y., Margenov, S., Iliev, O.: Finite volume discretization of equations describing nonlinear diffusion in li-ion batteries. In: Dimov, I., Dimova, S., Kolkovska, N. (eds.) Numerical Methods and Applications. Lecture Notes in Computer Science, vol. 6046, pp. 338–346. Springer, Berlin Heidelberg (2011)

# Coupling Free Flow and Porous Medium Flow Systems Using Sharp Interface and Transition Region Concepts

Iryna Rybak

**Abstract** Two different coupling approaches for isothermal single-phase free flow and isothermal single-fluid-phase porous medium systems are considered: sharp interface and transition region approach. The sharp interface concept implies the Beavers–Joseph–Saffman velocity jump condition together with restrictions that arise due to mass conservation and balance of normal forces across the fluid-porous interface. The transition region model is derived by means of the thermodynamically constrained averaging theory (TCAT). The equations are averaged over the thickness of the transition zone in the direction normal to the free flow and porous medium domains being joined. Coupling conditions are the mass conservation, the momentum balance and a generalization of the Beavers–Joseph condition. Two model formulations are compared and numerical simulation results are presented. For discretization of the coupled problem the finite volume method on staggered grids is used.

## 1 Introduction

Coupled free flow and porous medium flow systems arise routinely in environmental settings and industrial applications such as overland flow interactions with groundwater aquifers, evaporation from soil influenced by wind, fluid flow through filters, and water-gas management in fuel cells. Two different models are typically applied to describe physical processes in the flow domains, and these models should be coupled at the interface in the proper way. In the free flow region, the (Navier)–Stokes equations are usually considered to describe momentum conservation while Darcy's law is used to approximate the momentum balance in the porous medium.

I. Rybak (✉)

Institute of Applied Analysis and Numerical Simulation, University of Stuttgart,
Pfaffenwaldring 57, 70569  Stuttgart, Germany
e-mail: rybak@ians.uni-stuttgart.de

**Fig. 1** Coupled single-fluid-phase porous medium and free flow systems at the microscale (*left*) and the macroscale (*middle*, *right*)

Transitions between the free flow and porous medium domains can be modeled by the sharp interface approach imposing the appropriate coupling conditions at the fluid-porous interface [3, 7, 9] or by considering a transition zone between these flow regions and developing a transition region model [6]. The Beavers–Joseph–Saffman condition [1, 10] is a common practice to couple the free flow and porous medium domains, in conjunction with restrictions resulting from the conservation of mass and balance of normal forces across the sharp interface. This approach is restricted to flat interfaces and flows mainly parallel to the porous medium, unlike the transition region approach. When a transition zone is considered between the flow domains, the conservation equations are averaged over the transition region thickness and the resulting model varies in two spatial dimensions [6]. This approach resolves transfer of thermodynamic properties in tangential directions, unlike sharp interface approximations, and is a generalization of the sharp interface concept.

The goal of this work is to compare, both theoretically and numerically, the sharp interface and the transition region coupling concepts for isothermal single-fluid-phase porous medium and free flow systems. The transition region model is derived via TCAT approach in a similar way as for a two-fluid-phase porous medium [6].

## 2 Flow System Description

The system of interest contains a free flow region $\Omega_{\text{ff}}$ composed of a single fluid phase and a porous medium $\Omega_{\text{pm}}$ fully saturated with the same fluid (Fig. 1, left). These flow regions can be separated at the macroscale by a sharp interface $\Gamma$ (Fig. 1, middle) or by a transition region $\Omega_{\text{tr}}$ of a positive thickness $b > 0$ (Fig. 1, right).

In this work, we do not model compositional effects, therefore consider each phase consisting of a single chemical species. In addition, the coupled system is assumed to be isothermal, the fluid phase is incompressible and the solid is rigid.

The mass conservation equation in the free flow domain reads

$$\nabla \cdot \mathbf{v} = \mathbf{0} \quad \text{in } \Omega_{\text{ff}}, \tag{1}$$

where $\mathbf{v}$ is the fluid velocity.

Neglecting convective acceleration, considering the gravity to be the only external force and using Newton's law, the conservation of momentum for steady-state flows can be written as the Stokes equation

$$\nabla \cdot (p\mathbf{I} - 2\mu\mathbf{D}(\mathbf{v})) - \rho\mathbf{g} = \mathbf{0} \quad \text{in } \Omega_{\text{ff}}, \tag{2}$$

where $p$ is the pressure, $\mathbf{I}$ is the identity tensor, $\mu$ is the dynamic viscosity, $\mathbf{D}(\mathbf{v}) = \frac{1}{2}\left(\nabla\mathbf{v} + (\nabla\mathbf{v})^{\text{T}}\right)$ is the rate of strain tensor, $\rho$ is the fluid density, and $\mathbf{g}$ is the gravity.

Fluid flows through the porous medium are usually described by Darcy's law $\mathbf{v} = -\frac{\mathbf{K}}{\mu}(\nabla p - \rho\mathbf{g})$, which is together with the conservation of mass equation yields the porous medium flow formulation

$$-\nabla \cdot \left(\frac{\mathbf{K}}{\mu}(\nabla p - \rho\mathbf{g})\right) = 0 \quad \text{in } \Omega_{\text{pm}}, \tag{3}$$

where $\mathbf{K}$ is the intrinsic permeability tensor.

The free flow model (1), (2) and the porous medium model (3) can be coupled directly at the sharp interface $\Gamma$ or through a transition region, considering a model in $\Omega_{\text{tr}}$ and coupling it with the two flow models at the transition region boundaries. Problem (1)–(3) is also subject to boundary conditions at the external boundary of the coupled free flow and porous medium domains.

# 3 Sharp Interface Concept

We consider a sharp flat interface between the flow domains (Fig. 1, middle) that has no thickness and cannot store and transfer mass and momentum. In this case, the coupling conditions are well established [3, 7, 9], and they are algebraical jump conditions. The *mass conservation* across the interface reads

$$[\mathbf{v}\cdot\mathbf{n}]^{\text{ff}} = -[\mathbf{v}\cdot\mathbf{n}]^{\text{pm}} \quad \text{on } \Gamma, \tag{4}$$

where $\mathbf{n}$ is the unit normal vector at the interface (Fig. 1, middle), $\mathbf{n}^{\text{ff}} = \mathbf{n}, \mathbf{n}^{\text{pm}} = -\mathbf{n}$. The *balance of normal forces* is given by

$$\left[\mathbf{n}\cdot(p\mathbf{I} - 2\mu\mathbf{D}(\mathbf{v}))\cdot\mathbf{n}\right]^{\text{ff}} = [p]^{\text{pm}} \quad \text{on } \Gamma. \tag{5}$$

The *Beavers–Joseph–Saffman* interface condition [1, 10] can be written as

$$\left[\mathbf{v}\cdot\boldsymbol{\tau}_i + \frac{2\sqrt{\mathbf{K}}}{\alpha_{\text{BJ}}}\mathbf{n}\cdot\mathbf{D}(\mathbf{v})\cdot\boldsymbol{\tau}_i\right]^{\text{ff}} = 0, \quad i = 1, \ldots, d-1 \quad \text{on } \Gamma, \tag{6}$$

where $\boldsymbol{\tau}_i$ are the unit tangent vectors to the interface, $\alpha_{\mathrm{BJ}} > 0$ is the Beavers–Joseph parameter, and $d$ is the number of space dimensions.

## 4 Transition Region Concept

There are several possibilities to develop a transition region model for single-fluid-phase systems, e.g. considering the Brinkman equation which is a superposition of the Stokes equation and Darcy's law [2, 4, 8]. In this case, the equations in the transition zone are full dimensional. However, this model cannot be extended to more than one fluid phase and the definition of the flow parameters is not trivial.

To formulate the equations that describe coupled flow between the free flow domain and the porous medium we apply the TCAT approach [5, 6]. This technique is not restricted to the number of fluid phases and flow direction, and allows to derive general models. The transition region is averaged in the direction normal to the boundaries of the flow domains being joined (Fig. 1, right) that leads to the reduction of spatial dimensionality, and the macroscale equations are restricted to the two-dimensional surface. The detailed derivation of a general transition region model is presented in [6]. The objective of this work is to couple the transition region model with the free flow and the porous medium domains for single-fluid-phase systems.

Assuming the velocity of the transition region is zero, the mass conservation equation for the fluid phase can be written as

$$[\varepsilon\rho\mathbf{v}]^{\mathrm{top}} \cdot\mathbf{N} + \nabla^{\backprime}\cdot(\mathrm{b}\varepsilon\rho\mathbf{v}) = [\varepsilon\rho\mathbf{v}]^{\mathrm{bot}} \cdot\mathbf{N} \quad \text{in} \quad \Omega_{\mathrm{tr}}, \tag{7}$$

where the superscripts top and bot determine physical quantities averaged over the top and bottom boundaries $\Gamma_{\mathrm{top}}$ and $\Gamma_{\mathrm{bot}}$ of the transition region, $b > 0$ is the transition region thickness, $\varepsilon$ is the porosity, $\mathbf{N}$ is the unit vector tangent to the axis corresponding to the megascopic dimension (Fig. 1, right), and $\nabla^{\backprime}$ is the macroscale surfical del operator, $\nabla^{\backprime} = \nabla - \mathbf{N}\mathbf{N}\cdot\nabla$.

Equation (7) needs to be closed by specifying the values at the transition region boundaries $\Gamma_{\mathrm{top}}$ and $\Gamma_{\mathrm{bot}}$. These values come from the free flow and porous medium domains accordingly, and serve as the source terms. Considering sharp interfaces between the top of the transition region and the free flow domain, and between the bottom of the transition region and the porous medium, we get

$$[\rho\mathbf{v}]^{\mathrm{ff}} \cdot\mathbf{N} = [\varepsilon\rho\mathbf{v}]^{\mathrm{top}} \cdot\mathbf{N} \quad \text{on} \quad \Gamma_{\mathrm{top}}, \quad \text{and} \quad [\varepsilon\rho\mathbf{v}]^{\mathrm{bot}} \cdot\mathbf{N} = [\varepsilon\rho\mathbf{v}]^{\mathrm{pm}} \cdot\mathbf{N} \quad \text{on} \quad \Gamma_{\mathrm{bot}}.$$

When the transition region thickness $b = 0$, Eq. (7) reduces to the classical condition of mass conservation across the sharp interface

$$[\rho\mathbf{v}]^{\mathrm{ff}} \cdot\mathbf{N} = [\rho\mathbf{u}]^{\mathrm{pm}} \cdot\mathbf{N},$$

where Darcy's velocity $\mathbf{u}$ is the product of the averaged velocity $\mathbf{v}$ and porosity $\varepsilon$. In case of the same fluid in both domains, we get the interface condition (4). Therefore, Eq. (7) can be considered as a generalization of the classical jump condition.

We do not need to consider the solid phase mass conservation equation since the solid is assumed to be non-deformable and rigid.

Under the assumption of slow flow through the transition region, the momentum conservation for the fluid phase can be written as

$$\varepsilon \nabla^\backslash p - \hat{\mathbf{r}}_w^{\text{top}} \cdot \left(\mathbf{v}^{\overline{\text{top}}} - \mathbf{v}\right) - \hat{\mathbf{r}}_w^{\text{bot}} \cdot \left(\mathbf{v}^{\overline{\text{bot}}} - \mathbf{v}\right) = -\varepsilon^2 \hat{\mathbf{R}}_w \cdot \mathbf{v} \quad \text{in} \ \ \Omega_{\text{tr}}, \qquad (8)$$

where $\mathbf{v}$ is the fluid velocity averaged over the transition region thickness, $\mathbf{v}^{\overline{\text{top}}}$ and $\mathbf{v}^{\overline{\text{bot}}}$ are the velocities averaged over the top and bottom boundaries of the transition region, $\hat{\mathbf{r}}_w$ and $\hat{\mathbf{R}}_w$ are the resistance tensors, which depend on the morphology of the transition region. We do not model the solid phase momentum balance equation because the solid phase is rigid.

Again, considering sharp interfaces between the top of the transition region and the free flow domain as well as between the bottom of the transition region and the porous medium, we close Eq. (8). Momentum conservation at the boundary between the free flow domain and the transition region can be written as

$$\left[p\mathbf{I} - 2\mu\mathbf{D}\right]^{\text{ff}} \cdot \mathbf{N} = p^{\text{top}}\mathbf{N} - b\hat{\mathbf{r}}_w^{\text{top}} \cdot \left(\mathbf{v}^{\overline{\text{top}}} - \mathbf{v}\right) \qquad \text{on} \ \ \Gamma_{\text{top}}, \qquad (9)$$

and at the boundary between the transition region and the porous medium domain

$$p^{\text{bot}}\mathbf{N} + b\hat{\mathbf{r}}_w^{\text{bot}} \cdot \left(\mathbf{v}^{\overline{\text{bot}}} - \mathbf{v}\right) = [p]^{\text{pm}}\,\mathbf{N} \qquad \text{on} \ \ \Gamma_{\text{bot}}. \qquad (10)$$

We decompose the momentum conservation Eqs. (9) and (10) into the megascale and tangential components. The normal component of Eq. (9) reads

$$\left[p - 2\mu\mathbf{N}\cdot\mathbf{D}\cdot\mathbf{N}\right]^{\text{ff}} = p^{\text{top}} - b\hat{\mathbf{r}}_w^{\text{top}} \cdot \left(\mathbf{v}^{\overline{\text{top}}} - \mathbf{v}\right) \cdot \mathbf{N} \qquad \text{on} \ \ \Gamma_{\text{top}},$$

and the normal component of Eq. (10) is given by

$$p^{\text{bot}} + b\hat{\mathbf{r}}_w^{\text{bot}} \cdot \left(\mathbf{v}^{\overline{\text{bot}}} - \mathbf{v}\right) \cdot \mathbf{N} = [p]^{\text{pm}} \qquad \text{on} \ \ \Gamma_{\text{bot}},$$

that is combined with the transition region momentum conservation (8) yields

$$\left[p - 2\mu\mathbf{N}\cdot\mathbf{D}\cdot\mathbf{N}\right]^{\text{ff}} = [p]^{\text{pm}} + \left(p^{\text{top}} - p^{\text{bot}}\right) - b\varepsilon^2 \hat{\mathbf{R}}_w \cdot \mathbf{v}\cdot\mathbf{N}. \qquad (11)$$

When $b = 0$, Eq. (11) is the balance of normal forces across the sharp interface

$$\left[p - 2\mu\mathbf{N}\cdot\mathbf{D}\cdot\mathbf{N}\right]^{\mathrm{ff}} = [p]^{\mathrm{pm}}.$$

If $b > 0$, we need to define the pressure and the normal component of velocity at the top and bottom boundaries of the transition region

$$p^{\mathrm{top}} = [p]^{\mathrm{ff}}, \quad \mathbf{v}^{\overline{\mathrm{top}}}\cdot\mathbf{N} = [\mathbf{v}]^{\mathrm{ff}}\cdot\mathbf{N} \quad \text{on} \quad \Gamma_{\mathrm{top}},$$

$$p^{\mathrm{bot}} = [p]^{\mathrm{pm}}, \quad \mathbf{v}^{\overline{\mathrm{bot}}}\cdot\mathbf{N} = [\mathbf{v}]^{\mathrm{pm}}\cdot\mathbf{N} \quad \text{on} \quad \Gamma_{\mathrm{bot}}.$$

The tangential component of Eq. (9) can be written as

$$- \left[2\mu\mathbf{D}^{\backslash}\right]^{\mathrm{ff}}\cdot\mathbf{N} = -b\hat{\mathbf{r}}_w^{\mathrm{top}}\cdot\left(\mathbf{v}^{\overline{\backslash\mathrm{top}}} - \mathbf{v}^{\backslash}\right), \tag{12}$$

where $\mathbf{D}^{\backslash} = \frac{1}{2}\left(\nabla^{\backslash}\mathbf{v} + \left(\nabla^{\backslash}\mathbf{v}\right)^{\mathrm{T}}\right)$ and $\mathbf{v}^{\backslash} = \mathbf{v} - \mathbf{N}\mathbf{N}\cdot\mathbf{v}$ are restricted to two dimensions in the transition region. The tangential component of Eq. (10) is given by

$$b\hat{\mathbf{r}}_w^{\mathrm{bot}}\cdot\left(\mathbf{v}^{\overline{\backslash\mathrm{bot}}} - \mathbf{v}^{\backslash}\right) = 0. \tag{13}$$

The tangential component of velocity at the transition region boundaries can be defined as

$$\mathbf{v}^{\overline{\backslash\mathrm{top}}} = \left[\mathbf{v}^{\backslash}\right]^{\mathrm{ff}}, \quad \mathbf{v}^{\overline{\backslash\mathrm{bot}}} = \left[\mathbf{v}^{\backslash}\right]^{\mathrm{pm}}. \tag{14}$$

Equation (12) together with condition (14) can be considered as the generalization of the Beavers–Joseph condition at the boundary between the free flow and transition region. Combining Eq. (12) and (13) together with the transition region momentum conservation Eq. (8), we get

$$- \left[2\mu\mathbf{N}\cdot\mathbf{D}^{\backslash}\right]^{\mathrm{ff}} = -b\varepsilon\left[\nabla^{\backslash}p - \varepsilon\hat{\mathbf{R}}_w\cdot\mathbf{v}^{\backslash}\right]^{\mathrm{tr}}.$$

## 5 Numerical Experiments

We consider flow domains $\Omega_{\mathrm{ff}} = [0, 5\mathrm{m}] \times [1, 2\mathrm{m}]$ and $\Omega_{\mathrm{pm}} = [0, 5\mathrm{m}] \times [0, 1\mathrm{m}]$ with the sharp interface $\Gamma = (0, 5\mathrm{m}) \times \{1\mathrm{m}\}$ and the transition region $\Omega_{\mathrm{tr}} = [0, 5\mathrm{m}] \times [0.98, 1\mathrm{m}]$, which partially occupies the porous medium layer. The fluid is water with density $\rho = 10^3 \left[\mathrm{kg/m}^3\right]$ and dynamic viscosity $\mu = 10^{-3}$ [Pa s]. The soil is isotropic with permeability $k = 10^{-7} \left[\mathrm{m}^2\right]$. The Beavers–Joseph coefficient is $\alpha_{\mathrm{BJ}} = 1$. The gravitational effects are neglected.

The boundary conditions are described in Fig. 2, where the inflow condition at the left boundary of the free flow domain reads $\mathbf{v} = (u, v) = (10(y - 1)(2 - y), 0)$

**Fig. 2** Coupled domain, locations of the interface and transition region, boundary conditions

[m/s], the no-flow conditions at the left and bottom boundaries of the porous medium domain are given by $\partial p / \partial \mathbf{n} = 0$, and the outflow condition is $\partial \mathbf{v} / \partial \mathbf{n} = \mathbf{0}$.

Second order in space finite volume schemes on staggered grids are considered in both flow domains [[11], Chaps. 4.4, 6.2, 6.3]. The fluid pressure is computed in the centers of the control volumes, and in addition at the interface and the external boundary of the porous medium domain. The velocities are computed in the centers of the control volume faces. The method is locally mass conservative and does not require any stabilization. In the porous medium domain, the pressure is the primary variable and the velocities are computed at the post-processing stage. The computational grids are uniform and conforming at the interfaces between the domains. For the numerical simulations of the steady-state coupled problem, the monolithic approach is applied: the systems of linear algebraical equations resulting from the discretization of the flow models are built together with the interface conditions into one matrix and solved simultaneously.

To compare the sharp interface and the transition region models, we plot the horizontal component of the velocity at the cross-section $x = 2.5$ [m] through the coupled domain for the sharp interface and transition region approach (Fig. 3). The porous medium velocity is of order $10^{-3}$ [m/s]. The horizontal component of the velocity computed through the sharp interface concept has a jump at the fluid-porous interface resulting from the Beavers–Joseph condition.

The transition region approach is a composition of three models: the free flow model, the porous medium model, and the transition region model. In addition to the sharp interface model, it contains the mass and momentum conservation equations of codimension one, therefore the CPU time for both models is essentially the same.

The advantage of the transition region approach is that the model is not restricted to the flow direction and the interface can be curved. It is especially important for modeling filtration processes where the flow is mainly perpendicular to the porous layer. The sharp interface concept is based on the Beavers–Joseph interface condition which is derived for flows parallel to the porous medium. In both models, the parameters should be estimated: in the sharp interface concept it is Beavers–Joseph coefficient and in the transition region model these are the resistance coefficients.

**Fig. 3** Velocity profiles in the coupled domain at $x = 2.5$ [m] for the sharp interface and transition region models

The numerical simulation results presented in Fig. 3 demonstrate the velocity jump for the sharp interface concept according to the Beavers–Joseph condition. The transition region velocity profile is smooth at the fluid-porous interface due to the considered transition zone.

## 6 Conclusions

In this work, we considered two coupling approaches (sharp interface, transition region) for isothermal single-fluid-phase porous medium and free flow systems. The proposed transition region model is a generalization of the well established sharp interface concept based on the Beavers–Joseph condition. Numerical simulation results demonstrate the velocity profiles in the coupled domain for both models.

Many extensions to this work are possible such as considering deformable porous materials, modeling species and energy transport, considering compressible fluids, taking into account moving interfaces and multiple fluid phases in porous media.

## References

1. Beavers, G., Joseph, D.: Boundary conditions at a naturally permeable wall. J. Fluid Mech. **30**, 197–207 (1967)
2. Cimolin, F., Discacciati, M.: Navier-Stokes/Forchheimer models for filtration through porous media. Appl. Numer. Math. **72**, 205–224 (2013)

3. Discacciati, M., Miglio, E., Quarteroni, A.: Mathematical and numerical models for coupling surface and groundwater flows. Appl. Num. Math. **43**, 57–74 (2002)
4. Goyeau, B., Lhuillier, D., Gobin, D., Velarde, M.: Momentum transport at a fluid-porous interface. Int. J. Heat Mass Transf. **46**, 4071–4081 (2003)
5. Gray, W., Miller, C.: Thermodynamically constrained averaging theory approach for modeling flow and transport phenomena in porous medium systems: 3. Single-fluid-phase flow. Adv. Water Res. **29**, 1745–1765 (2006)
6. Jackson, A., Rybak, I., Helmig, R., Gray, W., Miller, C.: Thermodynamically constrained averaging theory approach for modeling flow and transport phenomena in porous medium systems: 9. Transition region models. Adv. Water Res. **42**, 71–90 (2012)
7. Layton, W., Schieweck, F., Yotov, I.: Coupling fluid flow with porous media flow. SIAM J. Numer. Anal. **40**, 2195–2218 (2003)
8. Le Bars, M., Worster, M.: Interfacial conditions between a pure fluid and a porous medium: implications for binary alloy solidification. J. Fluid Mech. **550**, 149–173 (2006)
9. Mosthaf, K., Baber, K., Flemisch, B., Helmig, R., Leijnse, A., Rybak, I., Wohlmuth, B.: A coupling concept for two-phase compositional porous-medium and single-phase compositional free flow. Water Resour. Res. **47**, W10,522 (2011)
10. Saffman, R.: On the boundary condition at the surface of a porous medium. Stud. Appl. Math. **50**, 93–101 (1971)
11. Versteeg, H., Malalasekra, W.: An introduction to computational fluid dynamics: The finite volume method. Prentice Hall, NJ (2007)

# Convergence Analysis of a FV-FE Scheme for Partially Miscible Two-Phase Flow in Anisotropic Porous Media

**Bilal Saad and Mazen Saad**

**Abstract** We study the convergence of a combined finite volume nonconforming finite element scheme on general meshes for a partially miscible two-phase flow model in anisotropic porous media. This model includes capillary effects and exchange between the phase. The diffusion term, which can be anisotropic and heterogeneous, is discretized by piecewise linear nonconforming triangular finite elements. The other terms are discretized by means of a cell-centered finite volume scheme on a dual mesh. The relative permeability of each phase is decentred according the sign of the velocity at the dual interface. The convergence of the scheme is proved thanks to an estimate on the two pressures which allows to show estimates on the discrete time and compactness results in the case of degenerate relative permeabilities. A key point in the scheme is to use particular averaging formula for the dissolution function arising in the diffusion term. We show also a simulation of $CO_2$ injection in a water saturated reservoir and nuclear waste management. Numerical results are obtained by in-house numerical code.

## 1 Introduction

In nuclear waste management, an important quantity of hydrogen can be produced by corrosion of the steel engineered barriers (carbon steel overpack and stainless steel envelope) of radioactive waste packages. A direct consequence of this production is

B. Saad (✉)
Division of Computer, Electrical and Mathematical Sciences and Engineering,
King Abdullah University of Science and Technology, 4700 KAUST,
Thuwal 23955-6900, Kingdom of Saudi Arabia
e-mail: bilal.saad@kaust.edu.sa

M. Saad
Ecole Centrale de Nantes, 1, rue de la Noé, BP 92101, 44321 Nantes, France
e-mail: Mazen.Saad@ec-nantes.fr

the growth of hydrogen pressure around alveolus which can affect all the functions allocated to the canisters, waste forms, backfill, host rock. Host rock safety function may be threatened by over pressurisation leading to opening fractures of the domain, inducing groundwater flow and transport of radionuclides.

In this work, we address the construction and convergence analysis of a combined finite volume nonconforming finite element scheme, based on a two pressures formulation, for two–phase two–component flow in porous media where the dissolution of the non-wetting phase can occur in different engineering application (e.g. nuclear storage and $CO_2$ storage). The convergence analysis is done in the degenerate case and for the general model including capillarity and gravity effects.

## 2 Mathematical Formulation of the Continuous Problem

We consider herein a porous medium saturated with a fluid composed of two phases (liquid and gas) and a mixture of two components (water and hydrogen). The water is supposed only present in the liquid phase (no vapor of water due to evaporation). Let $T > 0$, let be $\Omega$ a bounded open subset of $\mathbb{R}^d$ ($d \geq 1$) and we set $Q_T = (0, T) \times \Omega$. We write the *mass conservation* of each component

$$\Phi \partial_t \left( \rho_l^w s_l \right) + \mathrm{div} \left( \rho_l^w \mathbf{V}_l \right) = f_w, \tag{1}$$

$$\Phi \partial_t \left( \rho_l^h(p_g) s_l + \rho_g^h(p_g) s_g \right) + \mathrm{div} \left( \rho_l^h(p_g) \mathbf{V}_l + \rho_g^h(p_g) \mathbf{V}_g \right)$$
$$- \mathrm{div} \left( \phi s_l \rho_l D_l^h \nabla X_l^h \right) = f_g, \tag{2}$$

where $\Phi(x)$, $s_\alpha(t, x)$ ($s_l + s_g = 1$), $p_\alpha(t, x)$, $\rho_l^h(p_g)$, $\rho_g^h(p_g)$, $\rho_\alpha = \rho_\alpha^h + \rho_\alpha^w$, $X_l^h = \rho_l^h / \rho_l$ ($X_l^h + X_l^w = 1$) and $D_l^h$ represent respectively the (given) porosity of the medium, the saturation of the $\alpha$ phase ($\alpha = l, g$), the pressure of the $\alpha$ phase, the density of dissolved hydrogen, the density of the hydrogen in the gas phase, the density of the $\alpha$ phase, the mass fraction of the hydrogen in the liquid phase, the diffusivity coefficient of the dissolved gas phase in the liquid phase. The velocity of each fluid $\mathbf{V}_\alpha$ is given by the Darcy law

$$\mathbf{V}_\alpha = -\mathbf{K} \frac{k_{r_\alpha}(s_\alpha)}{\mu_\alpha} \left( \nabla p_\alpha - \rho_\alpha(p_\alpha)\mathbf{g} \right),$$

where $\mathbf{K}(x)$ is the intrinsic (given) permeability tensor of the porous medium, $k_{r_\alpha}$ the relative permeability of the $\alpha$ phase, $\mu_\alpha$ the constant $\alpha$-phase's viscosity, $p_\alpha$ the $\alpha$-phase's pressure and $\mathbf{g}$ the gravity. For detailed presentation of the model we refer to the presentation of the benchmark Couplex-Gaz [4].

To define the hydrogen densities, we use the ideal gas law and the Henry law $\rho_g^h = \frac{M^h}{RT} p_g$, $\rho_l^h = M^h H^h p_g$, where the quantities $M^h$, $H^h$, $R$ and $T$ represent respectively the molar mass of hydrogen, the Henry constant for hydrogen, the

universal constant of perfect gases and $T$ the temperature. To close the system, we introduce the capillary pressure law which links the jump of pressure of the two phases to the saturation

$$p_c(s_l) = p_g - p_l, \tag{3}$$

the application $s_l \mapsto p_c(s_l)$ is decreasing. This model also corresponding to the application of $CO_2$ storage when hydrogen is replaced by $CO_2$.

Let $T > 0$ be the final time fixed, let be $\Omega$ a bounded open subset of $\mathbb{R}^d$ $(d \geq 1)$ where $\partial\Omega$ is $\mathscr{C}^1$. We set $\Sigma_T = (0, T) \times \partial\Omega$ and we note $\Gamma_l$ the part of the boundary of $\Omega$ where the liquid saturation is imposed to one and $\Gamma_n = \Gamma \backslash \Gamma_l$. The chosen mixed boundary conditions on the pressures are

$$\begin{cases} p_g(t, x) = p_l(t, x) = 0 & \text{on } (0, T) \times \Gamma_l, \\ \mathbf{V}_l \cdot \mathbf{n} = \mathbf{V}_g \cdot \mathbf{n} = \phi s_l \rho_l D_l^h \nabla X_l^h \cdot \mathbf{n} = 0 & \text{on } (0, T) \times \Gamma_n, \end{cases}$$

where $\mathbf{n}$ is the outward normal to $\Gamma_n$. The initial conditions are defined on pressures

$$p_\alpha(t = 0) = p_\alpha^0 \text{ in } \Omega, \text{ for } \alpha = l, g. \tag{4}$$

Next we introduce a classically physically relevant assumptions on the coefficients of the system.

(H1) **Degeneracy.** The functions $M_l = \frac{k_{r_l}}{\mu_l}$ and $M_g = \frac{k_{r_g}}{\mu_g} \in \mathscr{C}^0([0, 1], \mathbb{R}^+)$, $M_\alpha(s_\alpha = 0) = 0$ and there is a positive constant $m_0 > 0$ such that for all $s_l \in [0, 1]$,
$$M_l(s_l) + M_g(s_g) \geq m_0.$$

(H2) **Density Bounded.** The density $\rho_l^h$ is in $\mathscr{C}^1(\mathbb{R})$, increasing and there exists two positive constants $\rho_m > 0$ and $\rho_M > 0$ such that $0 < \rho_m \leq \rho_l^h(p_g) \leq \rho_M$.

(H3) The capillary pressure fonction $p_c \in \mathscr{C}^1([0, 1]; \mathbb{R}^+)$ and there exists $\underline{p_c} > 0$ such that $\frac{dp_c}{ds_l} \leq -\underline{p_c} < 0$.

(H4) The functions $f_w$, $f_g \in L^2(Q_T)$and $f_w$, $f_g \geq 0$ a.e. for all $(t, x) \in Q_T$.

(H5) $D_l^h$ is a possibly null positive constant.

This problem renews the mathematical and numerical interest in the equation describing multiphase multicomponent flows through porous media. Existence of weak solutions for the two compressible, partially miscible flow in porous media, under various assumptions on physical data, we refer to [1]. In [2] and [7] the authors study respectively the convergence of a combined FV-FE scheme of the Keller-Segel model and of a immiscible compressible two phase flows un porous media. Study of the convergence of a finite volume scheme for a model of miscible two-phase flow in porous media under non-degeneracy and regularization of the physical situation on

**Fig. 1** Triangles $K,L \in \mathcal{T}_h$ and dual volumes $D,E \in \mathcal{D}_h$ associated with edges $\sigma_D, \sigma_E \in \mathcal{E}_h$

the relative permeability of each phase which physically vanishes when its saturation goes to zero, we refer to [3].

## 3 Combined Finite Volume–Nonconforming Finite Element Scheme

### 3.1 Primal and Dual Meshes

We perform a triangulation $\mathcal{T}_h$ of the domain $\Omega$ such that $\overline{\Omega} = \cup_{K \in \mathcal{T}_h} K$. We denote by $\mathcal{E}_h$ the set of all sides, by $\mathcal{E}_h^{\text{int}}$ the set of all interior sides, by $\mathcal{E}_h^{\text{ext}}$ the set of all exterior sides, and by $\mathcal{E}_K$ the set of all the sides of an element $K \in \mathcal{T}_h$. We define $h := max\{\text{diam}(K), K \in \mathcal{T}_h\}$. We assume the following shape regularity: there exists a positive constant $\kappa_T$ such that

$$\min_{K \in \mathcal{T}_h} \frac{|K|}{\text{diam}(K)^d} \geq \kappa_T. \tag{5}$$

We also use a dual partition $\mathcal{D}_h$ of $\Omega$ such that $\overline{\Omega} = \cup_{D \in \mathcal{D}_h} D$. There is one dual element $D$ associated with each side $\sigma_D \in \mathcal{E}_h$. We construct it by connecting the barycenters of every $K \in \mathcal{T}_h$ that contains $\sigma_D$ through the vertices of $\sigma_D$. We denote by $Q_D$ the barycenter of the side $\sigma_D$. As for the primal mesh, we set $\mathfrak{F}_h, \mathfrak{F}_h^{\text{int}}, \mathfrak{F}_h^{\text{ext}}$ and $\mathfrak{F}_D$ for the dual mesh sides. We denote by $\mathcal{D}_h^{\text{int}}$ the set of all interior and by $\mathcal{D}_h^{\text{ext}}$ the set of all boundary dual volumes. We finally denote by $\mathcal{N}(D)$ the set of all adjacent volumes to the volume $D$, $\mathcal{N}(D) := \{E \in \mathcal{D}_h; \exists \sigma \in \mathfrak{F}_h^{\text{int}} \text{ such that } \sigma = \partial D \cap \partial E\}$. For $E \in \mathcal{N}(D)$, we also set $d_{K|L} := |Q_E - Q_D|$, $\sigma_{K|L} := \partial D \cap \partial E$ and $K_{D|E}$ the element of $\mathcal{T}_h$ such that $\sigma_{K|L} \subset K_{D|E}$.

We consider a uniform step time $\delta t$, and define $t^n = n\delta t$ for $n \in [0, N]$. We define the following finite-dimensional spaces:

$$X_h := \{\varphi_h \in L^2(\Omega); \varphi_h|_K \text{ is linear } \forall K \in \mathcal{T}_h,$$
$$\varphi_h \text{ is continuous at the points } Q_D, D \in \mathcal{D}_h^{\text{int}}\},$$

we equip $X_h$ with the seminorm $\|u_h\|_{X_h}^2 := \sum_{K \in \mathscr{T}_h} \int_K |\nabla u_h|^2 \, dx.$

## 3.2 The Combined Scheme

For clarity and simplicity, we restrict the theoretical demonstration to a horizontal field, i.e. we neglect the gravity effect. The Henry law combined to the ideal gas law, to obtain that the density of hydrogen gas is proportional to the density of hydrogen dissolved $\rho_g^h = \mathscr{C}_1 \rho_l^h$ where $\mathscr{C}_1 = \frac{1}{H_h RT}$. Remark that the density of water $\rho_l^w$ in the liquid phase is constant and from the Henry law, we can write $\rho_l \nabla X_l^h = X_l^w \nabla p_g$, where $\mathscr{C}_2$ is a constant equal to $H^h M^h$.

**Definition 1** (*Combined scheme*) The fully implicit combined finite volume-non conforming finite element scheme for the problem (1)–(2) reads: find the values $p_{\alpha,D}^n$, $D \in \mathscr{D}_h$, $n \in \{1, \cdots, N\}$, such that

$$p_{\alpha,D}^0 = \frac{1}{|D|} \int_D p_\alpha^0(x) dx, \quad s_{\alpha,D}^0 = \frac{1}{|D|} \int_D s_\alpha^0(x) dx, \quad \text{for all } D \in \mathscr{D}_h^{\text{int}}, \quad (6)$$

$$|D| \phi_D \frac{s_{l,D}^n - s_{l,D}^{n-1}}{\delta t} - \sum_{E \in \mathscr{N}(D)} M_l(s_{l,D|E}^n) \, \Lambda_{D,E} \, \delta_{D|E}^n(p_l) = \frac{f_{w,D}^n}{\rho_l^w}, \quad (7)$$

$$
\begin{aligned}
|D| \phi_D &\frac{\rho_l^h(p_{g,D}^n) m(s_{l,D}^n) - \rho_l^h(p_{g,D}^{n-1}) m(s_{l,D}^{n-1})}{\delta t} \\
&- \sum_{E \in \mathscr{N}(D)} (\rho_l^h)_{D|E}^n \, M_l(s_{l,D|E}^n) \, \Lambda_{D,E} \, \delta_{D|E}^n(p_l) \\
&- \mathscr{C}_1 \sum_{E \in \mathscr{N}(D)} (\rho_l^h)_{D|E}^n \, M_g(s_{l,D|E}^n) \, \Lambda_{D,E} \, \delta_{D|E}^n(p_g) \\
&- \mathscr{C}_2 \sum_{E \in \mathscr{N}(D)} \phi_D s_{l,D|E}^n (X_l^w)_{D|E}^n \, D_l^h \, \delta_{D|E}^n(p_g) = f_{g,D}^n, \quad (8)
\end{aligned}
$$

$$p_c(s_{l,D}^n) = p_{g,D}^n - p_{l,D}^n. \quad (9)$$

Where $m(s_l) = s_l + \mathscr{C}_1 s_g$. We refer to the matrix $\Lambda$ of elements $\Lambda_{D,E}$, $D, E \in \mathscr{D}_h^{\text{int}}$, as the diffusion matrix. The stiffness matrix of the nonconforming finite element method, is defined as follow

$$\Lambda_{D,E} := -\sum_{K \in \mathcal{T}_h} (\Lambda(x)\nabla\varphi_E, \nabla\varphi_D)_{0,K} \quad D, E \in \mathcal{D}_h. \tag{10}$$

The mean value of the density of each phase on interfaces is not classical since it is given as

$$\frac{1}{(\rho_l^h(p_g))_{D|E}^n} = \begin{cases} \dfrac{1}{p_{g,E}^n - p_{g,D}^n} \displaystyle\int_{p_{g,D}^n}^{p_{g,E}^n} \dfrac{1}{\rho_l^h(\zeta)}\,d\zeta & \text{if } p_{g,D}^n \neq p_{g,E}^n, \\[4mm] \dfrac{1}{\rho_l^h(p_{g,D}^n)} & \text{otherwise,} \end{cases} \tag{11}$$

this choice is crucial to obtain estimates on discrete pressures.

This scheme consists in a finite volume method together with a phase-by-phase upstream scheme. The implicit finite volume scheme satisfies industrial constraints of robustness and stability. In comparison with incompressible fluid, compressible fluids requires more powerful techniques. We show that the proposed scheme satisfy, a discrete energy estimate on the pressures and a function of the saturation that denote capillary terms, that allow us to derive the convergence of a subsequence to a weak solution of the continuous equations as the size of the discretization tends to zero. The treatment of the degeneracy needs the introduction of powerful technics to link the velocities to the global pressure and the capillary pressure on the discrete form [6].

### *3.3 A Priori Estimates and Convergence*

We summarize the main estimates:

**Proposition 1**    *1.  (Maximum principle). Let $(s_{\alpha,D}^0)_{D \in \mathcal{T}} \in [0, 1]$. Then, the saturation $(s_{l,D}^n)_{D \in \mathcal{T}, n \in \{0,\dots,N\}}$ is positive.*
  *2.  Assume that all transmissibilities are non-negative, i.e. $\Lambda_{D,E} \geq 0 \;\; \forall D \in \mathcal{D}_h^{int}, E \in \mathcal{N}(D)$. Then, the solution of the combined scheme satisfies*

$$\sum_{n=1}^{N} \delta t \sum_{D \in \mathcal{D}_h} \sum_{E \in \mathcal{N}(D)} \Lambda_{D,E} M_\alpha(s_{\alpha,D|E}^n)|p_{\alpha,E}^n - p_{\alpha,D}^n|^2 \leq C, \tag{12}$$

  *3.  The discrete global pressure satisfies*

$$\sum_{n=1}^{N} \delta t \, \|p_h\|_{X_h}^2 \leq C, \tag{13}$$

*where $p = p_g + \tilde{p}(s_l)$, and $\tilde{p}(s_l) = -\int_0^{s_l} \dfrac{M_l(z)}{M(z)} p_c^{'}(z)\mathrm{d}z.$*

**Table 1** Parameter values for the porous medium and fluid characteristics used in test case 1

| Porous medium | | Fluid characteristics | |
|---|---|---|---|
| Parameter | Value | Parameter | Value |
| $\Phi$ [-] | 0.15 | $D_l^h$ [m$^2$ s$^{-1}$] | $3 \times 10^{-9}$ |
| $\mathbf{K}$ [m$^2$] | $5.10^{-20}$ | $\mu_l$ [Pa s] | $1 \times 10^{-3}$ |
| $p_r$ [Pa] | $2 \times 10^6$ | $\mu_g$ [Pa s] | $9 \times 10^{-6}$ |
| $n$ [-] | 1.54 | $H^h$ [mol Pa$^{-1}$ m$^{-3}$] | $7.65 \times 10^{-6}$ |
| $s_{lr}$ [-] | 0.4 | $M^h$ [Kg mol$^{-1}$] | $2 \times 10^{-3}$ |
| $s_{gr}$ [-] | 0 | $\rho_l^w$ [Kg mol$^{-3}$] | $10^3$ |

To prove the estimate (12), we multiply (7) by $\mathscr{C}_1 p_{l,D}^n - p_{g,D}^n$ and (8) by the nonlinear function $g_g(p_{g,D}^n) = \int_0^{p_g} \frac{1}{\rho_l^h(z)} dz$, then summing the resulting equation over $D \in \mathscr{D}_h$ and $n \in \{1, \cdots, N\}$ to deduce the estimates on velocities. The estimates (13) is a consequence of the proof done in [6], the authors prove this property on primal mesh satisfying the orthogonal condition. This proof use only two neighbors elements and it is based only on the definition of the global pressure. Thus, the estimate (13) remains valid on the dual mesh, that allow us, based on the use of the Kolmogorov relative compactness theorem, to derive the convergence of these approximation to a weak solution of the continuous problem in this paper provided the mesh size and the time step tend to zero.

The main result of this paper is the following theorem.

**Theorem 1** *There exists an approximate solutions $(p_{\alpha,D}^n)_{n,D}$ corresponding to the system (7)–(8), which converges in $L^2(Q_T)$(up to a subsequence) to a weak solution $p_\alpha$ of the system (1)–(2).*

## 4 Numerical Results: Gas Phase (Dis)appearance (Quasi-1D)

In this section, we evaluate numerically the finite volume-nonconforming finite element method derived in the Sect. 3 on a test case dedicated to gas-phase (dis) appearance (see the Couplex-Gas benchmark [4] for more details). The method has been implemented into in-house Fortran code.

The porous medium and fluid characteristics are presented in [4] and summarized in Table. 1.

Initial conditions are $p_l(t = 0) = 10^6$ Pa and $p_g(t = 0) = 0$ Pa. For boundary conditions on the left, the hydrogen flow rate is given $q_h = 5.57 \times 10^{-6}\chi_{[0,T_{\text{inj}}]}(t)$ kg/m$^2$/year, where $\chi_{[0,T_{\text{inj}}]}$ denote the characteristic function of the set $[0, T_{\text{inj}}]$ and we impose a zero water flow rate $q_w = 0$. The Dirichlet boundary conditions for the outflow boundary are the same as the initial conditions.

A structured grid with $200 \times 20$ cells was used for the computations and we used a constant time step of 10 years. Figure 2 show the phase pressures, with respect to time (years) during and after injection. For $0 < t < 14 \times 10^3$ years, the gas saturation is

**Fig. 2** Liquid and gas pressures $p_l$ (*left*) and $p_g$ (*right*) at the (0, 10) with respect to time (years).

**Fig. 3** $CO_2$ phase saturation, color scale ranges from $s_\ell = 0$ (*blue*) to $s_\ell = \max(s_\ell)$ (*red*)



zero and the liquid pressure stay constant; the whole domain is saturated with water. For $14 \times 10^3 \le t \le 1.6 \times 10^5$ years, the gas phase appears. For $t > 5 \times 10^5$ years, the gas saturation decreases and after a while, the gas phase disappears. At the end of the simulation the system reaches a stationary state and the liquid pressure gradient goes to zero.

## 5 CO$_2$ Injection in a Fully Water Saturated Domaine

The Fig. 3 shows the $CO_2$ phase saturation at different time. $CO_2$ is injected into the lower left part of a rectangular geometry ($200 \times 50$ m) with a flux of $4.10^{-2}$ kg m$^{-2}$ s$^{-1}$. Densities, viscosities and all other parameters are chosen as suggested in [5]. In this example, we used the Brooks-Corey model for the soil water characteristic and relative permeabilities. The $CO_2$ migrates upwards until it reaches the top of the domain with the nonflux conditions and is then driven to the right by advective forces.

# References

1. Caro, F., Saad, B., Saad, M.: Study of degenerate parabolic system modeling the hydrogen displacement in a nuclear waste repository. Discrete Continuous Dyn. Syst. Ser. S **7**, 191–205 (2014)
2. Chamoun, G., Saad, M., Talhouk, R.: Monotone combined edge finite volume-finite element scheme for anisotropic Keller-Segel model. Numer. Methods Partial Differ. Eqn. 26 (2014). doi:10.1002/num.21858
3. Eymard, R., Schleper, V.: Study of a numerical scheme for miscible two-phase flow in porous media. hal-00741425, version 3 (2013)
4. MOMAS: http://www.gdrmomas.org
5. Neumann, R., Bastian, P., Ippisch, O.: Modeling and simulation of two-phase two-component flow with disappearing nonwetting phase. Comput. Geosci. **17**, 139–149 (2013)
6. Saad, B., Saad, M.: Study of full implicit petroleum engineering finite volume scheme for compressible two phase flow in porous media. Siam J. Numer. Anal. **51**(1), 716–741 (2013)
7. Saad, B., Saad, M.: A combined finite volume-nonconforming finite element scheme for compressible two phase flow in porous media. Numer. Math. (in revision) (2014)

# Piecewise Linear Transformation in Diffusive Flux Discretizations

**D. Vidović, M. Dotlić, B. Pokorni, M. Pušić and M. Dimkić**

**Abstract** A piecewise linear transformation that allows interpolation of diffused concentration over material discontinuities is presented. It may be used either to evaluate concentration values at auxiliary points, or to approximate face fluxes directly. It does not violate the discrete minimum and maximum principles, so it can be used to construct discretization schemes that preserve solution positivity or discrete minimum and maximum principles. The method has been demonstrated to produce second-order accurate interpolated concentration values and first-order accurate fluxes even when interpolation nodes lie at opposite sides of a discontinuity.

## 1 Introduction

Second-order terms play a role in a variety of partial differential equations. They are used to represent a number of unrelated physical phenomena such as molecular diffusion, heat conduction, dispersion, flow through porous media etc. For the sake of study we put aside the possible complexity of the physical system and consider the simplest linear diffusion equation, obtained by substituting Fick's law

D. Vidović (✉) · M. Dotlić · B. Pokorni · M. Dimkić
Jaroslav Černi Institute, Belgrade, Serbia
e-mail: dragan.vidovic@jcerni.co.rs

M. Dotlić
e-mail: milandotlic@gmail.com

B. Pokorni
e-mail: bpokorni.jci@gmail.com

M. Dimkić
e-mail: jdjcerni@jcerni.co.rs

M. Pušić
Faculty of Mining and Geology, University of Belgrade, Belgrade, Serbia
e-mail: mpusic@ptt.rs

$$\mathbf{u} = -\mathbb{D}\nabla C \tag{1}$$

into the continuity equation

$$\nabla \cdot \mathbf{u} = g, \tag{2}$$

where $\mathbf{u}$ is the diffusive velocity, $C$ is the concentration, $\mathbb{D}$ is the diffusion tensor which may be anisotropic and discontinuous, and $g$ is a source term. Portions of domain $\Omega$ in which $\mathbb{D}$ is continuous are called *material zones* and interfaces between them are called *material interfaces*.

We consider three types of boundary conditions:

$$C = g_D \quad \text{on } \Gamma_D, \tag{3}$$

$$\mathbf{u} \cdot \mathbf{n} = g_N \quad \text{on } \Gamma_N, \tag{4}$$

$$\mathbf{u} \cdot \mathbf{n} = \Psi(C - g_R) \quad \text{on } \Gamma_R, \tag{5}$$

where $\Gamma_D \cup \Gamma_R = \overline{\Gamma_D \cup \Gamma_R}$, $\Gamma_D \cup \Gamma_R \neq \emptyset$, $\Gamma_D \cup \Gamma_N \cup \Gamma_R = \partial\Omega$, and $\Gamma_D$, $\Gamma_N$, and $\Gamma_R$ are mutually disjoint.

We assume that the domain $\Omega$ is divided into polyhedral control volumes (cells) such that each cell is entirely contained in a single material zone. With each cell $T$ we associate one *collocation point* $\mathbf{x}_T$ (for example the centroid or the circumcenter) and one concentration value $C_T$.

Finite volume discretization is performed by integrating Eq. (2) over each cell and applying the divergence theorem:

$$\oint_T \mathbf{u} \cdot \mathbf{n} dS \equiv \sum_f \chi_{T,f} u_f = \int_T g dT, \qquad u_f = \int_f \mathbf{u} \cdot \mathbf{n}_f dS, \tag{6}$$

where $\mathbf{n}$ is the outward unit normal vector, the sum runs over all faces $f$ of polyhedron $T$, $\mathbf{n}_f$ is a fixed unit normal vector associated with face $f$, and $\chi_{T,f} = 1$ if $\mathbf{n}_f$ points outside of $T$, or $\chi_{T,f} = -1$ otherwise.

Further discretization requires that flux $u_f$ is represented using concentration values in some of the surrounding cells, and this is where finite volume schemes start to differ. For second order accuracy one wants to use linear interpolation to obtain the concentration gradient. This leads to

$$u_f \approx \sum_i \alpha_i (C_T - C_{T_i}), \tag{7}$$

Coefficients $\alpha_i$ can be found from the system

$$- \mathbb{D}\mathbf{n}_f = \sum_i \alpha_i (\mathbf{x}_T - \mathbf{x}_{T_i}). \tag{8}$$

Examples of second-order accurate schemes include the diamond scheme of which a review is found in [6] and schemes in which fluxes obtained using the linear interpolation are combined in a non-linear way to preserve the solution positivity [4, 7, 9, 10, 12, 14–17] or the discrete minimum and maximum principles [2, 3, 5, 8, 11, 13]. In this paper we do not present another finite volume scheme, but only a single building block—the interpolation—that may be used in these or other discretization schemes.

Boundary conditions can be used in this interpolation. Dirichlet boundary condition can be evaluated at any point $\mathbf{x}_f \in \Gamma_D$ and the obtained value $C_f$ can be used in (7). Neumann conditions can be resolved by introducing auxiliary collocation points at Neumann faces. If $f$ is such a face then auxiliary concentration value $C_f$ is explicitly computed from

$$u_f = \sum_i \alpha_i (C_f - C_{T_i}) \tag{9}$$

and used in (7) to compute fluxes through other faces. Robin condition can be treated similarly. An alternative to introducing auxiliary collocation points is to use the Neumann and Robin conditions directly, as this is done in this paper.

Difficulties arise when face $f$ is close to a material interface. Concentration gradient is discontinuous at the interface, so if concentration values associated with cells in different material zones are used as interpolation nodes then the accuracy is reduced. On the other hand schemes preserving the solution positivity or the discrete minimum and maximum principles require that coefficients $\alpha_i$ in (7) are non-negative, i.e. that vector $-\mathbb{D}\mathbf{n}$ is a conical combination of differences $\mathbf{x}_T - \mathbf{x}_{T_i}$. Sometimes this does not hold for any combination of collocation points within the same material zone.

One cure is to use harmonic averaging points introduced in [1]. These are special points at the material interface where any piecewise linear solution can be exactly represented as a convex combination of concentration values in only two collocation points at the opposite sides of the interface.

In this paper we present a unified method to treat discontinuities and include boundary conditions in concentration interpolation. An advantage over the harmonic points is that the interpolation can be performed over multiple material interfaces and that Neumann and Robin conditions in multiple material zones can be used as well. The interpolation method and the the piecewise linear interpolation it is based on were introduced in [16]. We also suggest two ways to use this interpolation method, one that was presented in [16], and another one explained here in Sect. 3.2.

The paper is organized as follows. The piecewise linear transformation is described in Sect. 2. Two alternative ways to use the transformation in flux discretization are presented in Sect. 3. A numerical example is given in Sect. 4.

**Fig. 1** Domain is assumed to have a locally layered structure



## 2 Piecewise Linear Transformation

We assume that the domain consists of layers with smooth interfaces as shown in Fig. 1. Even though authors usually do not state this requirement explicitly, corners in material interfaces generally introduce singularities with respect to the solution differentiability which reduce the accuracy of any scheme we are aware of. This is also the case with the presented interpolation method.

In each material zone $\Omega_i$ we represent the concentration locally as a linear function

$$C(\mathbf{x}) = C_i + \mathbf{G}_i \cdot (\mathbf{x} - \mathbf{x}_i), \tag{10}$$

where $\mathbf{x}_i$ are arbitrary nearby points chosen at material interfaces as in Fig. 1, and $\mathbf{G}$ is an unknown vector. Function $C(\mathbf{x})$ must satisfy two conditions:

1. It must be continuous at the interfaces;
2. Fluxes through interfaces must be continuous.

For each interface the first condition determines three degrees of freedom in function (10) and the second condition determines yet another degree of freedom. Thus the whole piecewise linear function $C(\mathbf{x})$ has only four remaining degrees of freedom and it can be reformulated as

$$C(\mathbf{x}) = C_0 + \mathbf{G}_0 \cdot F(\mathbf{x}), \tag{11}$$

where $F(\mathbf{x})$ is a piecewise linear transformation defined by conditions 1 and 2 and explicitly derived in [16]. Transformation $F(\mathbf{x})$ is completely determined by the geometry and diffusion tensors, i.e. it does not depend on the concentration.

If a material zone interface is not smooth at some point (for example because more then two zones meet there) the presented interpolation method is still applicable, but the accuracy is reduced, as demonstrated in [16]. However, points $\mathbf{x}_i$ must be chosen

such that the interface is smooth in $\mathbf{x}_i$ because a normal vector is associated with this point. Therefore in practice for $\mathbf{x}_i$ we take mesh face centroids.

## 3 Usage in Flux Discretization

Function $C(\mathbf{x})$ can be used in discretization of flux $u_f$ either to produce a concentration value in an auxiliary point used in (7) and (8), or directly to compute the flux $u_f = -\mathbb{D}\nabla C(\mathbf{x}_f)$.

### 3.1 Evaluation in Auxiliary Points

This case was explained in [16]. Free parameters $C_0$ and $\mathbf{G}_0$ are determined by imposing additional requirements that $C(\mathbf{x})$ matches concentration values at some collocation points or that it satisfies boundary conditions at some faces. Four equations of form

$$C(\mathbf{x}_{\mathscr{T}}) = C_{\mathscr{T}} \quad \text{where } \mathscr{T} \text{ is a mesh cell;} \tag{12}$$

$$C(\mathbf{x}_d) = g_D(\mathbf{x}_d) \quad \text{if } \mathbf{x}_d \in \Gamma_D; \tag{13}$$

$$-\mathbf{n}_f^T \mathbb{D}(\mathbf{x}_f)\nabla C(\mathbf{x}_f) = g_N(\mathbf{x}_f) \quad \text{if } \mathbf{x}_f \in \Gamma_N; \tag{14}$$

$$-\mathbf{n}_f^T \mathbb{D}(\mathbf{x}_f)\nabla C(\mathbf{x}_f) = \Psi(\mathbf{x}_f)\left(C(\mathbf{x}_f) - g_R(\mathbf{x}_f)\right) \quad \text{if } \mathbf{x}_f \in \Gamma_R \tag{15}$$

are chosen to form a linear system that determines $C_0$ and $\mathbf{G}_0$. For $\mathbf{x}_f$ we choose the centroid of face $f$.

When this system is solved, coefficients of linear function $C(\mathbf{x})$ are represented as linear combinations of concentration values and boundary fluxes. When $C(\mathbf{x})$ is evaluated at an auxiliary point, the concentration in this point is represented as a linear combination of the same. To satisfy the minimum and maximum principles, the coefficients of this linear combination must be non-negative. Thus one chooses from Eqs. (12)–(15) such that the value at the auxiliary point is evaluated as a convex combination of concentration values used in (12)–(15), with possible addition of Neumann faces contribution. As noted in [16], finding such equations can be a difficult task on distorted meshes, in particular if we choose to evaluate $C(\mathbf{x})$ in a mesh node, because in some cases collocation points that form a convex combination lie several cells away from this node. Nevertheless, it is possible to find such equations for any auxiliary point in most meshes.

## 3.2 Direct Usage

Let cells $T_1$ and $T_2$ share face $f$. To approximate the flux $u_f$ from $T_1$ to $T_2$ we take $\mathbf{x}_0 = \mathbf{x}_{T_1}$ and thus $C_0 = C_{T_1}$. To determine the remaining three degrees of freedom of $\mathbf{G}_0$ we form a linear system

$$A\mathbf{G}_0 = \mathbf{b} \tag{16}$$

by picking three equations of form (12)–(15), where $\mathbf{b}$ has components among concentration differences $C_{T_0} - C_\mathcal{G}$, $C_{T_0} - g_D(\mathbf{x}_d)$, $C_{T_0} - g_R(\mathbf{x}_f)$, and prescribed boundary fluxes $g_N(\mathbf{x}_f)$. The flux is computed as

$$u_f = -|f|\mathbf{n}_f^T \mathbb{D}\mathbf{G}_0 = -|f|\mathbf{n}_f^T \mathbb{D}A^{-1}\mathbf{b}. \tag{17}$$

In schemes preserving the positivity or the discrete minimum and maximum principles it is required that each component of $-|f|\mathbf{n}_f^T \mathbb{D}A^{-1}$ is non-negative. It must also be required that these components do not exceed a certain prescribed maximal value, otherwise an ill-conditioned matrix $A$ may degrade the interpolation accuracy.

This approach has advantages over auxiliary points. Cases when collocation points leading to non-negative coefficients are more than three cells away were not encountered here. On the contrary, collocation points necessary to represent a node value as a convex combination may be more than ten cells away on distorted grids. In addition, linear systems are $3 \times 3$, while with the auxiliary points they are $4 \times 4$. Thus this approach may require considerably less computational effort, especially if distorted grids are used.

Note that (17) is not a complete flux discretization and it should not be used directly in (6) because it is not conservative. When building a finite volume discretization, one-sided fluxes of form (17) are combined as in [2–5, 7–17] to yield a conservative scheme.

## 4 Example

Unit cube is divided in zones $\Omega_1 = \{(x, y, z)|x < 0.5\}$ and $\Omega_2 = \{(x, y, z)|x \geq 0.5\}$. Diffusion tensor is

$$\mathbb{D} = \begin{cases} \mathbb{D}_1 & \text{in } \Omega_1, \\ \mathbb{D}_2 & \text{in } \Omega_2, \end{cases} \tag{18}$$

$$\mathbb{D}_1 = \begin{bmatrix} 3 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbb{D}_2 = \begin{bmatrix} 10 & 3 & 0 \\ 3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \tag{19}$$

The exact concentration field is

$$C = \begin{cases} 1 - 2x^2 + 4xy + 2y + 6x & \text{in } \Omega_1, \\ 3.5 - 2y^2 + 2xy + 3y + x & \text{in } \Omega_2. \end{cases} \tag{20}$$

**Table 1** Errors of interpolated node values and face fluxes

| $h$ | 1/10 | 1/20 | 1/40 | 1/80 |
|---|---|---|---|---|
| $\varepsilon_2^C$ | $7.64 \times 10^{-4}$ | $1.84 \times 10^{-4}$ | $4.14 \times 10^{-5}$ | $9.83 \times 10^{-6}$ |
| $\varepsilon_{max}^C$ | $1.43 \times 10^{-2}$ | $4.61 \times 10^{-3}$ | $5.33 \times 10^{-4}$ | $1.34 \times 10^{-4}$ |
| $\varepsilon_2^{\mathbf{u}}$ | $3.70 \times 10^{-2}$ | $1.51 \times 10^{-2}$ | $7.48 \times 10^{-3}$ | $3.62 \times 10^{-3}$ |
| $\varepsilon_{max}^{\mathbf{u}}$ | 0.594 | 0.249 | 0.203 | 0.0961 |

We specify the exact solution at $z = 0$ and $z = 1$, the exact flux at $y = 0$ and $y = 1$, and Robin condition at the remaining boundary with $g_R = 0$ if $x = 0$, $g_R = 10$ if $x = 1$, and $\Psi$ chosen accordingly.

We use four independently generated unstructured tetrahedral grids. The mesh parameter $h$ is proportional to the longest edge length. In each mesh node $N$ the solution $C_N$ was represented as a convex combination of surrounding cell values. Flux through each face $f$ was represented as a conical combination of concentration differences and prescribed boundary fluxes. Cell concentration values were obtained from the analytic solution. The errors of the interpolated concentration

$$\varepsilon_2^C = \left[ \frac{\sum_N (C(\mathbf{x}_N) - C_N)^2}{\sum_N (C(\mathbf{x}_N))^2} \right]^{1/2}, \quad \varepsilon_{max}^C = \frac{\max_N |C(\mathbf{x}_N) - C_N|}{\left[ \sum_N (C(\mathbf{x}_N))^2 \right]^{1/2} / \sum_N 1}. \quad (21)$$

and of the interpolated flux

$$\varepsilon_2^{\mathbf{u}} = \left( \frac{\sum_f |f|(\mathbf{u}(\mathbf{x}_f) \cdot \mathbf{n}_f - u_f / |f|)^2}{\sum_f |f|(\mathbf{u}(\mathbf{x}_f) \cdot \mathbf{n}_f)^2} \right)^{1/2},$$

$$\varepsilon_{max}^{\mathbf{u}} = \frac{\max_f |\mathbf{u}(\mathbf{x}_f) \cdot \mathbf{n}_f - u_f / |f||}{\left( \sum_f |f|(\mathbf{u}(\mathbf{x}_f) \cdot \mathbf{n}_f)^2 / \sum_f |f| \right)^{1/2}} \quad (22)$$

are shown in Table 1.

The reported errors demonstrate that the proposed interpolation method generates second-order accurate node values and first-order accurate fluxes even when interpolation nodes belong to different material zones.

## 5 Conclusion

We have presented a piecewise linear interpolation method to be used in discretization of diffusive fluxes in discontinuous anisotropic environment. The method is second order accurate even when interpolation nodes are found at opposite sides of a material discontinuity. It does not violate the discrete minimum and maximum principles, so it can be used to construct schemes that preserve the positivity or the discrete maximum

and minimum principles. Unlike the harmonic points, this interpolation method can use boundary conditions in different material zones and interpolate over multiple discontinuities.

# References

1. Agelas, L., Eymard, R., Herbin, R.: A nine-point finite volume scheme for the simulation of diffusion in heterogeneous media. Comptes rendus de l'Académie des Sciences Mathématique **374**(11–12), 673–676 (2009)
2. Bertolazzi, E.: Discrete conservation and discrete maximum principle for elliptic pdes. Math. Mod. Meth. Appl. Sci. **8**(4), 685–711 (1998)
3. Bertolazzi, E., Manzini, G.: A second-order maximum principle preserving finite volume method for steady convection-diffusion problems. SIAM J. Numer. Anal. **43**(5), 2172–2199 (2005)
4. Danilov, A., Vassilevski, Y.: A monotone nonlinear finite volume method for diffusion equations on conformal polyhedral meshes. Russ. J. Numer. Anal. Math. Model. **24**(3), 207–227 (2009)
5. Droniou, J., Le Potier, C.: Construction and convergence study of schemes preserving the elliptic local maximum principle. SIAM J. Numer. Anal. **49**(2), 459–490 (2011)
6. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. Handb. Numer. Anal. **7**, 713–1018 (2000)
7. Le Potier, C.: Schéma volumes finis monotone pour des opérateurs de diffusions fortement anisotropes sur des maillages de triangle non structurés. C.R. Math. Acad. Sci. Paris **341**, 787–792 (2005)
8. Le Potier, C.: A nonlinear finite volume scheme satisfying maximum and minimum principles for diffusion operators. Int. J. Finite Vol. **6**(2), 1–20 (2009)
9. Lipnikov, K., Shashkov, M., Svyatskiy, D., Vassilevski, Y.: Monotone finite volume schemes for diffusion equations on unstructured triangular and shape-regular polygonal meshes. J. Comp. Phys. **227**(1), 492–512 (2007)
10. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: Interpolation-free monotone finite volume method for diffusion equations on polygonal meshes. J. Comp. Phys. **228**(3), 703–716 (2009)
11. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: Minimal stencil finite volume scheme with the discrete maximum principle. Russ. J. Numer. Anal. Math. Modelling **27**(7), 369–385 (2012)
12. Lipnikov, K., Svyatskiy, D., Vassilevski, Y.: Anderson acceleration for nonlinear finite volume scheme for advection-diffusion problems. SIAM J. Sci. Comput. **35**(2), A1120–A1136 (2013)
13. Sheng, Z., Yuan, G.: The finite volume scheme preserving extremum principle for diffusion equations on polygonal meshes. J. Comp. Phys. **230**(7), 2588–2604 (2011)
14. Vassilevski, Y., Kapyrin, I.: Two splitting schemes for nonstationary convection-diffusion problems on tetrahedral meshes. Comput. Math. Math. Phys. **48**(8), 1349–1366 (2008)
15. Vidović, D., Dimkić, M., Pušić, M.: Accelerated non-linear finite volume method for diffusion. J. Comp. Phys. **230**(7), 2722–2735 (2011)
16. Vidović, D., Dotlić, M., Dimkić, M., Pušić, M., Pokorni, B.: Convex combinations for diffusion schemes. J. Comp. Phys. **264**, 11–27 (2013)
17. Yuan, A., Sheng, Z.: Monotone finite volume schemes for diffusion equations on polygonal meshes. J. Comp. Phys. **227**(12), 6288–6312 (2008)

# Comparison of Two Approaches for Treatment of the Interface Conditions in FV Discretization of Pore Scale Models for Li-Ion Batteries

**Shiquan Zhang, Oleg Iliev, Sebastian Schmidt and Jochen Zausch**

**Abstract**  Pore scale models of Li-ions transport allow to gain insight into the details of the charge and discharge processes in Li-ion batteries. These models are diffusion type PDE-systems with very complex, nonlinear interface conditions on the interfaces between the active particles in the porous electrodes and the electrolyte. In this work, we discuss two approaches for the treatment of these interface conditions in conjunction with a cell-centered Finite Volume (FV) discretization of the governing equations. The first approach treats exactly the fluxes on the interface, but approximates the Butler-Volmer flux. This approach requires less memory because it does not introduce unknowns on the interface. The second approach introduces unknowns on the interface and discretizes the fluxes, but the Li-ion Butler-Volmer flux is evaluated exactly on the interface. Our numerical results show that the two approaches give very close results when the current rate is low. However, when the current rate becomes higher, the second approach is more accurate than the first one.

## 1 Introduction

Li-ion batteries used for technical applications are based on porous insertion electrodes. In most applications, the porous electrodes are random structures of active particles. During discharging Li-ions are de-intercalated from the anode particles into

S. Zhang (✉)
School of Mathematics, Sichuan University, Chengdu 610064, China
e-mail: shiquanz3@gmail.com

O. Iliev · S. Schmidt · J. Zausch
Fraunhofer ITWM, Fraunhofer-Platz 1, 67663  Kaiserslautern, Germany
e-mail: sebastian.schmidt@itwm.fraunhofer.de

O. Iliev
e-mail: oleg.iliev@itwm.fraunhofer.de

J. Zausch
e-mail: jochen.zausch@itwm.fraunhofer.de

**Fig. 1** Current density distribution of electronic current flowing through the solid phase (simulated with *BEST* [1])

the electrolyte and transported through the electrolyte to the porous cathode. There they intercalate into the cathode particles and are then transported via diffusion into their interior. It is well understood that the microstructure (e.g. size and arrangement of the active particles in the porous electrodes) significantly influence the performance of the battery. Going beyond porous structures, it has even been shown that specifically designed electrodes, can achieve a much larger power density [9], but still a lot of research is needed in order to quantitatively evaluate the influence of 3D structures. Available 3D microscopic models include ion transport in the electrolyte and in the solid particles, coupled with an equation for the potential [2, 3, 8]. Solving these models is only possible on cuts through the whole cell covering nevertheless the whole cathode anode direction [4, 7]. An example of pore scale simulation results for the isothermal part of the model from [2, 3], based on the discretization from [6], can be seen on Fig. 1.

The coefficients in the governing equations [2, 3, 6] usually experience jump on interface between the solid particles and the electrolyte. What is even more important, the interface conditions, that are modelled by the so-called Butler-Volmer expression, are highly nonlinear. The discretization of these interface conditions in [6] is done in a way that the Butler-Volmer flux is not evaluated exactly on the interface, instead, values of the concentration and the potential at the nearest volume are used. The discretization from [6] is successfully used in solving a number of academic and industrial problems. However, it is not theoretically investigated, and numerical investigation of its accuracy is desirable. In this work we compare the discretization from [6] with a new discretization approach, based on introducing new unknowns for the concentration and the potential on the interface.

The work is organized as follows. In Sect. 2, we describe the governing equations and interface conditions. In Sect. 3, we discuss the two ways to treat the interface conditions. Numerical results are given in Sect. 4 to compare the two methods and discuss their differences.

**Fig. 2** Our computational domain consists of an anode region $\Omega_a$, electrolyte region $\Omega_e$, and a cathode region $\Omega_c$ of thicknesses $L_a$, $L_e$, $L_c$, respectively



**Table 1** Coefficients

| Sub-domain | $\alpha$ | $\beta$ | $\lambda$ | $\gamma$ |
|---|---|---|---|---|
| $\Omega_a$ | $D_{s,a}$ | $0$ | $0$ | $\sigma_a$ |
| $\Omega_e$ | $D_e - \frac{t_+}{F}\kappa\frac{1-t_+}{F}RT\frac{1}{c}$ | $\frac{t_++\kappa}{F}$ | $\kappa\frac{1-t_+}{F}RT\frac{1}{c}$ | $\kappa$ |
| $\Omega_c$ | $D_{s,c}$ | $0$ | $0$ | $\sigma_c$ |

## 2 Microscopic Model

The real porous electrode geometry is like the one shown on Fig. 1, but for the needs of this paper it is enough to consider a simple geometry, e.g., assuming that the anode and the cathode consist of one particle each. Figure 2 is a simple sketch of such simplified geometry of lithium-ion battery, where $\Omega_a$ and $\Omega_c$ are anode and cathode respectively, and $\Omega_e$ denotes the electrolyte. To understand the charge and discharge dynamics of the battery, one needs to know the evolution of the concentration of lithium ions $c$ and the distribution of the potential $\phi$ in the whole battery domain. The governing equations for the microscopic model can be expressed generally as [2, 3, 6]

$$\frac{\partial c}{\partial t} + \nabla \cdot \mathbf{N} = 0, \tag{1}$$

$$\nabla \cdot \mathbf{j} = 0, \tag{2}$$

where the fluxes can be written as

$$\mathbf{N} = -(\alpha(c, \phi)\nabla c + \beta(c, \phi)\nabla\phi), \tag{3}$$

$$\mathbf{j} = -(\lambda(c, \phi)\nabla c + \gamma(c, \phi)\nabla\phi). \tag{4}$$

Here $\alpha$, $\beta$, $\lambda$ and $\gamma$ are, in general, nonlinear coefficients which have different form and values in different subregions. In the engineering and physics literature usually the model is written in a way which directly reflects the physics, but for us the above compact formulation is more convenient. The coefficients are given in Table 1, and for their physical meaning we refer to [2, 3, 6] and references therein.

The initial conditions for concentration are piecewise constant, i.e. they are $c_{a,0}$, $c_{e,0}$ and $c_{c,0}$ for anode, electrolyte and cathode respectively. The boundary conditions are given in Table 2. The Butler-Volmer interface conditions [5] are given as follows

**Table 2** Boundary conditions

| Boundary | Concentration | Potential |
|---|---|---|
| Upper | $\nabla c \cdot \hat{n} = 0$ | $\nabla \phi \cdot \hat{n} = 0$ |
| Lower | $\nabla c \cdot \hat{n} = 0$ | $\nabla \phi \cdot \hat{n} = 0$ |
| $\Gamma_a$ | $-D_{s,a} \nabla c \cdot \hat{n} = 0$ | $\phi = \phi_{fix}$ |
| $\Gamma_c$ | $-D_{s,a} \nabla c \cdot \hat{n} = 0$ | $-\sigma_c \nabla \phi \cdot \hat{n} = I$ |

$$\mathbf{N}_{\Gamma_+} \cdot \hat{n} = \mathbf{N}_{\Gamma_-} \cdot \hat{n} = i_{se}(c_+, c_-, \phi_+, \phi_-)/F, \quad \Gamma \in \{\Gamma_{ae}, \Gamma_{ec}\} \qquad (5)$$

$$\mathbf{j}_{\Gamma_+} \cdot \hat{n} = \mathbf{j}_{\Gamma_-} \cdot \hat{n} = i_{se}(c_+, c_-, \phi_+, \phi_-), \quad \Gamma \in \{\Gamma_{ae}, \Gamma_{ec}\}. \qquad (6)$$

In the above interface conditions, the direction of $\hat{n}$ is always pointing into the electrolyte, $c_+, c_-$ are the concentration values on the interface taken in the electrolyte and in the solid particle, respectively. Similar notation convention applies to potential $\phi$, particle flux $\mathbf{N}$ and current density $\mathbf{j}$. The subscript $e$ denotes electrolyte and $s$ denotes solid phase, the latter being anode or cathode in our simplified geometry. Furthermore, $i_{se}$ denotes the lithium ion flux across the interface and it has the following expression:

$$i_{se} = 2k\sqrt{c_+ c_-(c_{max} - c_-)} \sinh\left[\frac{F}{2RT}(\phi_- - \phi_+ - U_0)\right]. \qquad (7)$$

Here $U_0$ represents the open circuit potential, it is a function of the state of charge (SOC). Denoting $\theta := \frac{c_s}{c_{max}}$, the functional $U_0$ and the constant parameters are given in Table 3.

## 3 Discretization of the Interface Conditions

For the discretization of Eqs. (1) and (2), a standard lowest order cell centered finite volume method is adopted for spatial space, and first order backward Euler is used with uniform time steps. Note that our example is a simple 2D domain, uniform spatial mesh of rectangles are used, and any part of interfaces can only be the common edge of two rectangles. Omitting the obvious details, and using 1D notations for convenience, integration of (1) and (2) over a finite volume with center "$i$", after applying the divergence theorem, will result in

$$h\frac{c_i^{new} - c_i^{old}}{\Delta t} = \mathbf{N}_{i+0.5} - \mathbf{N}_{i-0.5}, \quad i = 1, 2, ...n_x, \qquad (8)$$

$$-\mathbf{j}_{i+0.5} + \mathbf{j}_{i-0.5} = 0, \quad i = 1, 2, ...n_x, \qquad (9)$$

where $h$ is the spatial size and $n_x$ is the number of finite volumes (grid cells) in $x-$ direction. In above, $\mathbf{N}_{i+0.5} = \overline{\alpha}\frac{c_{i+1} - c_i}{h} + \overline{\beta}\frac{\phi_{i+1} - \phi_i}{h}$ and $\mathbf{j}_{i+0.5} = \overline{\lambda}\frac{c_{i+1} - c_i}{h} + \overline{\gamma}\frac{\phi_{i+1} - \phi_i}{h}$,

**Table 3** Parameters

| Parameter | Value | Unit | Parameter | Value | Unit |
|---|---|---|---|---|---|
| $L_a$ | 0.0009 | cm | $L_c$ | 0.0009 | cm |
| $L_e$ | 0.0012 | cm | $H$ | 0.0024 | cm |
| $c_{a,0}$ | 2639e-6 | mol/cm$^3$ | $c_{c,0}$ | 20574e-6 | mol/cm$^3$ |
| $c_{a,max}$ | 24681e-6 | mol/cm$^3$ | $c_{c,max}$ | 23671e-6 | mol/cm$^3$ |
| $c_{c,0}$ | 1200e-6 | mol/cm$^3$ | $\phi_{fix}$ | $U_{0,a}(c_{a,0}/c_{a,max})$ | V |
| $F$ | 96487 | As/mol | $k_a$ | 0.002 | Acm$^{2.5}$/mol$^{1.5}$ |
| $R$ | 8.314 | J/mol/K | $k_c$ | 0.2 | Acm$^{2.5}$/mol$^{1.5}$ |
| $\sigma_a$ | 10 | S/cm | $\sigma_c$ | 0.38 | S/cm |
| $T$ | 298 | K | $t_+$ | 0.39989 | – |
| $D_{s,a}$ | 1e-10 | cm$^2$/s | $D_{s,c}$ | 1e-10 | cm$^2$/s |
| $D_e$ | 1.622e-6 | cm$^2$/s | $\kappa$ | 0.02 | S/cm |

$U_{0,a} = -0.132 + 1.41 \exp(-3.52\theta)$
$U_{0,c} = 4 + 0.07 \tanh(-22\theta + 12) - 0.1((1.002 - \theta)^{-0.37} - 1.6)$
$\qquad -0.045 \exp(-72\theta^8) + 0.01 \exp(-200(\theta - 0.19))$

where $\overline{\alpha}$ is the harmonic average of $\alpha(c_{i+1}, \phi_{i+1})$ and $\alpha(c_i, \phi_i)$, and similarly for $\overline{\beta}$, $\overline{\lambda}$ and $\overline{\gamma}$. The definition of $\mathbf{N}_{i-0.5}$ and $\mathbf{j}_{i-0.5}$ are similarly.

The discretization of the interface conditions is the main focus of this paper. For the treatment of the interface conditions (5) and (6), we consider two different ways. The first way is using the nearby volume's center values $c_s, c_e, \phi_s, \phi_e$ in the function $i_{se}$, see (7). This is the approach used in [6]. The second way is introducing new variables $c_-, c_+, \phi_-, \phi_+$ on the interfaces and discretizing the fluxes on the interface. As the interface conditions (5) and (6) are basically one dimensional along the normal direction of interfaces, we just show this treatment in the 1D case, the extension to multi dimensional case is straightforward.

Method 1: *exact fluxes, approximate evaluation of* $i_{se}$

Let us suppose that for some fixed $i$, an interface $\Gamma$ is located at $x_{i+0.5}$. In this case, in (8) the interface condition has to be used to determine $\mathbf{N}_{i+0.5}$. We use (5), but instead of evaluating $i_{se}$ using $c_-, c_+, \phi_-, \phi_+$, we evaluate it using $c_s, c_e, \phi_s, \phi_e$. In this particular case the subscripts $e, s$ stand for the centers of the near-interface finite volumes. The above explanation means e.g.,

$$\mathbf{N}_{\Gamma_-} \cdot \widehat{n} = \mathbf{N}_{i+0.5} \cdot \widehat{n} = i_{se}(c_+, c_-, \phi_+, \phi_-)/F \approx i_{se}(c_e, c_s, \phi_e, \phi_s)/F. \quad (10)$$
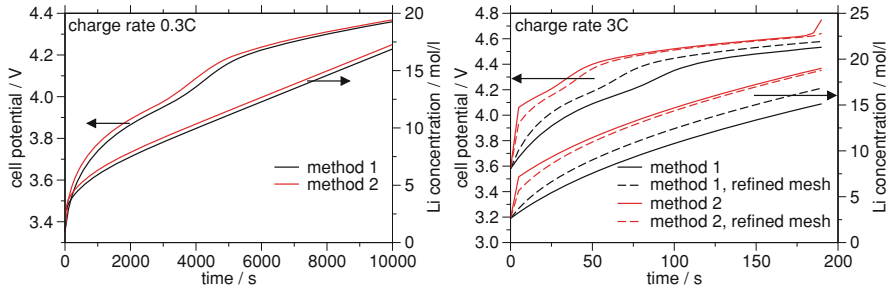
The flux $\mathbf{j}_{i+0.5}$ is treated in a similar way.

Method 2: *exact evaluation of* $i_{se}$, *approximate fluxes*

Here we introduce four unknowns on the interface, namely $c_-, c_+, \phi_-, \phi_+$. In this case we discretize the fluxes in the interface conditions (5) and (6) by using these new unknowns, i.e.
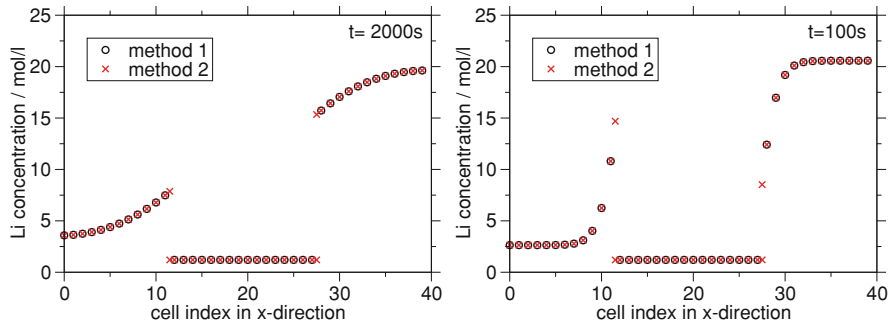
$$\begin{cases} -(\alpha(c_-,\phi_-)\frac{c_--c_s}{h/2} + \beta(c_-,\phi_-)\frac{\phi_--\phi_s}{h/2}) = i_{se}(c_+,c_-,\phi_+,\phi_-)/F, \\ -(\lambda(c_-,\phi_-)\frac{c_--c_s}{h/2} + \gamma(c_-,\phi_-)\frac{\phi_--\phi_s}{h/2}) = i_{se}(c_+,c_-,\phi_+,\phi_-), \\ -(\alpha(c_+,\phi_+)\frac{c_e-c_+}{h/2} + \beta(c_+,\phi_+)\frac{\phi_e-\phi_+}{h/2}) = i_{se}(c_+,c_-,\phi_+,\phi_-)/F, \\ -(\lambda(c_+,\phi_+)\frac{c_e-c_+}{h/2} + \gamma(c_+,\phi_+)\frac{\phi_e-\phi_+}{h/2}) = i_{se}(c_+,c_-,\phi_+,\phi_-). \end{cases} \tag{11}$$

## 4 Numerical Results

In this section, we use numerical results to discuss the accuracy of the above methods. Let us remark that generally the default number of volumes in $x$ direction are respectively 12, 16, 12 in $\Omega_a$, $\Omega_e$, $\Omega_c$ and 32 in $y$ direction, and the refined mesh is to double all these four numbers. To test the behavior at different values of the current density, we consider two examples. The first one is to apply a lower current density, i.e. $I = -0.0001$ A/cm$^2$, the other one is higher current density, i.e. $I = -0.001$ A/cm$^2$. These current densities correspond to charge rates of 0.3 and $3C$, respectively (with the usual definition of a $1C$-current corresponding to cell capacity divided by 1 h). For the lower current density, we take the total charge time to be $T_{charge} = 10000$ s and time step to be $\triangle t = 100$ s. For the higher current density, we take the total charge time to be $T_{charge} = 190$ s and time step to be $\triangle t = 5$ s. The results of these two examples are given in Fig. 3 where cell potential and solid concentration at the anode interface are displayed. The cell potential is the potential difference between cathode and anode boundary (since we fix the anode boundary potential it is just a point value at a special position, which is most important in the battery simulation). From the figures we can see that when the applied current density is low the results obtained with both methods are very close. However, at 3C the cell-potential is pronouncedly different for both methods which is consistent with the increased difference in solid interface concentration. Of course, due to the simplicity of the geometric model presented here, comparisons to experiment with the aim to judge which method is more exact are not meaningful. However, based on the following reasoning we can expect that method 2 is more exact: A snapshot of the concentration along the $x$-direction of the battery cell reveals that the concentration profile basically does not change with the two methods no matter whether the current is high or low (Fig. 4). The only difference is the additional concentration value directly on the interface in case of method two. This additional value continues the concentration trend according to the solid diffusion. Since our galvanostatic simulation setting basically fixes $i_{se}$ to the applied current density $I$ a difference in the concentration profile is not to be expected. However, the potential $\phi$ has to adjust such that the required $i_{se}$ is obtained. Since the potential depends on the open-circuit-potential $U_0(\theta)$ of the respective electrode, the potential indirectly depends on the solid interface concentration. Now at 0.3C the interface concentration does not differ much from the concentration of the nearest solid voxel. This is not true anymore at 3C since the concentration gradient is much steeper and therefore the

**Fig. 3** Cell potential and anode concentration used in interface conditions as function of time for charging rates 0.3C (*left*) and 3C (*right*). In the latter case comparison with a refined mesh is included (*dashed lines*)



**Fig. 4** Concentration snapshot at charging rates 0.3C (*left*) and 3C (*right*)

difference between interface and voxel concentration is much more pronounced. Via $U_0$ this couples back to the potential distribution and leads to the observed difference in cell-potential. From this we conclude that the discretization with additional interface unknowns can lead to more accurate results whenever a strong concentration gradient within the active particles is present (e.g. by slow diffusion or high applied current density).

In some sense this method is similar to increasing the spatial resolution. Therefore we have compared simulations of different mesh resolutions (right panel of Fig. 3). Indeed, higher resolution has the same effect on the cell potential as adding extra interface unknowns. As expected the effect of a higher resolution on the other hand is smaller when compared to the case with extra interface nodes.

In this context we also looked at the convergence behavior of both methods. To this end, we fix the number of volumes in $y$ direction to be 4, and double the number of volumes in $x$ direction in five steps, beginning with 12, 16, 12 volumes (for anode, electrolyte and cathode, respectively) and ending with 384, 512, 384 volumes. The meshes are denoted mesh 1 to mesh 6, respectively. We take the solution of "method 2" at the finest mesh (i.e. mesh 6) as the "exact" solution to evaluate the error, see

**Table 4** Maximal relative error of cell potential with respect to mesh 6 (with 1280 volumes in $x$-direction) for the two methods and charging rates

|  | Mesh 1 | Mesh 2 | Mesh 3 | Mesh 4 | Mesh 5 | Average order |
|---|---|---|---|---|---|---|
| Number of volumes in $x$ | 40 | 80 | 160 | 320 | 640 |  |
| Method 1 (0.3C) | 0.0097 | 0.0052 | 0.0027 | 0.0014 | 7.0128e-4 | 0.9476 |
| Method 1 (3C) | 0.0587 | 0.0366 | 0.0206 | 0.0107 | 0.0055 | 0.8540 |
| Method 2 (0.3C) | 0.0020 | 5.1991e-4 | 1.2758e-4 | 2.6692e-5 | 6.5943e-6 | 2.0608 |
| Method 2 (3C) | 0.0420 | 0.0150 | 0.0042 | 0.0011 | 2.1265e-4 | 1.9064 |

Table 4. Note that just for the maximal (at all time steps) relative error of cell potential, when the mesh is not too coarse, method 1 seems to be first order, and method 2 to be second order accurate. In the near future we will study the numerical behaviour of both methods more deeply. Furthermore the interplay between geometry, spatial resolution and additional interface unknowns will be investigated. The current study suggests already, that the additional technical complexity of method 2 may be justified by gaining a second order method also at the interface.

# References

1. Best—battery and electrochemistry simulation tool (2013). http://www.itwm.fraunhofer.de/best
2. Latz, A., Zausch, J.: Thermodynamic consistent transport theory of li-ion batteries. J. Power Sources **196**, 3296–3302 (2011)
3. Latz, A., Zausch, J., Iliev, O.: Modeling of species and charge transport in li-ion batteries based on non-equilibrium thermodynamics. Lect. Notes Comput. Sci. **6046**, 329–337 (2011)
4. Less, G.B., Seo, J.H., Han, S., Sastry, A.M., Zausch, J., Latz, A., Schmidt, S., Wieser, C., Kehrwald, D., Fell, S.: Micro-scale modeling of li-ion batteries: parameterization and validation. J. Electrochem. Soc. **159**(6), A697–A704 (2012)
5. Newman, J., Thomas-Alyea, K.E.: Electrochemical Systems. Wiley, New York (2004)
6. Popov, P., Vutov, Y., Margenov, S., Iliev, O.: Finite volume discretization of equations describing nonlinear diffusion in li-ion batteries. In: Dimov, I., Dimova, S., Kolkovska, N. (eds.) NMA 2010. LNCS, vol. 6046, pp. 338–346. Springer (2011)
7. Wang, C.W., Sastry, A.M.: Mesoscale modeling of li-ion polymer cell. J. Electrochem. Soc. **154**, A1035–A1047 (2007)
8. Wang, C.Y., Gu, W.B., Liaw, B.Y.: Micro-macroscopic coupled modeling of batteries and fuel cells. i. model development. J. Electrochem. Soc. **145**, 3407–3417 (1998)
9. Zhang, H., Yu, X., Braun, P.V.: Three-dimensional bicontinuous ultrafast-charge and -discharge bulk battery electrodes. Nat. Nanotechnol. **6**, 277–281 (2011)

# Part IV
# Hyperbolic Problems

# A Finite Volume Method for Large-Eddy Simulation of Shallow Water Equations

**Rajaa Abdellaoui, Fayssal Benkhaldoun, Imad Elmahi and Mohammed Seaid**

**Abstract** We present a robust finite volume method for large-eddy simulation of shallow water flows. The governing equations are derived from the Navier-Stokes equations with assumptions of shallow water flows including bed frictions and eddy viscosity. The turbulence effects are incorporated in the system by considering the Smagorinsky model. The numerical fluxes are reconstructed using a modified Roe's scheme that incorporates, in its reconstruction, the sign of the Jacobian matrix of the convective part of the large-eddy shallow water equations. The diffusion terms are discretized using a Green-Gauss diamond reconstruction. The proposed method is verified for the benchmark problem of flow around a circular cylinder.

## 1 Introduction

The description of the evolution of water flows in terms of water height and water velocity has proven to be very successful. Obviously, this description cannot be valid for very small scales at which molecular nature of the medium has to be taken into account. In the present work we consider the large-eddy simulation (LES) to model these small scales and also to analyse the subgrid errors. The basic idea of LES is to compute a space averaged flow field accurately. To achieve this, each flow variable

R. Abdellaoui · I. Elmahi
ENSAO Complex Universitaire, P.O 473, 60000 Oujda, Morocco
e-mail: ielmahi@ensa.univ-oujda.ac.ma

F. Benkhaldoun (✉)
LAGA, Université Paris 13, 99 Av J.B. Clement, 93430 Villetaneuse, France
e-mail: fayssal@math.univ-paris13.fr

M. Seaid
School of Engineering and Computing Sciences, University of Durham, South Road,
Durham  DH1 3LE, UK
e-mail: m.seaid@durham.ac.uk

$\omega$ is decomposed into a large-scale component $\overline{\omega}$ and a subgrid scale component $\omega'$. The large-scale component is obtained by the application of a filter operator, see for example [4]. We also introduce the mass-weighted (Favre) filtering of $\omega$ as

$$\widetilde{\omega} = \frac{\overline{h\omega}}{\overline{h}}.$$

By applying the filter operator to the standard shallow water equations, one obtains the depth-averaged equations

$$\partial_t \overline{h} + \partial_x (\overline{h}\widetilde{u}) + \partial_y (\overline{h}\widetilde{v}) = 0,$$

$$\partial_t (\overline{h}\widetilde{u}) + \partial_x \left( \overline{h}\widetilde{u}^2 + \frac{1}{2}g\overline{h}^2 \right) + \partial_y (\overline{h}\widetilde{u}\widetilde{v}) = -\frac{\overline{\tau}_{bx}}{\rho} + \nabla \cdot \left( (\nu + \nu_t) \nabla (\overline{h}\widetilde{u}) \right), \quad (1)$$

$$\partial_t (\overline{h}\widetilde{v}) + \partial_x (\overline{h}\widetilde{u}\widetilde{v}) + \partial_y \left( \overline{h}\widetilde{v}^2 + \frac{1}{2}g\overline{h}^2 \right) = -\frac{\overline{\tau}_{by}}{\rho} + \nabla \cdot \left( (\nu + \nu_t) \nabla (\overline{h}\widetilde{v}) \right),$$

where $t$ is the time variable, $\mathbf{x} = (x, y)^T$ the space coordinates, $\rho$ the water density, $g$ is the gravitational acceleration, $\nu$ is the kinematic viscosity of water, $h(t, x, y)$ is the water depth, $u(t, x, y)$ and $v(t, x, y)$ are the depth-averaged velocities in the $x$- and $y$-direction, respectively. In (1), $\nabla = (\partial_x, \partial_y)^T$ denotes the gradient operator and

$$\nu_t = \frac{(c_s \delta)^2}{\overline{h}} \sqrt{\partial_x (\overline{h}\widetilde{u})^2 + \frac{1}{2} \left( \partial_y (\overline{h}\widetilde{u}) + \partial_x (\overline{h}\widetilde{v}) \right) + \partial_y (\overline{h}\widetilde{v})^2},$$

where $c_s$ is a model constant which has to be chosen *a priori*, and $\delta$ is the grid filter width. Numerical tests in the literature with the Smagorinsky model use a Smagorinsky constant $c_s$ ranging from 0.01 to 0.1, see for example [6].

For free-surface water flows, the shear stresses are commonly modelled by the following quadratic friction law [5],

$$\overline{\tau}_{bx} = \rho C_f \frac{\widetilde{u}\sqrt{\widetilde{u}^2 + \widetilde{v}^2}}{2}, \qquad \overline{\tau}_{by} = \rho C_f \frac{\widetilde{v}\sqrt{\widetilde{u}^2 + \widetilde{v}^2}}{2}, \qquad (2)$$

where $C_f$ is the friction coefficient assumed to satisfy the semi-empirical law [5],

$$\frac{1}{\sqrt{C_f}} = -4\log \left( \frac{1.25}{4Re\sqrt{C_f}} \right), \qquad (3)$$

with $Re$ denotes the Reynolds number. In the current study we are interested in developing an unstructured finite volume method for solving the LES of shallow water Eq. (1). Numerical fluxes are reconstructed using the techniques used in [2] whereas, the finite volume discretization of the elliptic part in (1) is dealt with using a Green-Gauss diamond reconstruction. The performance of the present method is examined for the test problem of turbulent shallow water flows around a circular cylinder.

## 2 The Finite Volume Method

For simplicity in presentation, the LES shallow water system (1) can be rearranged in a conservative form as

$$\partial_t \mathbf{W} + \partial_x \left( \mathbf{F}(\mathbf{W}) - \tilde{\mathbf{F}}(\mathbf{W}) \right) + \partial_y \left( \mathbf{G}(\mathbf{W}) - \tilde{\mathbf{G}}(\mathbf{W}) \right) = \mathbf{S}(\mathbf{W}), \qquad (4)$$

where $\mathbf{W}$ and $\mathbf{S}$ are the vectors of conserved variables and source term, $\mathbf{F}$ and $\mathbf{G}$ are the convection tensor fluxes, $\tilde{\mathbf{F}}$ and $\tilde{\mathbf{G}}$ are the diffusion tensor fluxes

$$\mathbf{W} = \begin{pmatrix} h \\ hu \\ hv \end{pmatrix}, \quad \mathbf{F}(\mathbf{W}) = \begin{pmatrix} hu \\ hu^2 + \frac{1}{2}gh^2 \\ huv \end{pmatrix}, \quad \mathbf{G}(\mathbf{W}) = \begin{pmatrix} hv \\ huv \\ hv^2 + \frac{1}{2}gh^2 \end{pmatrix},$$

$$\mathbf{S}(\mathbf{W}) = \begin{pmatrix} 0 \\ -\dfrac{\tau_{bx}}{\rho} \\ -\dfrac{\tau_{by}}{\rho} \end{pmatrix}, \quad \tilde{\mathbf{F}}(\mathbf{W}) = \begin{pmatrix} 0 \\ (v + v_t)\,\partial_x\,(hu) \\ (v + v_t)\,\partial_x\,(hv) \end{pmatrix}, \quad \tilde{\mathbf{G}}(\mathbf{W}) = \begin{pmatrix} 0 \\ (v + v_t)\,\partial_y\,(hu) \\ (v + v_t)\,\partial_y\,(hv) \end{pmatrix},$$

where the "overbar" and "Favre", used to refer to filtered variables, has been omitted for ease in notation. The integral form of the Eq. (4) over a fixed volume $V$ is given by

$$\partial_t \int_V \mathbf{W}\, dV + \int_V \left( \partial_x \left( \mathbf{F}(\mathbf{W}) - \tilde{\mathbf{F}}(\mathbf{W}) \right) + \partial_y \left( \mathbf{G}(\mathbf{W}) - \tilde{\mathbf{G}}(\mathbf{W}) \right) \right) dV = \int_V \mathbf{S}(\mathbf{W})\, dV,$$

that, using the divergence theorem for the second integral leads to

$$\partial_t \int_V \mathbf{W}\, dV + \int_{\partial V} \mathscr{F}(\mathbf{W}; \mathbf{n})\, d\sigma - \int_{\partial V} \tilde{\mathscr{F}}(\mathbf{W}; \mathbf{n})\, d\sigma = \int_V \mathbf{S}(\mathbf{W})\, dV, \quad (5)$$

where

$$\mathscr{F}(\mathbf{W}; \mathbf{n}) = \mathbf{F}(\mathbf{W})n_x + \mathbf{G}(\mathbf{W})n_y, \qquad \tilde{\mathscr{F}}(\mathbf{W}; \mathbf{n}) = \tilde{\mathbf{F}}(\mathbf{W})n_x + \tilde{\mathbf{G}}(\mathbf{W})n_y,$$

and $\partial V$ is the surface surrounding the volume $V$.

Research on numerical solution of shallow water equations has received considerable attention during the last decades and several finite volume methods have been developed in the literature. In the current work, we consider a finite volume method based on the sign matrix developed and analyzed in [2] among others. The main advantages of this method lies on its implementation on unstructured triangular meshes and preserving conservation properties of the equations. Hence, using the

**Fig. 1** A generic control
volume and notations. The
co-volume $coV_{ij}$ is limited
by the *gray* area, $\mathbf{W}_i$ and $\mathbf{W}_j$
are the solution vectors at the
control volume $\mathscr{T}_i$ and $\mathscr{T}_j$,
respectively



control volume depicted in Fig. 1, a finite volume discretization of (5) yields

$$\mathbf{W}_i^{n+1} = \mathbf{W}_i^n - \frac{\Delta t}{|\mathscr{T}_i|} \sum_{j \in N(i)} \int_{\Gamma_{ij}} \mathscr{F}(\mathbf{W}^n; \mathbf{n}) \, d\sigma + \frac{\Delta t}{|\mathscr{T}_i|} \sum_{j \in N(i)} \int_{\Gamma_{ij}} \tilde{\mathscr{F}}(\mathbf{W}^n; \mathbf{n}) \, d\sigma$$

$$+ \frac{\Delta t}{|\mathscr{T}_i|} \int_{\mathscr{T}_i} \mathbf{S}(\mathbf{W}^n) \, dV, \tag{6}$$

where $N(i)$ is the set of neighboring triangles of the cell $\mathscr{T}_i$, $\Gamma_{ij}$ is the interface
between the two control volumes $\mathscr{T}_i$ and $\mathscr{T}_j$, $\mathbf{W}_i^n$ is an average value of the solution
$\mathbf{W}$ in the cell $\mathscr{T}_i$ at time $t_n$,

$$\mathbf{W}_i = \frac{1}{|\mathscr{T}_i|} \int_{\mathscr{T}_i} \mathbf{W} \, dV,$$

where $|\mathscr{T}_i|$ denotes the area of $\mathscr{T}_i$ and $\partial V$ is the surface surrounding the control
volume $V$. Here, $\mathbf{n} = (n_x, n_y)^T$ denotes the unit outward normal to the surface $\partial V$.
Following the formulation in [2], the proposed finite volume scheme consists of a
predictor stage and corrector stage as

$$\mathbf{W}_{ij}^n = \frac{1}{2} \left( \mathbf{W}_i^n + \mathbf{W}_j^n \right) - \frac{1}{2} \operatorname{sgn} \left[ \nabla \mathscr{F} \left( \overline{\mathbf{W}}_{ij}^n; \mathbf{n}_{ij} \right) \right] \left( \mathbf{W}_j^n - \mathbf{W}_i^n \right),$$

$$\tag{7}$$

$$\mathbf{W}_i^{n+1} = \mathbf{W}_i^n - \frac{\Delta t}{|\mathscr{T}_i|} \sum_{j \in N(i)} \mathscr{F} \left( \mathbf{W}_{ij}^n; \mathbf{n}_{ij} \right) |\Gamma_{ij}| + \Delta t \mathbf{S}_i^n,$$

with sgn [**A**] denotes the sign matrix of **A** and $\overline{\mathbf{W}}_{ij}^n$ is the Roes average state. A detailed
formulation of the sign matrix can be found in [2] and it will not be repeated here.

To discretize the diffusion fluxes in (6) we adapt a Green–Gauss diamond recon-
struction, see for example [1] and further references are therein. This method has
been selected because it can be applied on general unstructured grids, it does not
require serious restrictions on the angles of triangles, and it can be easily incorpo-
rated in our finite volume scheme. Hence, a co-volume, $coV_{ij}$, is first constructed by
connecting the barycentres of the elements that share the edge $\Gamma_{ij}$ and its endpoints

as shown in Fig. 1. Then, in the $x$-direction, diffusion fluxes are evaluated at an inner edge $\Gamma_{ij}$ as

$$\int_{\Gamma_{ij}} \nu \partial_x (hu) n_x \, d\sigma = \frac{\bar{\nu}|_{\Gamma_{ij}}}{|coV_{ij}|} \sum_{\varepsilon \in \partial coV_{ij}} \frac{(hu)_{N_1} + (hu)_{N_2}}{2} \int_{\varepsilon} n_{x\varepsilon} \, d\sigma, \quad (8)$$

where $N_1$ and $N_2$ are the nodes of the edge $\varepsilon$ on the surface $\partial coV_{ij}$, $(hu)_{N_1}$ and $(hu)_{N_2}$ are the values of the discharge $(hu)$ in the node $N_1$ and $N_2$, respectively. In (8), the diffusion coefficient $\bar{\nu}$ is defined by

$$\bar{\nu} = \frac{\nu_{N_1} + \nu_{N_2} + \nu_{N_3} + \nu_{N_4}}{4},$$

with $\nu_{N_k}$, $k = 1, \ldots, 4$, are values of the diffusion coefficient $\nu$ at the co-volume nodes $N_k$ approximated by linear interpolation from the values on the cells sharing the same vertex $N_k$.

## 3 Numerical Results

We present numerical results for the test example of LES shallow water flows over a circular cylinder. The main goals of this section are to illustrate the numerical performance of the finite volume method described above and to numerically verify its capability to solve turbulent shallow water problems. In all the computations reported herein, $c_s = 0.03$ and we used variable time stepsizes $\Delta t$ adjusted at each step according to

$$\Delta t = Cr \cdot \min (\Delta t_{\text{conv}}, \Delta t_{\text{diff}}),$$

where

$$\Delta t_{\text{conv}} = \min_{\Gamma_{ij}} \left( \frac{|\mathscr{T}_i| + |\mathscr{T}_j|}{2 |\Gamma_{ij}| \max_p |(\lambda_p)_{ij}|} \right), \quad \Delta t_{\text{diff}} = \min_{\Gamma_{ij}} \left( \frac{|\mathscr{T}_i|}{2 (\nu + \nu_t)_{ij}} \right)$$

with $\Gamma_{ij}$ is the edge between two cells $\mathscr{T}_i$ and $\mathscr{T}_j$, and $Cr$ is the Courant number set to 0.8 to ensure stability. The gravitational acceleration is fixed to $g = 9.81 m/s^2$. A schematic of the system considered in the present work is shown in Fig. 2. The system consists of a shallow water flow in a channel containing a circular cylinder. A similar test problem has been investigated in [3]. Here, the channel width is $L$, the channel height is $H$ and the diameter of cylinder is $D$. A water flow enters through the left boundary of channel with uniform velocity $u_\infty$. The Reynolds number for this problem is defined as $Re = Du_\infty/\nu$. Here, the governing equations (4) are solved on a computational domain $\Omega$ with smooth boundary $\partial\Omega = \Gamma_w \cup \Gamma_{\text{in}} \cup \Gamma_{\text{out}}$

**Fig. 2** Configuration of the flow around a cylinder



**Fig. 3** A zoom on Mesh 1 (*left plot*) and Mesh 2 (*right plot*) used in the simulation

shown in Fig. 2, and subject to the following boundary conditions

$$
\begin{aligned}
\mathbf{u}(t, \hat{\mathbf{x}}) &= \mathbf{0}, & \forall \, \hat{\mathbf{x}} \in \Gamma_{\mathrm{w}}, \\
u(t, \hat{\mathbf{x}}) &= u_{\infty}, & \forall \, \hat{\mathbf{x}} \in \Gamma_{\mathrm{in}}, \\
\mathbf{n}(\hat{\mathbf{x}}) \cdot \nabla \mathbf{u}(t, \hat{\mathbf{x}}) &= \mathbf{0}, & \forall \, \hat{\mathbf{x}} \in \Gamma_{\mathrm{out}},
\end{aligned}
\tag{9}
$$

for the flow and

$$
\mathbf{n}(\hat{\mathbf{x}}) \cdot \nabla h(t, \hat{\mathbf{x}}) = \mathbf{0}, \qquad \forall \, \hat{\mathbf{x}} \in \partial \Omega,
\tag{10}
$$

for the water height. Here, $\mathbf{u} = (u, v)^T$ and $\mathbf{n}(\hat{\mathbf{x}})$ denotes the outward unit normal in $\hat{\mathbf{x}}$ with respect to $\partial \Omega$. In all our computations we set $L = 1.25\,\mathrm{m}$, $H = 0.6\,\mathrm{m}$, $D = 6.3\,\mathrm{cm}$, $u_{\infty} = 14.3\,\mathrm{cm/s}$ and initially the water height is set to $h = 3.8\,\mathrm{cm}$. We perform computations with triangular finite volumes using the unstructured meshes shown in Fig. 3.

In Fig. 4 we display the snapshots of the $u$-velocity and $v$-velocity. In this figure we also include the transport of two passive tracers injected at the upstream side of the cylinder. The presented results indicate circulation zones moving downstream for both meshes. The results also indicate that refining the mesh, alters the flow features and also the tracer distribution past the cylinder. For instance, the size of the recirculation zones increases with the flow exhibiting eddies with different magnitudes and separating shear layers. We can see the small complex structures of the flow being captured by the proposed finite volume method.

Figure 5 presents cross-sections of the $u$-velocity and $v$-velocity at three different locations within the channel using the unstructured grids listed in Table 1. The vertical sections have been located upstream at $x = -0.1$, right behind the cylinder at $x = 0.04$ and downstream at $x = 0.1$. It is clear that the flow structures differ from

**Fig. 4** Snapshots of $u$-velocity (*first row*), $v$-velocity (*second row*) and tracer (*third row*). Results obtained using a coarse mesh Mesh 1 (*left column*) and using a fine mesh Mesh 3 (*right column*)



**Fig. 5** Cross-sections of the $u$-velocity (*first row*) and $v$-velocity (*second row*) at different locations in the channel. Here $x = -0.1$ (*first column*), $x = 0.04$ (*second column*) and $x = 0.1$ (*third column*)

**Table 1** Comparison results for the LES shallow water flows over a circular cylinder on three unstructured meshes

|        | # of elements | # of nodes | min $h$ | max $h$ | min $u$ | max $u$ |
|--------|---------------|------------|---------|---------|---------|---------|
| Mesh 1 | 19256         | 9853       | 0.241   | 0.245   | –0.069  | 0.229   |
| Mesh 2 | 41310         | 20980      | 0.219   | 0.223   | –0.082  | 0.242   |
| Mesh 3 | 81330         | 40550      | 0.213   | 0.218   | –0.095  | 0.256   |

one location to another and strongly depend on the mesh considered in the simulation. It is also evident that the finer mesh Mesh 3 would produce more accurate results that the coarse mesh Mesh 1. However, the results on the Mesh 3 and Mesh 2 demonstrate similar trends. This can clearly been seen in Table where minimum and maximum values of the water height and flow field are summarized.

## 4 Conclusions

In the present study, the turbulent fow past a circular cylinder is numerically solved by a robust finite volume method. The method uses shallow water assumptions and the Smagorinsky model in the governing equations and its belongs to the class of fractional step procedures where the convection part and diffusion part are discretized on separated control volumes. Conservative reconstruction of numerical fluxes is achieved thanks to the the sign of the Jacobian matrix of the convective part of the large-eddy shallow water equations. The numerical simulations are performed and comparisons are presented for simulations on different unstructured meshes. The presented results demonstrate the capability of the finite volume method that can provide insight into complex shallow water fow behaviors.

## References

1. Benkhaldoun, F., Elmahi, I., Seaid, M.: Well-balanced finite volume schemes for pollutant transport by shallow water equations on unstructured meshes. J. Comput. Phys. **226**, 180–203 (2007)
2. Benkhaldoun, F., Elmahi, I., Seaid, M.: A new finite volume method for flux-gradient and source-term balancing in shallow water equations. Comput. Methods Appl. Mech. Eng. **199**, 49–52 (2010)
3. Hinterberger, C., Fröhlich, J., Rodi, W.: Three-dimensional and depth-averaged large-eddy simulations of some shallow water flows. J. Hydraul. Eng. 857–872 (2007)
4. Sagaut, P.: Large Eddy Simulation for Incompressible Flows. Springer, Berlin (2001)
5. Schlichting, H.: Boundary Layer Theory. McGraw-Hill, New York (1968)
6. Smagorinsky, J.: General circulation experiments with the primitive equations. Mon. Weather Rev. **91**, 99–164 (1963)

# An Asymptotic Preserving Scheme for the Barotropic Baer-Nunziato Model

**Rémi Abgrall and Sophie Dallet**

**Abstract** We introduce in this paper a new scheme for obtaining approximations of solutions of the barotropic Baer-Nunziato (BN) model. This scheme is expected to provide relevant approximations when relaxation time scales embedded in pressure and velocity relaxation terms vanish. A brief recall of the BN model and the asymptotic model is first given. The scheme and its main properties are described and some numerical results are provided confirming that it behaves reasonably well.

## 1 Introduction

The mathematical and numerical modelling of two-phase flows is a widely debated topic. Depending on applications, the single-fluid or the two-fluid approach may be preferred. An advantage of the former is its simplicity and computational efficiency, whereas a probable drawback of the two-fluid formalism is that the use of high-order schemes is mandatory in order to obtain decent approximations of solutions. For flows involving non-negligible relative velocities, the two-fluid approach is mandatory. Actually, for some applications, a hybrid approach seems rather appealing, but this in turn requires the development of a consistent approach, which means retrieving at least the main patterns from the—simpler—single fluid model when some ade-

R. Abgrall
Institüt für Mathematik, Universität Zürich, Winterthurerstrasse 190,
CH-8057 Zürich, Switzerland

S. Dallet (✉)
Fluid Dynamics, Power Generation and Environment, EDF R&D, 6 quai Watier,
F78400 Chatou, France
e-mail: sophie.dallet@edf.fr

S. Dallet
Laboratoire d'Analyse Topologie Probabilités, UMR CNRS 7353, 39 rue Joliot Curie,
13453 Marseille cedex 13, France

quate parameters are tuned to 0 in the—expected more complex—two-fluid model. The two-fluid model investigated in this paper is a barotropic version of the original Baer-Nunziato (BN) model [2]. This five-equation model has been investigated quite recently in [7] and [3] for instance. The authors of these two references respectively propose a well-balanced scheme and a relaxation scheme in order to compute approximations of solutions of the barotropic BN model. In this paper, we wish to construct a scheme that preserves the asymptotic regime when the relaxation time scales that are active in so-called source terms tend to vanish (see [4]). Another aim is to obtain a sufficiently cheap and accurate algorithm. Hence, after a very brief description of the whole model, we will introduce the scheme and give its main properties; next, we will present some computational results in order to evaluate the capabilities of the scheme when simulating a Riemann problem and also to examine the behaviour of the scheme when some small parameter tends to 0.

## 2 The Barotropic Baer-Nunziato Model

We recall the governing equations of the barotropic Baer-Nunziato model. We denote as usual $\alpha_k$, $\rho_k$ and $u_k$ the statistical fraction, density and velocity within phase $k$, such that $\alpha_1 + \alpha_2 = 1$, $\alpha_k \in ]0; 1[$, and $m_k = \alpha_k \rho_k$ stands for the mass fraction in phase $k$. The model includes pressure and velocity relaxation terms, with corresponding relaxation time scales embedded in $\Theta(W) > 0$ and $\Lambda(W) > 0$. We also set $V_I = \beta u_1 + (1 - \beta)u_2$ and $P_I = (1 - \beta)P_1 + \beta P_2$, where $\beta = 0$ or 1, and define relative pressure $P_r = P_2 - P_1$ and velocity $u_r = u_2 - u_1$. In this barotropic formulation, the pressure $P_k$ is an increasing function of $\rho_k$: $P_k = \mathscr{P}_k(\rho_k)$. Thus the system reads:

$$
\begin{cases}
\dfrac{\partial \alpha_2}{\partial t} + V_I \dfrac{\partial \alpha_2}{\partial x} & = \Theta(W)(P_2 - P_1) \\[2mm]
\dfrac{\partial}{\partial t}(\alpha_k \rho_k) + \dfrac{\partial}{\partial x}(\alpha_k \rho_k u_k) & = 0 \\[2mm]
\dfrac{\partial}{\partial t}(\alpha_k \rho_k u_k) + \dfrac{\partial}{\partial x}(\alpha_k \rho_k u_k^2) + \dfrac{\partial}{\partial x}(\alpha_k P_k) - P_I \dfrac{\partial \alpha_k}{\partial x} & = (-1)^{k+1} \Lambda(W)|u_r|u_r
\end{cases}
\tag{1}
$$

The convective part of this system is hyperbolic. We recall that the eigenvalues are: $\lambda_1 = V_I$, $\lambda_{2-5} = u_k \pm c_k$; moreover the set of right eigenvectors spans the whole space $\mathbb{R}^5$ unless $|u_k - V_I| = c_k$. Fields associated with eigenvalues $\lambda_{2-5}$ are genuinely non linear; the 1−wave is linearly degenerated due to the specific choice of $V_I$. Regular solutions of system (1) comply with the following balance equation:

$$
\partial_t \Big( \sum_k m_k(\frac{u_k^2}{2} + f_k(\rho_k)) \Big) + \partial_x \Big( \sum_k (\alpha_k u_k (\rho_k \frac{u_k^2}{2} + \rho_k f_k(\rho_k) + P_k)) \Big) = -\Lambda |u_r|u_r^2 - \Theta P_r^2
$$

where the function $f_k(\rho_k)$ is such that: $f_k'(\rho_k) = \dfrac{\mathscr{P}_k(\rho_k)}{\rho_k^2}$.

We now assume that scalar functions $\Theta(W)$, $\Lambda(W)$ behave as: $\Theta(W) = \dfrac{\theta(W)}{\varepsilon^2}$ and: $\Lambda(W) = \dfrac{\lambda(W)}{\varepsilon^2}$ with respect to some small parameter $\varepsilon$. Hence, by using a Chapman-Enskog expansion, we know (see [1] and [5]) that system (1) may be rewritten in the asymptotic regime in the following modified form:

$$
\begin{cases}
\dfrac{\partial \rho}{\partial t} + \dfrac{\partial}{\partial x}(\rho u) & = 0 \\[2mm]
\dfrac{\partial \rho Y}{\partial t} + \dfrac{\partial}{\partial x}(\rho Y u + \rho Y(1-Y)u_r) & = 0 \\[2mm]
\dfrac{\partial \rho u}{\partial t} + \dfrac{\partial}{\partial x}(\rho u^2 + P + \rho Y(1-Y)u_r^2) & = 0 \\[2mm]
|u_r|u_r = \varepsilon^2 \dfrac{\rho Y(1-Y)}{\lambda}(1/\rho_2 - 1/\rho_1)\dfrac{\partial P}{\partial x} & \\[2mm]
P_r = 0 &
\end{cases}
$$

while neglecting $\mathscr{O}(\varepsilon^2)$ contributions -except for $u_r^2$ terms-, and noting $\rho$, $\rho Y$, $\rho u$, $P$ the total mass, the mass of species 2, the total momentum and the mean pressure. Actually, the latter correspond to: $(1-\alpha_2)\rho_1 + \alpha_2\rho_2$, $\alpha_2\rho_2$, $(1-\alpha_2)\rho_1 u_1 + \alpha_2\rho_2 u_2$, and $(1-\alpha_2)P_1 + \alpha_2 P_2$ respectively. In the asymptotic model, the new equation of state is obtained by setting $P_r = 0$ and thus solving: $\mathscr{P}_1(\dfrac{\rho(1-Y)}{1-\alpha_2}) - \mathscr{P}_2(\dfrac{\rho Y}{\alpha_2}) = 0$ with respect to $\alpha_2$, at any point $(x, t)$, which eventually provides $P = \mathscr{P}_1(\dfrac{\rho(1-Y)}{1-\alpha_2})$.

## 3 Numerical Scheme

We present below a semi-implicit scheme in order to compute approximations of the solutions of system (1), in such a way that no constraint would arise in the choice of relaxation time scales, i.e. with $\varepsilon$. Another objective is to have discrete pressure contributions such that the relative velocity would agree with the asymptotic situation. The scheme basically relies on the single-phase algorithm quite recently introduced in [6]. Thus a classic staggered mesh arrangement is used. Pressures, mass fractions and total mass are evaluated at the centre of Finite Volume cells, while velocities are estimated at cell boundaries. Convective contributions are accounted for explicitly; meanwhile, the discrete pressure terms are implicit. Thus, setting $\delta w = (w^{n+1} - w^n)/\Delta t^n$ whatever $w$ is, the time scheme is:

$$\delta m_k + \left[\dfrac{\partial}{\partial x}(m_k u_k)\right]^n = 0$$

$$\delta \alpha_2 + \left[V_I \dfrac{\partial \alpha_2}{\partial x}\right]^n = \left[\Theta(u)(P_2 - P_1)\right]^{n+1} \quad \text{setting:} \quad P_k^{n+1} = \mathscr{P}_k(\dfrac{m_k^{n+1}}{\alpha_k^{n+1}})$$

$$\delta(m_k u_k) + \left[\frac{\partial}{\partial x}(m_k u_k^2)\right]^n + \left[\frac{\partial}{\partial x}(\alpha_k P_k)\right]^{n+1} - \left[P_I \frac{\partial \alpha_k}{\partial x}\right]^{n+1} = \left[(-1)^{k+1} \Lambda |u_r| u_r\right]^{n+1}$$

We define, for $m = n, n+1$ the mean values $(\phi)_{i+\frac{1}{2}}^m$ around cell interfaces $x_{i+\frac{1}{2}}$ as:

$$(\phi)_{i+\frac{1}{2}}^m = 1/2(\phi)_i^m + 1/2(\phi)_{i+1}^m, \qquad for: \phi = m_k, P_I$$

**Step 1:** The mass fractions are first advanced in time using the cell scheme:

$$h\delta\left((m_k)_i\right) + (F_k)_{i+\frac{1}{2}}^n - (F_k)_{i-\frac{1}{2}}^n = 0$$

where $h$ is the length of each primal or dual cell, and:

$$(F_k)_{i+\frac{1}{2}}^n = (u_k)_{i+\frac{1}{2}}^n \begin{cases} (m_k)_i^n & if \ (u_k)_{i+\frac{1}{2}}^n > 0 \\ (m_k)_{i+1}^n & otherwise \end{cases}$$

**Step 2:** Once the latter have been computed, volume fractions $(\alpha_2)_i^{n+1}$ are then obtained using:

$$h\delta\left((\alpha_2)_i\right) + (V_I)_{i-\frac{1}{2}}^n \left((\alpha_2)_i^n - H_{i-\frac{1}{2}}^n\right) + (V_I)_{i+\frac{1}{2}}^n \left(H_{i+\frac{1}{2}}^n - (\alpha_2)_i^n\right) = h\frac{\theta_i^{n+1}}{\varepsilon^2}(P_r)_i^{n+1}$$

with $(P_r)_i^{n+1} = (P_2)_i^{n+1} - (P_1)_i^{n+1}$ and where the upwind flux $H_{i+\frac{1}{2}}$ and pressures $(P_k)_i^{n+1}$ are defined by:

$$(P_k)_i^{n+1} = \mathscr{P}_k(\frac{(m_k)_i^{n+1}}{(\alpha_k)_i^{n+1}}) \quad and: \quad H_{i+\frac{1}{2}}^n = \begin{cases} (\alpha_2)_i^n & if \ (V_I)_{i+\frac{1}{2}}^n > 0 \\ (\alpha_2)_{i+1}^n & otherwise \end{cases}$$

This second step requires solving a non-linear equation $g(y) = 0$ within each cell, with respect to $y = (\alpha_2)_i^{n+1}$; new values $(\alpha_1)_i^{n+1} = 1 - (\alpha_2)_i^{n+1}$ can then be deduced. When $\theta(W) = \theta_0$ or when $\theta(W) = \theta_0 \alpha_1 \alpha_2$, where $\theta_0$ is a positive constant, the function $g(y)$ is monotone in $]0, 1[$, and admits a unique solution $y_{sol} \in ]0, 1[$ provided that equations of state satisfy natural conditions:

$$\lim_{y \to +\infty} \mathscr{P}_k(y) = +\infty \qquad \lim_{y \to +0^+} \mathscr{P}_k(y) = a \in \mathbb{R}$$

**Step 3:** Approximate values of velocities $(u_k)_{i+\frac{1}{2}}$ are evaluated on the staggered mesh in the third step. Once again, discrete convective fluxes $(G_k)_i^n$ on the boundary of the staggered mesh are calculated using the upwind scheme:

$$(G_k)_i^n = (F_k)_i^n \widetilde{(u_k)_i}^n$$

noting: $(F_k)_i^n = 1/2(F_k)_{i-\frac{1}{2}}^n + 1/2(F_k)_{i+\frac{1}{2}}^n$ and

$$\widetilde{(u_k)_i}^n = \begin{cases} (u_k)_{i-\frac{1}{2}}^n & if \quad (F_k)_i^n > 0 \\ (u_k)_{i+\frac{1}{2}}^n & otherwise \end{cases}$$

$$h\delta\left((m_k)_{i+\frac{1}{2}}(u_k)_{i+\frac{1}{2}}\right) + (G_k)_{i+1}^n - (G_k)_i^n + \left[(\alpha_k)_{i+1}^{n+1}(P_k)_{i+1}^{n+1} - (\alpha_k)_i^{n+1}(P_k)_i^{n+1}\right]$$

$$-(P_I)_{i+\frac{1}{2}}^{n+1}\left[(\alpha_k)_{i+1}^{n+1} - (\alpha_k)_i^{n+1}\right] = (-1)^{k+1}h\frac{\lambda_{i+1/2}^{n+1}}{\varepsilon^2}(u_r)_{i+\frac{1}{2}}^{n+1}\left|(u_r)_{i+\frac{1}{2}}^{n+1}\right|$$

This third step only requires computing the root of a second-order polynomial. As soon as $(m_k)_{i+\frac{1}{2}}^{n+1} > 0$, the existence and uniqueness of a real root is obtained and the solution is known explicitly. The preservation of the total mass is also guaranteed on the staggered mesh.

**Property** *The mass fractions $m_k$ remain positive if the following CFL-like condition holds:*

$$\Delta t \left( \max\left((u_k)_{i+\frac{1}{2}}^n; 0\right) - \min\left((u_k)_{i-\frac{1}{2}}^n; 0\right) \right) < h$$

*Volume fractions $\alpha_k$ remain in $]0; 1[$ by construction.*

We now have the following result, for constant $\theta(W) = \theta_0$ or $\theta(W) = \theta_0\alpha_1\alpha_2$:

**Theorem 1** *We consider the expansion: $\phi(x, t) = \phi_0(x, t) + \varepsilon\phi_1(x, t) + \mathcal{O}(\varepsilon^2)$. We assume bounded initial conditions such that $\alpha_k(x, t = 0) \in ]0; 1[$ and $m_k(x, t = 0) > 0$. We also assume that boundary conditions do not depend on $\varepsilon$.*

- *Then, for $n \geq 0$ the scheme associated with steps $(1 - 3)$ admits a limit when $\varepsilon \to 0$: if $W^n$, a discrete solution obtained by the scheme at time $t^n$—on all cells—has a bounded limit when $\varepsilon$ tends to 0, then $W^{n+1}$ has a limit, bounded too, when $\varepsilon$ tends to 0.*
- *For $n \geq 1$, and for all cells indexed by i, we have:*

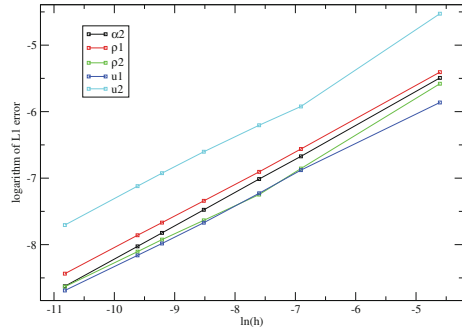$$((\alpha_k)_i^n)_0 \in ]0; 1[ \quad and: \quad ((m_k)_i^n)_0 > 0$$

*provided that the following CFL condition is ensured:*

$$\Delta t^n \left( \left[\max\left((u_k)_{i+\frac{1}{2}}^n; 0\right)\right]_0 - \left[\min\left((u_k)_{i-\frac{1}{2}}^n; 0\right)\right]_0 \right) < h.$$

- *We also have the asymptotic behaviour:*

$$((P_r)_i^n)_0 = 0 \quad and: \quad ((P_r)_i^n)_1 = 0.$$

**Fig. 1** Logarithm of the $L^1$ norm of the error $||w - w_h||_{L^1}$ as a function of $ln(h)$ for the 1D Riemann problem, for $w = \alpha_2, \rho_1, \rho_2, u_1, u_2$



*Moreover, the discrete relative velocity satisfies:*

$$\left((u_r)^n_{i+\frac{1}{2}}\right)_0 = 0$$

$$h\left((u_r)^{n+1}_{i+\frac{1}{2}}\right)_1 \left|\left((u_r)^{n+1}_{i+\frac{1}{2}}\right)_1\right| = \lim_{\varepsilon \to 0} \left(\left[(1/\rho_2 - 1/\rho_1)\frac{\partial P}{\partial x}\right]^{n+1}_{\tau, \ i+\frac{1}{2}} + R^n_{i+\frac{1}{2}}\right)$$

*where the estimates for the mean pressure gradient and the residual are:*

$$\left[(\tau_2 - \tau_1)\frac{\partial P}{\partial x}\right]^{n+1}_{\tau, \ i+\frac{1}{2}} = [\tau_{2,h} - \tau_{1,h}]\left[\sum_{k=1,2}(\alpha_k)^{n+1}_{i+1}(P_k)^{n+1}_{i+1} - \sum_{k=1,2}(\alpha_k)^{n+1}_i(P_k)^{n+1}_i\right]$$

*noting:* $\tau_k = 1/\rho_k$ *and* $\tau_{k,h} = ((\alpha_k)^{n+1}_{i+\frac{1}{2}})/((m_k)^{n+1}_{i+\frac{1}{2}})$

$$R^n_{i+\frac{1}{2}} = \sum_k \frac{(-1)^k}{(m_k)^{n+1}_{i+\frac{1}{2}}}\left[\left[(F_k)^n_i - (F_k)^n_{i+1}\right](u_k)^n_{i+\frac{1}{2}} + \left[(F_k)^n_{i+1}\widetilde{(u_k)}^n_{i+1} - (F_k)^n_i\widetilde{(u_k)}^n_i\right]\right]$$

**Idea of the proof:** We first show all these properties—except for the first-order's estimation on the relative velocity $(u_r)_1$—making a proof by mathematical induction on the number of time iterations. We consider each discretization successively—step 1, then step 2, and finally step 3—for all cells, and we deduce for each step that discrete unknowns have a limit, which is bounded, when $\varepsilon$ tends to 0, and meanwhile we get associated properties (positivity of partial masses—step 1—, positivity of statistical fractions and estimate for the relative pressure—step 2—, first estimate for the relative velocity $(u_r)_0$—step 3—) within each step. Eventually, we obtain the second estimate for the relative velocity $(u_r)_1$ using previous results.

**Fig. 2** $\alpha_2$, $\rho_k$ and $u_k$ for $\varepsilon = 10^{-4}$ (*red*), 1 (*green*), 2 (*pink*) and without source terms (*blue*)

## 4 Numerical Results

We first consider the system without source terms in order to compare the approximations obtained with the present scheme to the exact solution of a Riemann problem during mesh refinement. The source terms are then considered. The aim of this second test is to confirm the asymptotic behaviour of the solution.

In both cases the equations of state are given by: $\mathscr{P}_k(\rho_k) = \rho_k^{\gamma_k}$, with $\gamma_k > 1$. The time step complies with a classic CFL condition, setting $CFL = 0.5$ and the final time is $T = 0.1$ in all tests.

**Test 1: A Riemann problem for the convective part**
We set $\gamma_1 = 2$ and $\gamma_2 = 3$, $V_I = u_1$ and $P_I = P_2$. Initial conditions are:

$$(\alpha_2, \rho_1, \rho_2, u_1, u_2)_L = (0.7, 1, 0.8, 0.3, 0.4)$$

$$(\alpha_2, \rho_1, \rho_2, u_1, u_2)_R = (0.3, 0.88998555539, 0.5, 0.16183014405, 0.35732339488)$$

The solution of this Riemann problem contains a $(u_1 - c_1)$-shock wave, followed subsequently by a $u_1$-contact discontinuity and finally by a $(u_2 + c_2)$-shock wave. The finer mesh contains 50,000 cells, whereas the coarser mesh contains 100 cells. We observe on Fig. 1 a $h^{1/2}$ asymptotic rate of convergence as expected.

**Fig. 3** $ln(||u_r||_{L^1})$ (*left*) and $ln(||P_r||_{L^1})$ (*right*) as function of $ln(\varepsilon)$

**Test 2: Asymptotic behaviour**

In this test case, we set $\gamma_1 = 3$, $\gamma_2 = 1.5$, $V_I = u_2$ and $P_I = P_1$. We also set source terms coefficients: $\lambda = \theta = 1$. Initial conditions for the second test case are: $(\alpha_2, \rho_1, \rho_2, u_1, u_2)_L = (0.9, 0.8, 1, 0, 0)$ and $(\alpha_2, \rho_1, \rho_2, u_1, u_2)_R = (0.4, 1.2, 0.2, 0, 0)$.

A first order (respectively second-order) rate of convergence is retrieved for the relative velocity (respectively for the relative pressure). A mesh with 200 cells (Fig. 3) has been used for this test, although results on convergence rates are not mesh sensitive. The approximations obtained for several values of $\varepsilon$ with a mesh including 1,000 cells and the approximation obtained when source terms are not considered with 15,000 cells can be observed on Fig. 2.

# References

1. Ambroso, A., Chalons, C., Coquel, F., Galié, T., Godlewski, E., Raviart, P.A., Seguin, N.: The drift-flux asymptotic limit of barotropic two-phase two-pressure models. Commun. Math. Sci. **6**(2), 521–529 (2008)
2. Baer, M.R., Nunziato, J.W.: A two-phase mixture theory for the deflagration-to-detonation transition (ddt) in reactive granular materials. Int. J. Multiph. Flow **12**(6), 861–889 (1986)
3. Coquel, F., Hérard, J.M., Saleh, K., Seguin, N.: A robust entropy-satisfying finite volume scheme for the isentropic baer-nunziato model. ESAIM: Mathematical Modeling and Numerical Analysis (2013)
4. Cordier, F., Degond, P., Kumbaro, A.: An asymptotic preserving all-speed scheme for the euler and navier-stokes equations. J. Comput. Phys. **231**(17), 5685–5704 (2012)
5. Duval, F., Guillard, H.: A darcy law for the drift velocity in a two-phase flow model. J. Comput. Phys. **224**(1), 288–313 (2007)

6. Herbin, R., Latché, J.C., Nguyen, T.T.: Explicit staggered schemes for the compressible euler equations. ESAIM: Proceedings, vol. 40, pp. 83–102 (2013)
7. Thanh, M., Kröner, D., Nam, N.: Numerical approximation for a baer-nunziato model of two-phase flows. Appl. Numer. Math. **61**, 702–721 (2011)

# Numerical Simulations of a Fluid-Particle Coupling

**Nina Aguillon**

**Abstract** We present numerical simulations of a model of coupling between a inviscid compressible fluid and a pointwise particle. The particle is seen as a moving interface, through which interface conditions are prescribed. Key points are to impose those conditions at the numerical level, and to deal with the coupling between an ordinary and a partial differential equations.

## 1 The Model

We consider the following coupling, introduced in [2], between a pointwise particle of position $h$, and a fluid governed by the isothermal Euler equations, having density $\rho(t, x)$ and velocity $u(t, x)$ at time $t$ and point $x$:

$$
\begin{cases}
\partial_t \rho + \partial_x (\rho u) = 0, \\
\partial_t (\rho u) + \partial_x \left( \rho u^2 + c^2 \rho \right) = -D(\rho, \rho(u - h'(t)))\delta_{h(t)}(x), \\
m h''(t) = D(\rho(t, h(t)), \rho(u(t, h(t)) - h'(t))).
\end{cases}
\tag{1}
$$

Here, $c$ is the speed of sound. The fluid and the particle interact with each other through the drag force $D$, which applies only at the point where the particle is located. If $D$ has the same sign as $u - h'$, it formally tends to bring the velocities of the fluid and the particle closer to each other. Indeed the third line of (1), which is nothing else than Newton's law applied to the particle, yields that the particle accelerates if its velocity is smaller than the fluid's velocity. This system is a generalization of the coupling between a particle and an inviscid fluid introduced and studied in [4, 5, 15] (see references therein). In [6], one can find another model of coupling between

N. Aguillon (✉)
Université Paris Sud, 91405 Orsay, France
e-mail: nina.aguillon@math.u-psud.fr

a pointwise particle and a compressible inviscid fluid. The model is different, and local in time existence of solution for small subsonic data is proved, with tools develop in [8]. Let us start with two remarks about System (1). We denote by $H$ the Heaviside function. With the new unknown $w := H(x - h(t))$, which verifies $\partial_t w - h'(t)\partial_x w = 0$, we can write (1) as a non-conservative system of conservation laws. It is not strictly hyperbolic: its Jacobian matrix has eigenvalues $u - c$, $u + c$ and $h'$, and is not diagonalizable when $h' = u \pm c$. Moreover, as shocks appear in finite time in the solutions of the Euler equations, the right hand-sides of (1) are not well defined. However, it is possible to reformulate the System (1) as an interface problem. In the sequel we denote by $(\rho_-, u_-)$ and $(\rho_+, u_+)$ the traces of the fluid on the left and on the right of the particle: e.g. $\rho_-(t) = \lim_{x \to h(t)-} \rho(t, x)$. Interface conditions are imposed by saying that the traces must belong to a certain set. In the spirit of [3], we call that set the *germ* and we denote it by $\mathscr{G}_D(h')$.

**Definition 1** We denote by $F_\alpha$ an antiderivative of the function $\rho \mapsto \frac{\alpha^2/\rho + c^2\rho}{|D(\rho, \alpha)|}$. The germ $\mathscr{G}_D(h')$ is the set of $((\rho_-, u_-), (\rho_+, u_+))$ in $(\mathbb{R}_+ \times \mathbb{R})^2$ such that

1. $\rho_-(u_- - h') = \rho_+(u_+ - h')$. We denote by $\alpha$ this quantity;
2. Either $F_\alpha(\rho_-) - F_\alpha(\rho_+) = \text{sign}(\alpha)$, or there exists $\theta \in [0, 1]$ and $\rho_0 \leq \frac{|\alpha|}{c}$ such that

   a. $\rho_- \leq \frac{\alpha}{c} \leq \rho_+$ and $F_\alpha(\rho_-) - F_\alpha(\rho_0) = \theta$ and $F_\alpha(\frac{c^2}{\alpha^2\rho_0}) - F_\alpha(\rho_+) = (1 - \theta)$;

   b. $\rho_+ \leq \frac{-\alpha}{c} \leq \rho_-$ and $F_\alpha(\rho_+) - F_\alpha(\rho_0) = \theta$ and $F_\alpha(\frac{c^2}{\alpha^2\rho_0}) - F_\alpha(\rho_-) = (1 - \theta)$;

3. If $u_- > h'$ and $u_- - h' \leq c$, then $u_+ - h' \leq c$;
4. If $u_+ < h'$ and $u_+ - h' \geq -c$, then $u_- - h' \geq -c$.

This relation are obtained thanks to a thickening of the particle, where the Heaviside function $H$ is replaced by one of its regularization $H_\varepsilon$. It appears that the densities and velocities at the entry and at the exit of the particle are always linked by the relations of Definition 1, whatever the size $\varepsilon$ of the particle is, and which regularization is chosen (see [2] for more details). The "Riemann invariants" of the wave associated to eigenvalue $h'$ of System (1) are $\alpha = \rho(u - h')$ and $F_\alpha - H_\varepsilon$.

**Definition 2** A triplet $(\rho, u, h) \in L^\infty(\mathbb{R}_+ \times \mathbb{R}) \times L^\infty(\mathbb{R}_+ \times \mathbb{R}) \times W_{loc}^{2,\infty}(\mathbb{R}_+)$ is called an entropy solution of the problem (1) if:

1. The pair of functions $(\rho, u)$ is a weak entropy solution of the isothermal Euler equations on the sets $\{(t, x) \in \mathbb{R}_+^* \times \mathbb{R} : x > h(t)\}$ and $\{(t, x) \in \mathbb{R}_+^* \times \mathbb{R} : x < h(t)\}$;
2. For almost every $t > 0$, the traces around the particle exist and belong to the germ at speed $h'(t)$: $((\rho_-(t), u_-(t)), (\rho_+(t), u_+(t))) \in \mathscr{G}_D(h'(t))$;
3. For almost every $t > 0$, the particle is driven by the ODE:

$$mh''(t) = c^2(\rho_-(t) - \rho_+(t))\left(1 - \frac{(u_-(t) - h'(t))(u_+(t) - h'(t))}{c^2}\right). \quad (2)$$

The following Proposition, which is proven by simple computations, justifies the first point of Definition 1 and the reformulation (2) of the ODE.

**Proposition 1** *A solution $(\rho, u)$ of the Euler equation on the sets $\{x < h\}$ and $\{x > h\}$, with total bounded variations, conserves the total mass $\int_{\mathbb{R}} \rho dx$ if and only if for almost every time,*

$$\rho_-(u_- - h') = \rho_+(u_+ - h').$$

*In that case, it conserves the total impulsion $\int_{\mathbb{R}} \rho u dx + mh'$ if and only if for almost every time, the particle is driven by Eq. (2).*

*Proof* The proof consists in cutting integrals on $\mathbb{R}$ as integrals on $\{x < h\}$ and $\{x > h\}$. For the total impulsion, we obtain that $h$ must verify

$$mh''(t) = h'(\rho_+ u_+ - \rho_- u_-) + (\rho_- u_-^2 + c^2 \rho_-) - (\rho_+ u_+^2 + c^2 \rho_+).$$

When the mass is conserved, we express $u_\pm$ in terms of $\rho_\pm$ and $\alpha := \rho_\pm(u_\pm - h')$ to obtain (2).

The main result of [2] exhibits some conditions under which the Riemann problem for a motionless particle is well-posed.

**Theorem 1** *Consider a particle having a constant velocity equal to some real v. If the drag force D has the same sign as $\alpha := \rho(u - v)$, is an increasing function of $\alpha$ and if $|D|$ is a decreasing function of $\rho$, then for all $((\rho_L, u_L), (\rho_R, u_R))$ in $(\mathbb{R}_+ \times \mathbb{R})^2$, there exists a unique self similar solution to the Riemann problem*

$$\begin{cases} \partial_t \rho + \partial_x (\rho u) = 0, \\ \partial_t (\rho u) + \partial_x \left(\rho u^2 + c^2 \rho\right) = -D(\rho, \alpha)\delta_{vt}(x), \\ (\rho(0, x), u(0, x)) = (\rho_L, u_L)\mathbf{1}_{x<0} + (\rho_R, u_R)\mathbf{1}_{x>0} \end{cases} \quad (3)$$

The main difficulty is the non-hyperbolicity of the system. The Riemann problem has a more complicated structure than in the strictly hyperbolic case, and in particular uniqueness can be lost (see for example [11, 14]). This is the case for the drag force $D(\rho, \alpha) = \rho$ illustrated below, which violates the hypothesis of Theorem 1. Remark that this source term is similar to the source term in the shallow water equations with discontinuous topography. The Riemann problem (3) with $\rho_L = 0.7$, $\rho_R = 5$, $q_L = 5$, $q_R = 9$, $c = 2$ and $\lambda = 1.5$ admits three solutions, depicted on the right of Fig. 1. As in [7], this coexistence of solutions persists at the numerical level. We can see on the left of Fig. 1 two solutions selected by the Godunov scheme when replacing the Dirac measure by

$$x \mapsto \exp((x/\eta - \xi)^2)/(\eta\sqrt{2}),$$

with $\eta = 0.005$ and $\xi = -0.5$ or $\xi = 0.5$. We used a splitting between the fluid part and the regularized source term. The subsonic and supersonic solutions are obtained

**Fig. 1** *Left* solutions at time $T = 0.15$ given by the Godunov scheme for different regularizations of the Dirac measure. *Right* the three solutions of the Riemann problem

for large range of parameter $\xi$, with a very quick transition between the two passing through the mixed solution.

## 2 Finite Volume Schemes for the Coupled System

In the sequel, we adopt classical notation for finite volume schemes. We denote by $q = \rho u$ the momentum of the fluid. In particular, $U_j^n = (\rho_j^n, q_j^n)$ is an approximation of the solution at the $n$th iteration in time and in the $j$th cell, and $g$ is the numerical flux. Consider the case where the particle has a fixed constant velocity $v$, and denote by $j_0^n$ the cell where the particle lies at the $n$th iteration in time. The three points scheme

$$\begin{cases} U_j^{n+1/2} = U_j^n - \frac{\Delta t}{\Delta x}(g(U_j^n, U_{j+1}^n) - g(U_{j-1}^n, U_j^n)), \\ U_{j_0^n}^{n+1} = U_{j_0^n}^{n+1/2} - \frac{\Delta t}{\Delta x}\begin{pmatrix} 0 \\ D(\rho_{j_0^n}^{n+1/2}, \rho_{j_0}^{n+1/2}(u_{j_0^n}^{n+1/2} - v)) \end{pmatrix}, \end{cases}$$

corresponds to a splitting scheme between the evolution of the fluid (first line) and the influence of the particle (second line). This scheme does not converge toward the correct solution, even in the simplest case where $D(\rho, \rho(u - h')) = \lambda\rho(u - h')$ (which fulfills the hypothesis of Theorem 1) and the initial data belongs to the germ. It can be seen on Fig. 2. This failure to capture a small scale phenomenon recalls the difficulties encountered when approximating non-classical shocks (see for example [13]) or non-conservative systems (see for example [9, 16]). It illustrates that the reformulation as an interface problem of system (1) is necessary.

### 2.1 Schemes for a Motionless Particle

When the particle is motionless, we can easily implement schemes based on the exact resolution of the Riemann problem, which is constructed in the proof of Theorem 1.
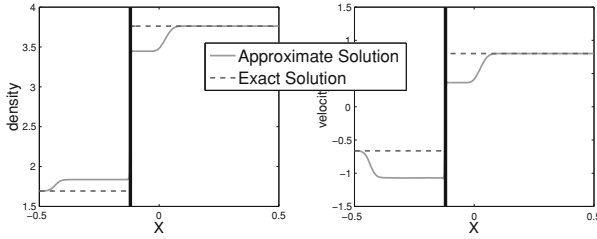
**Fig. 2** Solution at time $T = 0.04$ given by the fluid-particle splitting

Since the particle is not moving, the particle is a fixed interface that we place between cells numbered 0 and 1. We use a ghost-fluid approach (see [1, 10]) to write the scheme

$$\begin{cases} U_j^{n+1} = U_j^n - \frac{\Delta t}{\Delta x}(g(U_j^n, U_{j+1}^n) - g(U_{j-1}^n, U_j^n)) & \text{for } j \notin \{0, 1\}, \\ U_0^{n+1} = U_0^n - \frac{\Delta t}{\Delta x}(g(U_0^n, U_{\text{part},-}^n) - g(U_{-1}^n, U_0^n)) & \text{for } j = 0, \\ U_1^{n+1} = U_1^n - \frac{\Delta t}{\Delta x}(g(U_1^n, U_2^n) - g(U_{\text{part},+}^n, U_1^n)) & \text{for } j = 1. \end{cases} \quad (4)$$

Here, $U_{\text{part},-}^n = (\rho_{\text{part},-}^n, q_{\text{part},-}^n)$ and $U_{\text{part},+}^n = (\rho_{\text{part},+}^n, q_{\text{part},+}^n)$ are the values of the density and the momentum of the fluid on lines $x = 0^-$ and $x = 0^+$ of the unique self similar solution to (3), with

$$\rho_L = \rho_0^n, \ u_L = \frac{q_0^n}{\rho_0^n}, \ \rho_R = \rho_1^n \text{ and } u_R = \frac{q_1^n}{\rho_1^n}.$$

Remark that when $g$ is the Godunov flux, $U_0^{n+1}$ and $U_1^{n+1}$ are the averages of the exact solution *with particle* given by Theorem 1. In other words, it is the original Godunov scheme for the fluid/particle coupling. If we start with a Riemann problem belonging to $\mathscr{G}_D(0)$, i.e. verifying the relations of Defintion 1, we obtain for all $n$,

$$U_{\text{part},-}^n = U_L \text{ and } U_{\text{part},+}^n = U_R.$$

Adopting the vocabulary of [12], it follows that the scheme (4) is well balanced with respect to the whole germ $\mathscr{G}_D(0)$. We used this scheme to simulate a clogged organ pipe. The pipe is initially filled with a fluid at rest having density 5 kg/m, and we take c = 1 m/s. At time $t > 0$, a constant flow of 3 kg/s is imposed on the left entry of the pipe, while the gas exits freely on the right. The pipe is blocked in its middle by a porous particle that we model using the drag force $D(\rho, \rho u) = \rho u$. At time 0.041 s, the shock emitted by the left boundary condition hits the particle. The Riemann problem with the particle develops one shock on each side of the particle. Roughly speaking, most of the air is stuck in front of the particle, causing an elevation of its density and a decrease of its velocity. A small part of the fluid manages to pass through the particle, and has a large velocity on the exit of the particle by conservation

**Fig. 3** Two successive interactions between shocks and particle in a clogged organ pipe

of momentum through the particle. The shock on the left of the particle interacts with the left boundary at time 0.114 s, creating another shock that meets the particle at time 0.153 s. Asymptotically, the fluid has constant momentum all over the pipe, with high density and low velocity before the particle, and low density and high velocity afterwards. Shapes of the solution after the first two interactions of a shock with the particle are depicted on Fig. 3. This simulation illustrates the convergence of the ghost fluid scheme (4) on Riemann problem. We used the Godunov numerical flux but the results are similar with the Rusanov flux.

## 2.2 Dealing with a Moving Particle

We now focus on the case where the particle is free to move under the influence of the fluid. We saw in the introduction that it was necessary to treat the particle as an interface. Therefore, the particle must end up at an interface between two cells at the end of each time iteration. We could have used a mesh tracking the particle, but with in mind more complex applications (with numerous particles for example) we decided to use a fixed mesh and a Glimm's approach to replace the particle. At each time iteration, a real number $x_r$ is uniformly picked up in $[0, \Delta x]$. In the $j$th cell, the fluid is updated by the exact value of the solution at time $\Delta t$ and at point $x_r$ of the Euler equation with initial data

$$U^0(x) = U^n_{j-1}\mathbf{1}_{x<0} + U^n_j\mathbf{1}_{0<x<\Delta x} + U^n_{j+1}\mathbf{1}_{\Delta x<x}.$$

Under the classical CFL condition $\Delta t < \frac{\Delta x}{2\max_x |u(x)|+c}$, the solution consists in the juxtaposition of two Riemann problems. When $j$ corresponds to a neighbor cell of the particle, one of these Riemann problems takes the particle into account. The particle's position is updated in accordance to $x_r$. If the particle it at the interface $j^n_0 + 1/2$ at time $n$, and has speed $v^n$, then at time $n + 1$ we placed it:

**Fig. 4** *Left* velocity of the particle. Each discontinuity on its acceleration is caused by a shock hitting the marble. *Right* density of the fluid on the tube through time. We can see the shock with decreasing strength trapped between the marble and the bottom of the tube

1. at interface $j_0^n + 3/2$ if $v^n > 0$ and $x_r < v^n \Delta t$, in which case $j_0^{n+1} = j_0^n + 1$;
2. at interface $j_0^n - 1/2$ if $v^n < 0$ and $x_r > \Delta x + v^n \Delta t$, in which case $j_0^{n+1} = j_0^n - 1$;
3. at interface $j_0^n + 1/2$ otherwise, in which case $j_0^{n+1} = j_0^n$;

Eventually, we update the particle's velocity using (2) and the numerical traces. The following numerical simulation is inspired by [6]. A marble falls into a cylinder filled with a compressible inviscid gas, which is initially at rest and of density $1/225$ kg/m. Both the gas and the marble are subject to gravity and friction. The complete system writes

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0 \\ \partial_t(\rho u) + \partial_x(\rho u^2 + c^2 \rho) = -\lambda(u - h')\delta_{h(t)}(x) - \rho g - v_F(\rho, u) \\ mh''(t) = \lambda(u(t, h(t)) - h'(t)) - mg - mv_S(h'(t)) \end{cases}$$

where we take as in [6], $v_S(h') = 10^{-2}h'$, $v_F(\rho, u) = 10^{-8}\rho u|u|$, $c = 15$ m/s, $m = 0.004$ kg and $g = 9.81$ m/s$^2$. We took $\lambda = 5$ m$^2 \cdot$ kg/s. The first term of the ODE should be understood as in Eq. (2). At first, the marble compresses the gas beneath it, creating a shock, and its velocity decreases due to friction. At some time, the shock reflects on the closed bottom of the tube, and then hits the marble, creating a discontinuity in its acceleration and accelerating it. This can be seen on Fig. 4, on the left. When the shock interacts with the marble, it is somehow split in two: a part is reflected downward and a part passes through the marble and exits freely on the top end of the tube. Therefore the shock trapped between the marble and the bottom of the tube is of decreasing strength, as it can be seen on the plot of the density on the right of Fig. 4. The particle being very light, it is very sensible to the fluid's velocity, which is positive when the first shocks are moving upward. It causes the marble to climb back up for a while, then the gravity becomes predominant and the marble falls down again. The results are qualitatively the same as in [6]. However, they do not match perfectly, because the modeling is quite different. In particular in [6],

the friction between the fluid and the marble is taken into account via a source term $v_I = 5\left(h' - \frac{u_- + u_+}{2}\right)^2$, while it is modeled directly through the interface conditions of Definition 1 in the present work.

## 3 Perspectives

Let us start with some remarks on System (1), for which we proved in [2] existence and uniqueness to the Riemann problem when the particle is motionless, and give in this paper some qualitative properties and illustrative numerical simulations. Further theoretical study of System (1) seems difficult, as we have to deal with a system which is neither conservative nor hyperbolic. Even the extension of Theorem 1 to a freely moving particle is tricky, because the solution is not self-similar, and the traces around the particle constantly change. It is not difficult to extend the result to other pressure law, at least to $p(\rho) = a\rho^\gamma$, with $1 < \gamma \le 3$, $a > 0$ and where no vacuum appears. Therefore, this model could be used to model the influence of an obstacle into the shallow water equation. Similarly, the extension to the full Euler equations is interesting, and could take into account exchange of heat between the fluid and the particle.

## References

1. Abgrall, R., Karni, S.: Computations of compressible multifluids. J. Comput. Phys. **169**(2), 594–623 (2001)
2. Aguillon, N.: Riemann problem for a fluid-particule coupling (2014) (Submitted)
3. Andreianov, B., Karlsen, K.H., Risebro, N.H.: A study of $L^1$-dissipative solvers for scalar conservation laws with discontinuous flux. Arch. Ration. Mech. Anal. **201**, 27–86 (2011)
4. Andreianov, B., Lagoutière, F., Seguin, N., Takahashi, T.: Well-posedness for a one-dimensional fluid-particle interaction model. SIAM J. Math. Appl. **46**(2), 1030–1052 (2014)
5. Andreianov, B., Seguin, N.: Analysis of a Burgers equation with singular resonant source term and convergence of well-balanced schemes. Discrete Contin. Dyn. Syst. **32**, 1939–1964 (2012)
6. Borsche, R., Colombo, R.M., Garavello, M.: On the interactions between a solid body and a compressible inviscid fluid. To appear in Interfaces Free Bound (2014)
7. Boutin, B., Coquel, F., LeFloch, P.G.: Coupling techniques for nonlinear hyperbolic equations. III. The well-balanced approximation of thick interfaces. SIAM J. Numer. Anal. **51**, 1108–1133 (2013)
8. Colombo, R.M., Guerra, G.: On general balance laws with boundary. J. Diff. Equat. 1017–1043 (2010)
9. Dal Maso, G., Lefloch, P.G., Murat, F.: Definition and weak stability of nonconservative products. J. Math. Pures Appl. **74**(6), 483–548 (1995)
10. Fedkiw, R.P., Aslam, T., Merriman, B., Osher, S.: A non-oscillatory Eulerian approach to interfaces in multimaterial flows (the ghost fluid method). J. Comput. Phys. **152**, 457–492 (1999)
11. Goatin, P., LeFloch, P.G.: The Riemann problem for a class of resonant hyperbolic systems of balance laws. Ann. Inst. H. Poincaré Anal. Non Linéaire **21**, 881–902 (2004)

12. Greenberg, J.M., Leroux, A.Y.: A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. SIAM J. Numer. Anal. **33**(1), 1–16 (1996)
13. Hayes, B.T., Lefloch, P.G.: Nonclassical shocks and kinetic relations: strictly hyperbolic systems. SIAM J. Math. Anal. **31**(5), 941–991 (2000)
14. Isaacson, E., Temple, B.: Nonlinear resonance in systems of conservation laws. SIAM J. Appl. Math. **52**, 1260–1278 (1992)
15. Lagoutière, F., Seguin, N., Takahashi, T.: A simple 1D model of inviscid fluid-solid interaction. J. Diff. Equat. **245**, 3503–3544 (2008)
16. Parés, C.: Numerical methods for nonconservative hyperbolic systems: a theoretical framework. SIAM J. Numer. Anal. **44**, 300–321 (2006)

# A Simple Finite Volume Approach to Compute Flows in Variable Cross-Section Ducts

**Bruno Audebert, Jean-Marc Hérard, Xavier Martin and Ouardia Touazi**

**Abstract** In order to derive a simple one-dimensional approach that could handle fluid flows in smooth ducts as well as in ducts of discontinuous cross-section, we propose herein a Finite Volume approach that relies on an integral formulation of the multidimensional flow model. While focusing on Euler equations, we compare two-dimensional results with approximations obtained using the present approach, and also with the classical formulation for variable cross-sections using a well-balanced scheme. Numerical simulations confirm the ability of this integral method to provide approximations that compare well with 2D results. This method also enables to deal with all-even including vanishing-cross-section ducts. This approach may also be applied when considering other single-phase or multi-phase fluid flow models.

B. Audebert · J.-M. Hérard · X. Martin (✉) · O. Touazi
EDF R&D, Fluid Dynamics, Power Energy and Environment, 6 quai Watier,
78400 Chatou, France
e-mail: xavier-x.martin@edf.fr

J.-M. Hérard
e-mail: jean-marc.herard@edf.fr

B. Audebert
e-mail: bruno.audebert@edf.fr

O. Touazi
e-mail: ouardia.touazi@edf.fr

J.-M. Hérard and X. Martin
Laboratoire d'Analyse Topologie Probabilités, UMR CNRS 7353, 39 rue Joliot Curie,
13453 Marseille cedex 13, France

# 1 Introduction

Numerical tools devoted to the computation of single-phase or two-phase flows in ducts with variable cross sections are very useful in industry, because they enable to obtain a reasonable approximation of the true flow in unsteady situations, using standard computers. This is of particular interest for hydraulic circuits, as well as in some medical applications, however it requires the ability to deal with smooth or discontinuous cross sections. When neglecting viscous effects and external forces, the classical approach which is overwhelmingly retained consists of constructing numerical approximations of solutions of systems that take the form:

$$(S_1) \begin{cases} \dfrac{\partial S}{\partial t} = 0 \\[2mm] \dfrac{\partial \rho S}{\partial t} + \dfrac{\partial \rho u S}{\partial x} = 0 \\[2mm] \dfrac{\partial \rho u S}{\partial t} + \dfrac{\partial \rho u^2 S}{\partial x} + S \dfrac{\partial P}{\partial x} = 0 \\[2mm] \dfrac{\partial ES}{\partial t} + \dfrac{\partial u(E+P)S}{\partial x} = 0 \end{cases}$$

where $S(x)$ stands for the area of the cross section, and $\rho$, $u$, $P$, $E$ denote the density, velocity, pressure and total energy of the fluid. Several investigations of the problem that arises with discontinuous cross-sections have been published, among which we may cite [1, 6–8, 10], wherein authors focus either on the continuous or the discrete framework. Roughly speaking, most of the schemes that have emerged to cope with this problem rely on the well-balanced strategy [5]. The use of this strategy would even seem mandatory; otherwise approximate solutions can sometimes converge towards incorrect solutions (see [3, 6, 8]). Nonetheless, an inconvenience of this strategy is that it assumes that the Riemann invariants of the standing wave associated with $\lambda = 0$ are preserved, which of course makes sense for mass flux and total enthalpy flux, but is questionable in the case of the last Riemann invariant. This has been confirmed by numerical comparisons (see the work reported in [4] for instance), and it is actually quite a well-known problem, the classic treatment for which consists of the introduction of head losses using various empirical closure laws. This problem has motivated the present work, which aims at providing a somewhat different approach in order to eliminate the limitations and drawbacks of the classical approach. Another motivation will be discussed in the conclusion.

The current paper presents the main ideas and results of the work and is organised as follows: firstly, we present the modified one-dimensional approach; next we present a few numerical results, with a comparison with the two-dimensional approach, the classical approach $(S_1)$ and the modified one-dimensional formulation, using sufficiently fine and reliable meshes.

## 2 A Finite Volume Approach for One-Dimensional Flows

The one-dimensional formulation is obtained as follows. Starting with the three-dimensional governing equations, restricted here to the Euler framework, thus:

$$(S_2) \begin{cases} \dfrac{\partial \rho}{\partial t} + \nabla.(\rho \underline{u}) = 0 \\ \dfrac{\partial \rho \underline{u}}{\partial t} + \nabla.(\rho \underline{u} \otimes \underline{u}) + \nabla P = 0 \\ \dfrac{\partial E}{\partial t} + \nabla.\big(\underline{u}(E + P)\big) = 0 \end{cases}$$

where the total energy $E$ is $E = \rho((\underline{u})^2 + \varepsilon(P, \rho))/2$ and $\varepsilon(P, \rho)$ is the internal energy, we integrate over time—from time $t^n$ to $t^{n+1}$—and space using coarse control volumes as depicted on Fig. 1. At time $t = t_p$, we denote:

$$\Omega_i^\varphi \Phi_i^p = \int_{\Omega_i^\varphi} \Phi(\underline{x}, t_p) dv$$

for: $\Phi = \rho, \underline{Q}, E$ and also $\Omega_i^\varphi = S_i \times h_i$ the volume occupied by the fluid within the $i$-cell. Using previous definitions, and noting $\Gamma_i$ the boundary of control volume $\Omega_i$, straightforward calculations yield:

$$(S_3) \begin{cases} \Omega_i^\varphi \left(\rho_i^{n+1} - \rho_i^n\right) + \int_{[t^n, t^{n+1}]} \int_{\Gamma(i)} (\rho \underline{u}.\underline{n})(\underline{x}_\Gamma, t) d\Gamma dt = 0 \\ \Omega_i^\varphi \left(\underline{Q}_i^{n+1} - \underline{Q}_i^n\right) + \int_{[t^n, t^{n+1}]} \int_{\Gamma(i)} ((\rho \underline{u}.\underline{n})\underline{u} + P\underline{n})(\underline{x}_\Gamma, t) d\Gamma dt = 0 \ s \\ \Omega_i^\varphi \left(E_i^{n+1} - E_i^n\right) + \int_{[t^n, t^{n+1}]} \int_{\Gamma(i)} ((\rho \underline{u}.\underline{n})H)(\underline{x}_\Gamma, t) d\Gamma dt = 0 \end{cases}$$

where $\underline{Q} = \rho \underline{U}$ is the momentum and $H = (E + P)/\rho$ is the total enthalpy. Of course, viscous effects and gravity forces could also be included if required.

We may now introduce a simple explicit Finite Volume scheme FVCA (Finite-volumes for Variable Cross-section Applications) as follows:

$$(FVCA) \begin{cases} \Omega_i^\varphi \left(\rho_i^{n+1} - \rho_i^n\right) + \Delta t^n \sum_{j \in V(i)} (\rho \underline{u}.\underline{n})_{ij}^h \Gamma_{ij}^\varphi = 0 \\ \Omega_i^\varphi \left(\underline{Q}_i^{n+1} - \underline{Q}_i^n\right) + \Delta t^n \sum_{j \in V(i)} ((\rho \underline{u}.\underline{n})\underline{u} + P\underline{n})_{ij}^h \Gamma_{ij}^\varphi = 0 \\ \Omega_i^\varphi \left(E_i^{n+1} - E_i^n\right) + \Delta t^n \sum_{j \in V(i)} ((\rho \underline{u}.\underline{n})H)_{ij}^h \Gamma_{ij}^\varphi = 0 \end{cases}$$

where $(\psi)^h$ stands for some suitable flux scheme (exact or approximate Godunov scheme) associated with the continuous flux $\psi$, and setting $\Delta t^n = t^{n+1} - t^n$; $V(i)$ refers to the neighbouring cells of cell $i$ and to ghost "mirror" cells associated with the wall boundaries of cell $i$ (see Fig. 1).

**Fig. 1** Finite volume $\Omega_i$ with neighbouring cells, fluid interfaces and inner wall-boundaries



We now assume that the initial condition at time $t^n$ is such that the transverse velocity in the $y$-direction is null everywhere: $U_{y_i}^n = 0$. Using the exact Riemann solution for fluxes around all interfaces, and using the mirror technique for all wall boundaries, it may be easily checked that the scalar product of $\underline{e}_y$ with the discrete momentum equation in (*FVCA*) leads to: $(Q_{y_i}^{n+1} - Q_{y_i}^n) = 0$, and thus $Q_{y_i}^{n+1} = 0$ or $U_{y_i}^{n+1} = 0$. This simply means that the discrete flow remains 1D. We detail now mass and energy balance equations. These read:

$$\Omega_i^\varphi \left( \rho_i^{n+1} - \rho_i^n \right) + \Delta t^n \left( (\rho u_x)_{i+1/2}^h \Gamma_{i+1/2}^\varphi - (\rho u_x)_{i-1/2}^h \Gamma_{i-1/2}^\varphi \right) = 0 \quad (1)$$

and:

$$\Omega_i^\varphi \left( E_i^{n+1} - E_i^n \right) + \Delta t^n \left( (\rho H u_x)_{i+1/2}^h \Gamma_{i+1/2}^\varphi - (\rho H u_x)_{i-1/2}^h \Gamma_{i-1/2}^\varphi \right) = 0 \quad (2)$$

setting $\Gamma_{i+1/2}^\varphi = min(S_i, S_{i+1})$. Eventually, the discrete $x$-momentum balance for $Q_x = \rho u_x$ takes the final form:

$$\Omega_i^\varphi \left( Q_{x_i}^{n+1} - Q_{x_i}^n \right) + \Delta t^n \left( (\rho u_x^2 + P)_{i+1/2}^h \Gamma_{i+1/2}^\varphi - (\rho u_x^2 + P)_{i-1/2}^h \Gamma_{i-1/2}^\varphi \right)$$
$$+ \Delta t^n P_{i+\frac{1}{2},i}^* \left( S_i - \Gamma_{i+1/2}^\varphi \right) - \Delta t^n P_{i-\frac{1}{2},i}^* \left( S_i - \Gamma_{i-1/2}^\varphi \right) = 0 \quad (3)$$

where $P_{i\pm\frac{1}{2},i}^*$ is an estimation of the Riemann pressure on the wall boundaries $i \pm 1/2$.

Focusing for instance on perfect gas EOS, hence setting $P = (\gamma - 1)\rho\varepsilon(P, \rho)$, and using classical results (see [2] for example), we obtain when $S_i > S_{i+1}$:

- if $M_i = \frac{u_i^n}{c_i^n} < 0$, then: $P_{i+\frac{1}{2},i}^* = \begin{cases} P_i^n \left(1 + \frac{\gamma-1}{2}M_i\right)^{\frac{2\gamma}{\gamma-1}} & \text{if } 1 + \frac{\gamma-1}{2}M_i \geq 0 \\ \\ 0 & \text{otherwise} \end{cases}$

**Fig. 2** Experimental setup: 1D pipe with a sudden contraction and position of the initial membrane

- if $M_i = \frac{u_i^n}{c_i^n} > 0$, then: $P^*_{i+\frac{1}{2},i} = P_i^n \left( 1 + \gamma M_i \left( 1 + \frac{(\gamma+1)^2}{16} M_i^2 \right)^{1/2} + \frac{\gamma(\gamma+1)}{4} M_i^2 \right)$

The same technique is applied when $S_i < S_{i+1}$ in order to estimate $P^*_{i+\frac{1}{2},i+1}$.

On the whole, we can now compute mass, $x$-momentum and energy balance with the aid of (1–3), assuming that some standard explicit CFL condition holds for $\Delta t^n$. The counterpart of the latter expressions of $P^*_{i\pm\frac{1}{2},i}$ can be found for any EOS, using the mirror state and shock/rarefaction curves in GNL waves. Obviously, there are no intrinsic limits for cross-section values, even if $S_i = 0$. Depending on the choice of numerical fluxes at the fluid interfaces, CFL-like conditions must be introduced in order to guarantee positive discrete values of the density $\rho_i^n$. Further details can be found in [11].

## 3 Numerical Results

We present in this section a few results arising from a comparison of the three distinct approaches.

- A first approach simply consists of computing the complete set of equations ($S_2$) using the approximate Godunov scheme [2] on a fine enough two-dimensional mesh of about one million cells; the results will be called the reference solution;
- The second series of results were obtained with the classical well-balanced strategy applied to the set of one-dimensional equations ($S_1$), with the focus here on very fine meshes only; the well-balanced Rusanov scheme used in these computations is the one proposed in [8] and also used in [3] where the convergence towards the correct solution has been verified;
- The third series illustrates the numerical approximations obtained by computing the integral system (1–3) on fine one-dimensional meshes (called 1D+).

Actually, two slightly different ways of estimating the pressure on the wall boundaries will be applied to the set of formulas above, corresponding respectively to the exact Riemann solution and to the same approximation obtained by setting $M_i = 0$.

The experimental setup is the following: a one dimensional pipe contains a sudden cross-section contraction located at $x = 0.8$ (see Fig. 2). At the start of the simulation, a membrane at $x = 0.7$ separates two distinct states $(\rho_L, u_L, P_L) = (1, 0, 10^5)$ and $(\rho_R, u_R, P_R) = (0.125, 0, 10^4)$. Hence, at the beginning, a right-going shock wave followed by a contact discontinuity propagates, then "hits" the cross-section

**Fig. 3** Density profiles at $t = T_0$ in test case 1. *Dashed blue curve* 1D+ approach with 50,000 cells. *Dotted red* and *dashed green curves* 1D+ approach with 50,000 and 1,000 cells respe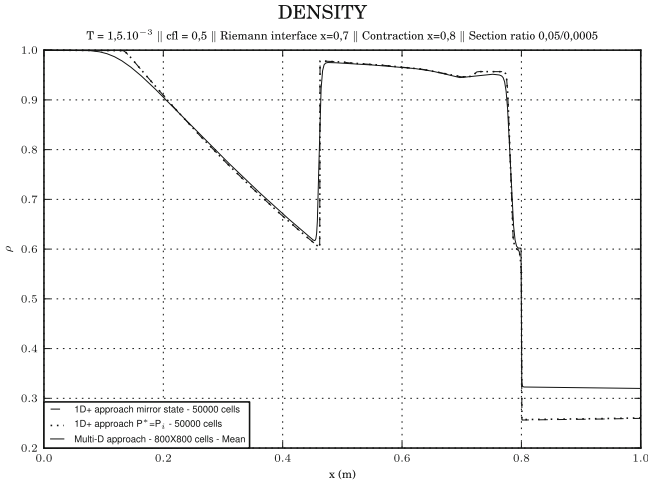ctively, assuming $M_i = 0$ in wall pressures. *Dotted-dashed magenta curve* well-valanced Rusanov scheme with 50,000 cells. *Black curve* $y-$averaging of 2D results

contraction; this results in a right-going transmitted wave and a left-going reflected wave. We have used a perfect gas EOS setting $\gamma = 7/5$. The fine one-dimensional meshes used for the classical and $1D+$ approaches contain 50,000 regular cells, and the CFL number has been set to 1/2. Two different cross-section ratios are considered, $S_l/S_r = 2$ (Fig. 3), and $S_l/S_r = 100$ (Figs. 4 and 5) in test cases 1 and 2 respectively.

**Test case 1:** $S_l/S_r = 2$: This corresponds to a rather classical situation arising in many practical simulations. We have plotted on Fig. 3 the density profiles at time $t = T_0 = 1.5 \times 10^{-3}$. As was expected in this particular case, the $1D+$ approximation where $M_i$ is set to 0 (dotted red for 50,000 cells and dashed green for one thousand cells) fits experimental "results" (in black) quite well, and performs better than the standard wall-pressure estimation (dashed blue with 50,000 cells). The former $1D+$ approach (setting $M_i = 0$ in the wall pressure formula) is also much more relevant than the classical approach (1) using the well-balanced Rusanov scheme ([8], magenta dashed dotted line in Fig. 3). Results of the $1D+$ approach are similar, whenever a coarse mesh (one thousand cells) or a fine mesh (50,000 cells) is used.

**Test case 2:** $S_l/S_r = 100$: Here, the well-balanced Rusanov scheme [8] fails to provide approximations on fine meshes, and a similar problem occurs when using the well-balanced approximate Godunov scheme [6]. Thus we were only able to compare results of the multi-dimensional approach to the results provided by the $1D+$ approach (see Figs. 4 and 5). Both estimations of $P^*_{i\pm\frac{1}{2},i}$ provide similar results, which again was expected, and the comparison with the multi-dimensional approach is even better in this case, which may be explained.

**Fig. 4** Density profiles at $t = T_0$ in test case 2. *Dashed curve* 1D+ approach. *Dotted curve* 1D+ approach assuming $M_i = 0$ in wall pressure estimations. *Black curve* $y$-averaging of two-dimensional results



**Fig. 5** Comparison of wall pressures in test case 2. *Dashed red curve* 1D+ approach. *Dotted* and *dashed green curves* 1D+ approach setting $M_i = 0$. *Dotted blue curve* multidimensional computation using $400^2$ cells. *Full black curve* multidimensional computation using $800^2$ cells

## 4 Conclusion and Further Work

The present $1D+$ approach is a very simple one relying on a straightforward integral formulation on particular Finite volumes, combined with an estimation of wall-pressure interactions. We have briefly presented a few of the results from among the

sixteen distinct situations that have been investigated up till now, where rarefaction or shock waves interact with eight contractions ($S_l/S_r = 10^{-2}/10^{-1}/0.5/0.9$ and $S_l/S_r = /(0.9)^{-1}/2/10/100$, see [11]). We would like to emphasize that:

- The present approach could be extended in order to take external forces, viscous contributions into account, without any loss of generality;
- The focus here has been on Euler equations but other (single phase or multiphase) fluid flow models could also be considered;
- A key point is that vanishing cross sections may occur in the duct; furthermore, it must be emphasized that numerical results depend continuously on the cross-section distribution. This can hardly be achieved with the classical approach, at least not when using well-balanced schemes that rely on approximate Godunov schemes. Moreover, even when the classical approach ($S_1$) is feasible, numerical results do not sufficiently match multi-dimensional results.

Another important point is that this method could be extended in order to improve the formulation that is currently applied in a particular three-dimensional porous framework widely used in the nuclear industry (see [9] for instance). We also plan to use the present results in order to improve the basic well-balanced strategy.

# References

1. Clain, S., Rochette, D.: First and second-order finite volume methods for the one-dimensional non-conservative euler system. J. Comput. Phys. **228**, 8214–8248 (2009)
2. Gallouët, T., Hérard, J.M., Seguin, N.: On the use of symetrizing variables for vacuums. Calcolo **40**(3), 163–194 (2003)
3. Girault, L., Hérard, J.M.: A two-fluid hyperbolic model in a porous medium. ESAIM: Math. Model. Numer. Anal. **44**(6), 1319–1348 (2010)
4. Girault, L., Hérard, J.M.: Multidimensional computations of a two-fluid hyperbolic model in a porous medium. Int. J. Finite Volumes **7**(1), 1–33 (2010)
5. Greenberg, J.M., Leroux, A.Y.: A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. SIAM J. Numer. Anal. **33**(1), 1–16 (1996)
6. Helluy, P., Hérard, J.M., Mathis, H.: A well-balanced approximate riemann solver for compressible flows in variable cross-section ducts. J. Comput. Appl. Math. **236**(7), 1976–1992 (2012)
7. Kröner, D., LeFloch, P., Thanh, M.D.: The minimum entropy principle for compressible fluid flows in a nozzle with discontinuous cross-section. ESAIM: Math. Model. Numer. Anal. **42**(3), 425–443 (2008)
8. Kröner, D., Thanh, M.D.: Numerical solutions to compressible flows in a nozzle with variable cross-section. SIAM J. Numer. Anal. **43**(2), 796–824 (2006)
9. Le Coq, G., Aubry, S., Cahouet, J., Lequesne, P., Nicolas, G., Pastorini, S.: The THYC computer code. A finite volume approach for 3 dimensional two-phase flows in tube bundles. Bulletin de la Direction des études et recherches-Electricité de France. Série A, nucléaire, hydraulique, thermique. In French **1**, 61–76 (1989)

10. LeFLoch, P., Thanh, M.D.: The riemann problem for fluid flows in a nozzle with discontinuous cross-section. Commun. Math. Sci. **1**, 763–797 (2003)
11. Martin, X.: Numerical modeling of flows in obstructed media. Ph.D. thesis (in preparation)

# A 1D Stabilized Finite Element Model for Non-hydrostatic Wave Breaking and Run-up

P. Bacigaluppi, M. Ricchiuto and P. Bonneton

**Abstract** We present a stabilized finite element model for wave propagation, breaking and run-up. Propagation is modelled by a form of the enhanced Boussinesq equations, while energy transformation in breaking regions is captured by reverting to the shallow water equations and allowing waves to locally converge into discontinuities. To discretize the system we propose a non-linear variant of the stabilized finite element method of (Ricchiuto and Filippini, *J.Comput.Phys.* 2014). To guarantee monotone shock capturing, a non-linear mass-lumping procedure is proposed which locally reverts the third order finite element scheme to the first order upwind scheme. We present different definitions of the breaking criterion, including a local implementation of the convective criterion of (Bjørkavåg and Kalisch, *Phys.Letters A* 2011), and discuss in some detail the implementation of the shock capturing technique. The robustness of the scheme and the behaviour of different breaking criteria is investigated on several cases with available experimental data.

## 1 Modelling Approach and Main Objectives

When arriving in the near shore region, waves are relatively long, with a ratio water-depth over wavelength $\sigma^2 \ll 1$. When approaching the shoreline the wave steepens and non-linear effects start to become dominating up to the moment in which the wave breaks ($\varepsilon = A/d \sim 1$), with important production of vorticity, and with potential
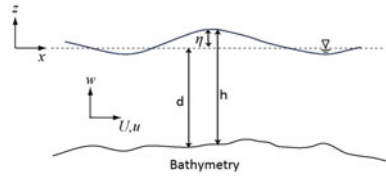
---

P. Bacigaluppi (✉) and M. Ricchiuto
Inria Bordeaux Sud-Ouest, 200 av. de la Vieille Tour, 33405 Talence Cedex, France
e-mail: paola.bacigaluppi@mail.polimi.it

M. Ricchiuto
e-mail: Mario.Ricchiuto@inria.fr

P. Bonneton
UMR CNRS EPOC, Bordeaux University, av. des Facultés, 33405 Talence, France
e-mail: p.bonneton@epoc.u-bordeaux1.fr

**Fig. 1** Depth averaged equations: notation



energy transformation and dissipation. The first phase of the process can be modelled by the a properly chosen set of non-hydrostatic equations, such as for example, the Boussinesq equations, or other type of dispersive models [2]. The treatment of wave breaking is more delicate. Several approaches exist, see [10] for a recent review. The extensive study of [4], indicate that the energy transformation in the breaking region can be modelled by the dissipation across nonlinear discontinuities of hyperbolic models such as the shallow water equations. This is confirmed by the numerical results of [5, 10, 13, 18, 19]. For this reason, we use a hybrid model reverting from the enhanced Boussinesq to the shallow water equations in properly defined breaking regions.

To discretize the equations, we start from the stabilized finite element approach of [16], which has a very interesting potential in terms of providing low dispersion errors and very high accuracy on unstructured adaptive meshes. Here, we propose a new nonlinear variant of the method. In our approach, the third order finite element scheme is reverted to the first upwind scheme across discontinuities via nonlinear mass-lumping procedure. The objective of this paper is to present the hybrid modelling approach, and in particular the definition of the breaking detection algorithm, and the discussion of the discontinuity capturing methodology, and in particular of the choice of the mass-lumping limiter. Concerning the first aspect, we consider the hybrid criteria of [19], and [10], and a novel local implementation of the convective criterion of [3]. The mass-lumping limiter is instead chosen based on the requirement that smooth extrema should be preserved, and is based on a smoothness sensor. The model obtained is extensively tested. The behaviour of different breaking models is studied on several cases allowing comparisons with experimental data.

## 2 Hybrid Equations for Wave Breaking Treatment

To simulate wave propagation, we start from the following system, based on the enhanced Boussinesq equations in the form proposed in [14] (cf. Fig. 1):

$$
\begin{cases}
\partial_t \eta + \partial_x q = 0 & \text{(1a)} \\[4pt]
\partial_t q + \partial_x(uq) + gh\partial_x\eta + ghC_f u = f_{\text{break}}(x, t)\mathscr{D}(\eta, q) & \text{(1b)} \\[4pt]
\mathscr{D}(\eta, q) = Bd^2\partial_{xxt}q + \beta gd^3\partial_{xxx}\eta + \dfrac{1}{3}d\partial_x d\partial_{xt}q + 2\beta gd^2\partial_x d\partial_{xx}\eta & \text{(1c)}
\end{cases}
$$

with $\eta = \eta(x, t)$ the wave height, $q = hu$ the discharge, $h = \eta + d$ the local height of the water column, $u = u(x, t)$ the depth-averaged velocity, and with $d = d(x)$ the depth w.r.t. an average still water level. The term $\mathscr{D}(\eta, q)$ represents the dispersive effects, with $B$ and $\beta$ obtained by optimizing the linear dispersion relation. The flag $f_{\text{break}}$ assumes the value 1 in the Boussinesq regions, and 0 in breaking fronts, and allows to revert to the hyperbolic shallow water equations. We consider here three breaking criteria.

The simplest, due to Tonelli and Petti [19], is based on a local measure of non-linearity. Breaking regions are denoted as those for which $\varepsilon = |\eta|/|d| \geq \varepsilon_{\text{cr}}$, with $\varepsilon_{\text{cr}} \approx 0.8$. Once a breaking front has been detected, its end (*de-breaking*) is located as the point in the flow direction where $\varepsilon$ is below $\approx 0.35$ (see [10, 19] for more).

The second criterion, proposed in [10], uses a hybrid condition involving vertical velocity and slope. A point is flagged as breaking if *either* $|\partial_t \eta| > \gamma \sqrt{gh}$ *or* $|\partial_x \eta| > \tan \phi_{\text{cr}}$. The values $\gamma$ and $\phi_{\text{cr}}$ may depend on the case simulated (see [10] for more).

Lastly, we consider a local implementation of the *convective* criterion of [3]. The idea is that breaking occurs when the free surface velocity is larger than the wave celerity. In [3] only simple cases have been considered for which at least the celerity is known a-priori. Here we proceed as follows:

1. Pre-flagging using the criterion of [10] with smaller $\gamma$ and $\phi_{\text{cr}}$;
2. For every front (set of neighbouring pre-flagged nodes) locate crest and trough;
3. For every front evaluate celerity $C_b$ and crest velocity $u_S$;
4. Final flagging: if $u_S \geq C_b$ set $f_{\text{break}} = 0$ for $x \in [x_{\min}, x_{\max}]$

Combining the relations $\partial_t \eta \approx -C_b \partial_x \eta$ and $\partial_t \eta = -\partial_x q$, we obtain $C_b \approx \partial_x q / \partial_x \eta$ which is implemented as $C_b = (q_{\text{crest}} - q_{\text{trough}})/(\eta_{\text{crest}} - \eta_{\text{trough}})$. To obtain $u_S$, vertical asymptotic expansions can be used to show that (see e.g. [3, 7]) $u_S = u - \alpha h^2 \partial_{xx} u$, with $\alpha = 1/3$ the analytical value. Here this constant is kept free, to account for the different wave shoaling provided by Boussinesq models, and to be able to correct wave under-shoaling [7]. The results reported are obtained with $\alpha = 2/3$. A parametric study is under way to understand the influence of this parameter for different Boussinesq equations. The definition of $[x_{\min}, x_{\max}]$, giving local position and width of the breaking region is the same used in [10, 18].

## 3 Discretization and Discontinuity Capturing

The numerical discretization follows the initial developments made in [16] where upwind stabilized residual based and finite element discretizations of the Boussinesq equations of [14] have been analyzed and tested on a large number of one and two-dimensional wave propagation problems. Already for $P^1$ interpolation, the results of [16] show a high potential of the approach in terms of providing low dispersion error and high accuracy with the flexibility of a natural unstructured mesh formulation.

Here we propose a discontinuity capturing method based on a nonlinear lumping of the mass matrix allowing to locally recover first order upwind flux differencing.

Set $\mathbf{W} = [\eta \ q]^T$, $F(\mathbf{W}) = [q \ (uq + g\frac{h^2}{2})]^T$, $S = [0 \ -gh\partial_x d]^T$, $D = [0 \ \mathscr{D}]^T$, $F_f = [0 \ -ghC_f u]^T$, and $A = \frac{\partial F(\mathbf{W})}{\partial \mathbf{W}}$ the shallow water flux Jacobian. Let also $d\mathbf{W}_i/dt$ be the (continuous) time derivative of the value of $\mathbf{W}$ at node $i$, $\Delta x$ the 1D mesh spatial spacing, $I_2$ the $2 \times 2$ identity matrix, and denote with superscripts $i \pm 1/2$ arithmetic cell-average values. The spatial discretization we propose reads:

$$\Delta x \frac{d\mathbf{W}_i}{dt} + \delta^{i-1/2}\{\frac{\Delta x}{6}[\frac{d\mathbf{W}_{i-1}}{dt} - \frac{d\mathbf{W}_i}{dt}] + \frac{\Delta x}{2}\text{sign}(A^{i-\frac{1}{2}})\frac{d\mathbf{W}^{i-\frac{1}{2}}}{dt}\}$$

$$+ \delta^{i+1/2}\{\frac{\Delta x}{6}[\frac{d\mathbf{W}_{i+1}}{dt} - \frac{d\mathbf{W}_i}{dt}] - \frac{\Delta x}{2}\text{sign}(A^{i+\frac{1}{2}})\frac{d\mathbf{W}^{i+\frac{1}{2}}}{dt}\}$$

$$+ \frac{I_2 + \text{sign}(A^{i-\frac{1}{2}})}{2}(F_i - F_{i-1} + \Delta x \, S^{i-\frac{1}{2}})$$

$$+ \frac{I_2 - \text{sign}(A^{i+\frac{1}{2}})}{2}(F_{i+1} - F_i + \Delta x \, S^{i+\frac{1}{2}}) = f_{\text{break}_i} D_i + \mathscr{F}_{f_i} \quad (2)$$

One can distinguish the terms associated to the Galerkin approximation, and the stabilization terms, multiplied by the sign of the Jacobian $A$. These terms have been simplified using the properties of the $P^1$ finite element approximation, as detailed in [16]. The right hand side contains the contributions of friction and dispersive terms, also involving centred and upwind biasing contributions, and requiring the evaluation of auxiliary variables necessary for the high order derivatives. These terms are quite complex and we refer to [16] for details. Note that if the right hand side is zero, for $\delta^{i\pm1/2} = 0$ we obtain the standard first order upwind flux difference scheme. Our implementation in this limit actually follows the well-balanced, positivity preserving upwind approach of e.g. [6], and it includes an entropy fix [9] to avoid problems in strongly accelerating regions with small water heights (cf. [1] for more). So, if $\delta_i = 0$ and $f_{\text{break}_i} = 0$ the scheme is locally first-order, it preserves the positivity of the depth, and it is well-balanced. Whenever $f_{\text{break}_i} = 1$, we automatically set $\delta_i = 1$. In this case, the resulting scheme is third-order accurate in space, as amply demonstrated in [16]. The main ingredient is the choice of the limiter $\delta(\mathbf{W})$. An extensive study and comparison of different limiters available in literature is provided in [1]. Many of these result in an over-dissipative method. An effective definition is based on the smoothness sensor

$$\sigma_i = \min(1, r_i), \quad r_i = \frac{\frac{|\eta_i - \eta_{i-1}|}{\Delta x} + \frac{|\eta_i - \eta_{i+1}|}{\Delta x}}{\frac{|\eta_{i+2} - 4\eta_{i+1} + 6\eta_i - 4\eta_{i-1} + \eta_{i-2}|}{12\Delta x^2}}$$

with $r_i$ the ratio between the magnitude of the first order derivative and the difference between a fourth and second order approximation of the second order derivative. In smooth regions, the denominator of $r_i$ is of $\mathscr{O}(\Delta x^2)$ while the numerator is bounded, resulting in $\sigma = 1$. On a discontinuity, while the numerator is of an order $\mathscr{O}(1/\Delta x)$, the denominator is of an order $\mathscr{O}(1/\Delta x^2)$, giving $\sigma = \mathscr{O}(\Delta x)$. Finally, we have set

$\delta_i = \sigma_i$ if $\sigma_i \le 1/2$, and $\delta_i = 1$ otherwise. The typical result obtained for a standard Riemann problem is reported on Fig. 2 where the sensor proposed is compared to the Superbee and to the Monotonized Central limiter [12]. In the tests that follow, as in [15] we pre-multiply $\delta$ by an exponential function smoothly switching off high order terms in vicinity of dry fronts. For all the tests considered, time integration has been performed with the non-dissipative second-order Crank-Nicholson scheme.

## 4 Numerical Validation

### 4.1 Periodic Wave over a Submerged Bar

We consider the experiment of plunging breaking periodic waves over a submerged bar of Beji and Battjes. This test has been first done by Dingemans to verify the Delft Hydraulics model HISWA, and then repeated by Beji and Battjes [7, 17]. To give an overview of the qualitative behaviour of the solution, wave profiles at different breaking instants are reported on Fig. 3 for the three tested breaking criteria. In the figure we report the wave profiles at the first breaking instance, at an intermediate time (same for all criteria), and at the last seen breaking instance for a given wave. The vertical lines delimit the breaking region in which the shallow water equations are used. The criterion of [10] provides the strongest and most regular breaking behaviour, with wave heights considerably decreasing along the plateau. The local implementation of the convective criterion proposed gives weaker breaking, and slightly higher waves. We have also observed numerically a more intermittent behaviour of the flag. Lastly, the criterion of [19] gives the weakest breaking, with wave heights only slightly decreasing.

These observations are confirmed by the temporal evolution of the wave height in four experimental gauges (respectively at the beginning and the end of the upward slope of the bar and in the middle and end of the plateau), reported in Fig. 4. The results obtained with the criterion of [10] show very good agreement with experiments, while the convective criterion is slightly worse in terms of wave heights. The non-linearity sensor of [19] fails to detect some wave breaking areas, at least on this level of mesh size. We mention that better results are obtained in [19] on much finer grids, and that the results of the convective criterion could be improved by increasing the value of the constant $\alpha$ in the definition of the free surface velocity (under investigation).

### 4.2 Run-up of a Periodic Wave

This test, known as the spilling breaking test of Hansen and Svendsen, involves the shoaling and breaking over a shore of a set of regular waves, and corresponds to the test 051041 described in [8]. A qualitative view of the wave profiles obtained is
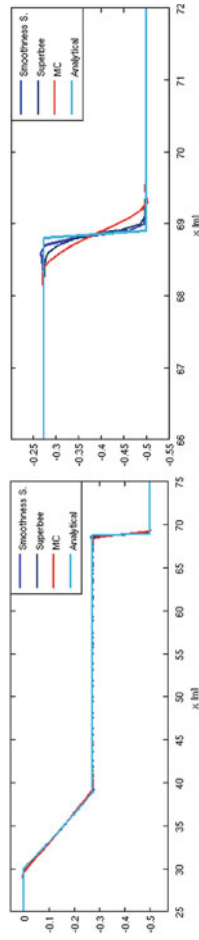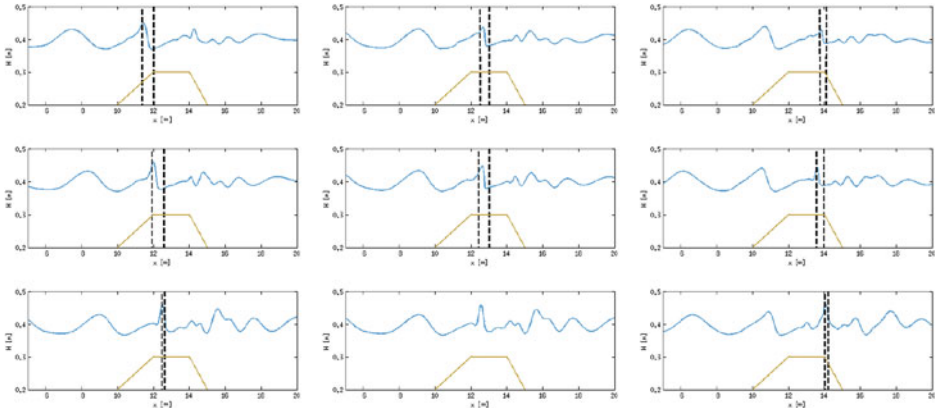
**Fig. 2** Riemann problem: smoothness sensor versus Superbee and Monotonized Central limiters

**Fig. 3** Plunging breaking test. Wave profiles corresponding to: first (*left*), intermediate (*center*), and last breaking instants. Breaking criterion: Kazolea, Delis and Synolakis (*top*), convective (*middle*), Tonelli and Petti (*bottom*). The *vertical lines* delimit the breaking (shallow water) region

reported on Fig. 5, showing the effects of wave shoaling and wave breaking over a constant slope bathymetry. On Fig. 6, instead, we report a quantitative comparison of the time-average of the wave height and of the wave set-up along the shore with experimental data. The numerical results are those obtained with the criterion of [10], and with the convective criterion. We can see from the change in slope in the computed results that wave breaking is predicted too early by the criterion of Kazolea, Delis and Synolakis, while the wave heights are under-predicted in both cases. This, according to [20], might be due to the use of a weakly nonlinear Boussinesq model for propagation. For this test, the convective criterion gives a better prediction of the breaking position. The wave set-up is predicted very well by both models.

## 5 Conclusions and Perspectives

We have presented a one-dimensional finite element model for non-hydrostatic wave propagation, breaking, and run-up. The model combines a weakly non-linear Boussinesq model with the hyperbolic shallow water equations. The blending is obtained via a wave breaking criterion based on physical arguments. We propose an enhancement of the stabilized finite element method of [16] consisting in a discontinuity capturing technique relying on a nonlinear lumping of the mass matrix. This allows the local treatment of discontinuous shallow water flows, and wetting/drying fronts. When combined with the hybrid dispersive-hyperbolic modeling approach, this method allows to provide an accurate description of the wave transformation in the near shore region.

**Fig. 4** Gauge data for the submerged bar test. *Top* gauge 1 (*left*) and 2 (*right*). *Bottom* gauge 3 (*left*) and 4 (*right*). Breaking criteria: Tonelli and Petti (local in the legend), Kazolea, Delis and Synolakis (hybrid in the legend), and convective (physical in the legend). Symbols: experiments

**Fig. 5** Hansen and Svendsen test 051041. Snapshots of the wave profiles at different instants

**Fig. 6** Hansen and Svendsen test 051041. Wave height (*left*) and mean water level (set-up) (*right*) for Kazolea, Delis and Synolakis (*top row*, hybrid in the legend), and convective criterion (*bottom row*, physical in the legend)

The numerical results, while confirming the stability and robustness of the numerics proposed, also provide an initial validation for the different breaking criteria tested. Our implementation of the convective breaking criterion of [3] shows some promise, even though the criterion of Kazolea et al. gives similar, and sometimes better, results, with a much simpler implementation. The very simple criterion of [19] is not able to provide similar results.

The work planned for the future involves a more systematic study of the definition of the free surface velocity used on the convective criterion, the implementation of the model in two dimensions and on unstructured meshes, following [10, 16], and the use of fully non-linear dispersive models, as in [5, 11].

## References

1. Bacigaluppi, P.: Upwind stabilized finite element modeling of non-hydrostatic wave breaking and run-up (2013). MSc Thesis, Aerospace Department, Politecnico di Milano
2. Barthélemy, E.: Nonlinear shallow water theories for coastal waves. Surv. Geophys. **25**, 315–337 (2004)
3. Bjørkavåg, M., Kalisch, H.: Wave breaking in boussinesq models for undular bores. Phys. Lett. A **375**(14), 1570–1578 (2011)
4. Bonneton, P.: Modelling of periodic wave transformation in the inner surf zone. Ocean Eng. **34**, 1459–1471 (2007)
5. Bonneton, P., Chazel, F., Lannes, D., Marche, F., Tissier, M.: A splitting approach for the fully nonlinear and weakly dispersive greennaghdi model. J. Comput. Phys. **230**(4) (2011)
6. Castro, M., Ferrero, A., García-Rodríguez, J., González-Vida, J., Macías, J., Pareés, C., Vázquez-Cendón, M.: The numerical treatment of wet/dry fronts in shallow flows: application to one-layer and two-layer systems. Math. Comput. Model. **42**, 419–439 (2005)
7. Dingemans, M.: Water Wave Propagation Over Uneven Bottoms: Linear wave propagation. Advanced series on ocean engineering. World Scientific Pub, Singapore (1997)
8. Hansen, J., Svendsen, I.: Regular waves in shoaling water: experimental data. Tech. Rep. 21, Technical Report, ISVA series paper (1979)
9. Harten, A., Hyman, J.: Self-adjusting grid methods for one-dimensional hyperbolic conservation laws. J. Comput. Phys. **50**(2) (1983)
10. Kazolea, M., Delis, A., Synolakis, C.: Numerical treatment of wave breaking on unstructured finite volume approximations for extended boussinesq-type equations. J. Comput. Phys. (2014). http://dx.doi.org/10.1016/j.jcp.2014.01.030
11. LeMétayer, O., Gavrilyuk, S., Hank, S.: A numerical scheme for the green-naghdi model. J. Comput. Phys. **229**(6) (2010)
12. LeVeque, R.: Finite-Volume Methods for Hyperbolic Problems. Cambridge Texts in Applied Mathematics. Cambridge University Press, Cambridge (2004)
13. Ma, G., Shi, F., Kirby, J.: Shock-capturing non-hydrostatic model for fully dispersive surface wave processes. Ocean Model **43–44**, 22–35 (2012)
14. Madsen, P., Sørensen, O.: A new form of the Boussinesq equations with improved linear dispersion characteristics. A slowly-varying bathymetry. Coast. Eng. **18**, 183–204 (1992)
15. Ricchiuto, M., Bollermann, A.: Stabilized Residual Distribution for shallow water simulations. J. Comput. Phys. **228**, 1071–1115 (2009)
16. Ricchiuto, M., Filippini, A.: Upwind residual discretization of enhanced boussinesq equations for wave propagation over complex bathymetries. J. Comput. Phys. (2014). http://dx.doi.org/10.1016/j.jcp.2013.12.048

17. Shiach, J., Mingham, C.: A temporally second-order accurate Godunov-type scheme for solving the extended Boussinesq equations. Coast. Eng. **56**, 32–45 (2009)
18. Tissier, M., Bonneton, P., Marche, F., Chazel, F., Lannes, D.: A new approach to handle wave breaking in fully non-linear boussinesq models. Coast. Eng. **67**, 54–66 (2012)
19. Tonelli, M., Petti, M.: Simulation of wave breaking over complex bathymetries by a Boussinesq model. J. Hydraulic Res. **49** (2011)
20. Veeramony, J., Svendsen, I.: The flow in surf-zone waves. Coast. Eng. **39**, 93–122 (2000)

# A Quasi-1D Model of Biomass Co-Firing in a Circulating Fluidized Bed Boiler

**Michal Beneš, Pavel Strachota, Radek Máca, Vladimír Havlena and Jan Mach**

**Abstract**  We introduce an outline of the mathematical model of combustion in circulating fluidized bed boilers. The model is concerned with multiphase flow of flue gas, bed material, and two types of fuels (coal and biomass) that can be co-fired in the furnace. It further considers phase interaction resulting in particle attrition, devolatilization and burnout of fuel particles, and energy balance between heat production and consumption (radiative and convective transfer to walls). Numerical solution by means of the finite volume method together with a Runge-Kutta class time integration scheme is mentioned only briefly as the used methods are generic and well documented elsewhere. Some representative results are also presented.

**MSC2010:** 65M08 · 76T30 · 80A25

M. Beneš · P. Strachota (✉) · R. Máca · J. Mach
Department of Mathematics, Faculty of Nuclear Sciences and Physical Engineering, Czech
Technical University in Prague, Prague, Czech Republic
e-mail: pavel.strachota@fjfi.cvut.cz

M. Beneš
e-mail: michal.benes@fjfi.cvut.cz

R. Máca
e-mail: radek.maca@fjfi.cvut.cz

J. Mach
e-mail: jan.mach@fjfi.cvut.cz

V. Havlena
Honeywell ACS Advanced Technology Laboratory in Prague, Prague, Czech Republic
e-mail: vladimir.havlena@honeywell.com

# 1 Introduction

Fluidized bed (FB) combustion [3, 5, 21] is a modern technology of industrial energy production from both fossil and renewable fuel sources. FB boilers are very flexible with respect to fuel properties. They also release limited amounts of pollutants into the atmosphere thanks to relatively low combustion temperatures and sulfur oxides being captured by the fluidized bed material. Mathematical modeling and numerical simulations of FB combustion are helpful in design and control of energy production facilities.

We propose a model of a circulating fluidized bed (CFB) [3] combustion chamber extended by a simple description of solid phase recirculation. It is able to capture the physical aspects of fluidized flow and combustion processes. Such model is sufficient to study the temperature and dynamic behavior of the fluidized bed which in turn is important for designing a control strategy leading to optimal operation conditions. Combustion chemistry is not considered at the moment.

# 2 Summary of Governing Equations

The model consists of equations of multiphase flow, fuel transformation, and heat transport in a variable cross section duct, originally based on two-phase flow of gas (air) and solid found in [10]. The conservation laws are formulated for *quasi-one-dimensional nozzle flow* [1, Chap. 7] along the vertical axis $x$ of the combustor which has a rectangular cross section $A(x)$.

The mixture comprises the flue gas phase and three solid phases: the bed material (granular *dolomite* or *limestone*) denoted by the index $s$, granular *coal char* ($c$), and granular *biomass char* ($b$). The volume fractions of the individual phases $i \in \{g, s, c, b\}$ are represented by the quantities $\varepsilon_i(t, x)$ satisfying the relation $\varepsilon_g + \varepsilon_s + \varepsilon_c + \varepsilon_b = 1$. We denote the velocity of phase $i$ by $u_i$ and its density by $\rho_i$. The remaining variables are the concentrations (mass fraction) of oxygen $Y_{O_2}$ and the *volatile matter* (VM, see Sect. 3) released from both fuels.

*Mass Balance* The *continuity equations* for the individual phases $i \in \{g, s, c, b\}$, and gaseous species $X \in \{O_2, VM_c, VM_b\}$ read

$$\frac{\partial (\rho_i \varepsilon_i)}{\partial t} + \frac{1}{A} \frac{\partial (A \rho_i \varepsilon_i u_i)}{\partial x} = \mathcal{M}_i(t, x) + \mu_i, \tag{1}$$

$$\frac{\partial (\rho_g \varepsilon_g Y_X)}{\partial t} + \frac{1}{A} \frac{\partial (A \rho_g \varepsilon_g u_g Y_X)}{\partial x} = \mathcal{M}_X(t, x) + \mu_X \tag{2}$$

where the bulk solid material densities $\rho_i$, $i \in \{s, c, b\}$ are assumed to be constant. The source terms $\mathcal{M}_g, \mathcal{M}_i$ $[\text{kg} \cdot \text{m}^{-3} \text{s}^{-1}]$ describe the injection of the secondary air together with "instantaneously" released *volatiles*, bed material, and solid fuel components into the furnace at the respective elevations. $\mathcal{M}_{O_2}$ $[\text{kg} \cdot \text{m}^{-3} \cdot \text{s}^{-1}]$

is the inflow of oxygen as part of the secondary and tertiary air. $\mathcal{M}_{VM_c}$, $\mathcal{M}_{VM_b}$ $\left[kg \cdot m^{-3} \cdot s^{-1}\right]$ is the inflow of VM as part of the fuel inflow. The terms $-\mu_c$, $-\mu_b$ account for mass loss due to char particle burnout. This mass converts into flue gas and *ash*, which becomes part of the bed material ($\mu_s$). The term $\mu_g$ is the production of flue gas due to combustion of char and $-\mu_{O_2}$ is the rate of oxygen consumption by combustion. The terms $\mu_i$ satisfy $\mu_g + \mu_s + \mu_c + \mu_b = 0$. Lastly, the terms $\mu_{VM_c}$, $\mu_{VM_b}$ account for the burnout of coal and biomass VM.

*Passive Particle Transport* Combustion rate of fuel particles and attrition dynamics depend on the size and shape of the solid particles. Given the number of particles per unit volume $n_i$ for the phase $i \in \{s, c, b\}$, the average mass and diameter of one (spherical) particle are equal to $m_i^{(1)} = \frac{\rho_i \varepsilon_i}{n_i}$, $d_{p,i} = \sqrt[3]{\frac{6\varepsilon_i}{\pi n_i}}$, respectively. The quantities $n_i$ are subject to *passive transport* described by the equations

$$\frac{\partial n_i}{\partial t} + \frac{1}{A} \frac{\partial (An_i u_i)}{\partial x} = \mathcal{M}_{n_i} (t, x) + \mu_{n_i} \tag{3}$$

for each $i \in \{s, c, b\}$. The source term $\mathcal{M}_{n_i} \left[m^{-3} \cdot s^{-1}\right]$ in (3) is nonzero where the solid phase enters the combustor. The term $\mu_{n_i}$ stands for particle number change due to *attrition* and *conversion of burnt out fuel mass into ash* (Sect. 3).

*Momentum Balance* The *momentum equations* for the individual phases $i \in \{g, s, c, b\}$ assume the form

$$\frac{\partial (\rho_i \varepsilon_i u_i)}{\partial t} + \frac{1}{A} \frac{\partial \left(A\rho_i \varepsilon_i u_i^2\right)}{\partial x} = -\frac{\partial P_i}{\partial x} + \sum_{k \in \{g,s,c,b\}} \beta_{ki} (u_k - u_i) - \frac{2 f_i \varepsilon_i \rho_i u_i^2}{D}$$
$$+ R_i \varepsilon_i g + \mathcal{P}_i (t, x) + \pi_i \tag{4}$$

where $P_g$ is the pressure of gas and for $i \in \{s, c, b\}$, we put $\frac{\partial P_i}{\partial x} = G\frac{\partial \varepsilon_i}{\partial x}$ where $G\left(\varepsilon_g\right)$ is the solids stress modulus [10]. The coefficients $f_i$ and $\beta_{ki}$ express the wall friction and inter-phase friction forces. In the gravity/buoyancy term $R_i \varepsilon_i g$, we have $R_g = \rho_g$ and $R_i = \rho_i - \rho_g$ for $i \in \{s, c, b\}$. The source terms $\mathcal{P}_i \left[kg \cdot m^{-2} \cdot s^{-2}\right]$ account for the density of vertical momentum of the injected material. The term $\pi_g$ represents the vertical momentum of the fuels burnt into flue gas. $\pi_s$ is the vertical momentum of the ash which remains after burning $\mu_g + \mu_c$ fuels. $\pi_c$ and $\pi_b$ are the vertical momenta of the burnt out fuels with masses $\mu_c$ and $\mu_b$, respectively. The equality $\pi_g + \pi_s + \pi_c + \pi_b = 0$ holds.

*Energy Balance* We assume local thermal equilibrium between all phases at each point. Extending the steps in [1, Chap. 7], we derive the equations for internal energy $e_i$ of each phase $i \in (g, s, c, b)$ and sum them to arrive at the single *equation for internal energy* of the phase mixture in terms of temperature $T$

$$\sum_{i \in \{g,s,c,b\}} \rho_i \varepsilon_i c_{p,i} \left( \frac{\partial T}{\partial t} + u_i \frac{\partial T}{\partial x} \right) = \frac{1}{A} \sum_{i \in \{g,s,c,b\}} \left( -P_i \frac{\partial (A u_i)}{\partial x} + \varepsilon_i \frac{\partial}{\partial x} \left( A \lambda_i \frac{\partial T}{\partial x} \right) \right)$$
$$+ \sum_{i \in \{g,s,c,b\}} \left( \mathscr{E}_i - u_i (\mathscr{P}_i + \pi_i) + (\mathscr{M}_i + \mu_i) \left( \frac{u_i^2}{2} - \int_0^T c_{p,i}(\tau) \, d\tau \right) \right) + \dot{Q}$$

(5)

where $\dot{Q} \left[ W \cdot m^{-3} \right]$ is the total heat source term consisting of heat production by combustion and consumption by radiative and convective heat transfer. The system is closed by the equation of state for ideal gas, i.e. $\rho_g = \frac{P}{R_{spec} T}$ where $R_{spec}$ $\left[ J \cdot kg^{-1} \cdot K^{-1} \right]$ is the specific gas constant.

## 3 Modeling the Particular Effects

*Inter-phase Momentum Transfer and Wall Friction* The drag between the gas phase $g$ and the phase $i \in \{s, c, b\}$ appears in (4) as the inter-phase friction coefficient $\beta_{ig}$ given by the empirical formula adopted from [10]. The formula for the solid–solid momentum transfer coefficient $\beta_{ki}$ between two solid phases $i, k \in \{s, c, b\}$ has been taken from [18] with additional parameter settings based on the data in [14]. The wall friction factors $f_g$, $f_i$ are given by the modified Hagen-Poiseuille expression and by the Konno-Saito correlation [10].

*Attrition.* *Attrition* of solid particles in the fluidized bed is a joint effect of abrasion and fragmentation [21]. However, it is generally believed that the contribution of fragmentation is negligible [12, 20]. Therefore, only abrasion is modeled.

For *bed material attrition*, we use the model of [19] with a single equation for particle mass loss. It is aimed primarily at circulating fluidized beds and overcomes the drawbacks of several previously published models. For *fuel char particles*, it is necessary to take into account *combustion-assisted attrition*. The formula proposed in [3, p. 116] and [4] is used.

For each solid phase $i \in \{s, c, b\}$, the particle mass attrition rate per unit volume $\dot{m}_i$ is converted to the rate of increase in the number of the particles by the relation $\mu_{n_i} = \dot{n}_i = -\frac{n_i}{m_i} \dot{m}_i$. This corresponds to a constant mass being distributed among an increasing number of particles with a decreasing average size.

*Fuel Devolatilization and Burnout* As the particle heats up, volatiles release and burn simultaneously (they form the *volatile flame*) [5]. It is not possible to track the age of each individual particle in the proposed model. As the time to complete
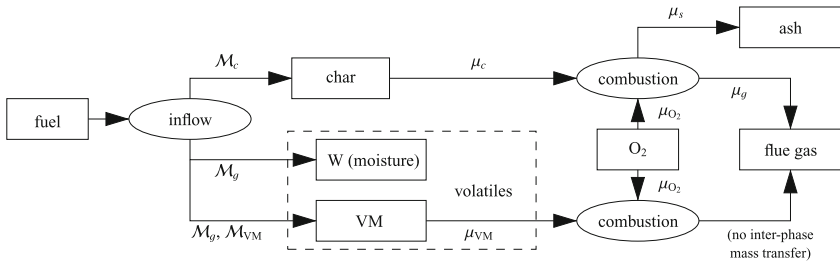
**Fig. 1** Simplified fuel conversion in the proposed CFB combustion model

devolatilization is an order of magnitude shorter than the time to complete char burnout [3]), we assume that the fuel particles enter the combustion chamber already fully devolatilized and we further treat the burnout of VM and the solid char particles separately. The mass exchange terms in Eqs. (1)–(2) involved in fuel conversion (devolatilization and burnout) are depicted in Fig. 1.

We assume that *VM burnout* is governed by the Arrhenian kinetic theory and depends primarily on temperature and oxygen concentration. The rate of *char burnout* is controlled by two independent processes: the transfer of oxygen to the particle surface from the ambient gas and by the reaction rate of pure carbon with oxygen [4]. Combining the correlations for both, one obtains the combustion rate $\mu_i$. *Oxygen consumption* can then be calculated using the ultimate analysis [5] of the fuels. Finally, the *lower heating value* of the fuel is used to calculate the *heat production*.

*Heat Transfer* Volumetric heat consumption is caused by radiative and convective transfer to walls. Currently, we employ empirical formulas that do not take into account the effect of the solid phase presence.

The *radiative heat transfer* to walls per unit volume is given by $\dot{Q}_R = \sigma \left( \varepsilon (T) T^4 - \varepsilon (T_{\text{wall}}) T_{\text{wall}}^4 \right)$ where $\sigma$ is the Boltzmann constant and $\varepsilon$ is the emissivity (and absorptivity) of the flue gas [13, p. 277]. $\varepsilon$ depends on the concentration of the radiant species $H_2O$ and $CO_2$.

The *convective heat transfer* occurs directly at the walls and is given by the term $\dot{Q}_C = \frac{4}{D} \alpha (T - T_{\text{wall}})$ where $\alpha$ $\left[ \text{W} \cdot \text{m}^{-2} \cdot \text{K}^{-1} \right]$ is the convective heat transfer coefficient per unit wall surface. The formula for $\alpha$ has been taken from [15].

## 4 Numerical Solution and Implementation

For numerical solution, the proven combination of a finite volume scheme [7] on a regular 1D mesh with the Runge-Kutta-Merson adaptive explicit time solver [6, 11] is employed. The system of model Eqs. (1)–(5) can be written in aggregate vector form

$$\frac{\partial \mathbf{W}}{\partial t} + \frac{\partial}{\partial x} (\mathbf{F}_c - \mathbf{F}_v) = \mathbf{Q}$$

**Fig. 2** Time progress of the cumulative quantities in the time range from $t = 0$ s to $t = 2000$ s

where $\mathbf{W}$ is the vector of state variables and $\mathbf{F}_c$, $\mathbf{F}_v$ the convective and viscous fluxes, respectively. The corresponding semidiscrete finite volume scheme for the unknown numerical solution $\mathbf{W}_K$ reads

$$\frac{\mathrm{d}\mathbf{W}_K}{\mathrm{d}t} = \frac{1}{|K|} \left[ \left( \mathbf{F}_{c,R} - \mathbf{F}_{v,R} + \mathbf{F}_{v,R} \right) - \left( \mathbf{F}_{c,L} - \mathbf{F}_{v,L} + \mathbf{F}_{v,L} \right) \right] - \mathbf{Q}_K$$

where the subscripts $R$, $L$ indicate the values of the respective fluxes at the left and right boundary of the cell $K$ with size $|K|$. Adjustable artificial diffusion $\mathbf{F}_v$ is used for stabilization.

The *boundary conditions* are implemented by means of auxiliary (ghost) cells beyond the computational domain. By default, zero Neumann b.c. (extrapolation) is imposed on all *physical* quantities (e.g. the conservative variable $\rho_i \varepsilon_i u_i$ in (4) consists of the physical quantities $\rho_i$, $\varepsilon_i$, and $u_i$). This is overriden at the inlet by Dirichlet b.c. for the volumetric fractions $\varepsilon_i$, particle diameters $d_{p,i}$, temperature $T$, and gas mass inflow $\rho_g \varepsilon_g u_g$. At the outlet, the Dirichlet b.c. for pressure is prescribed.

*Initial conditions* are set up to quickly reach the stationary state with the given boundary conditions. In brief, zero flow velocities, atmospheric pressure, zero fuel inventory, and gas pre-heating is assumed.

**Fig. 3** Spatial profiles of the selected quantities along the vertical axis of the combustor at the time $t = 2000\,\mathrm{s}$

The algorithm is parallel and has been written using C/C++ and the MPI library [17] for computation on an arbitrary number of CPU cores of a Linux HPC cluster.

## 5 Results

For the simulations, we currently use the data from the technical documentation of one particular CFB heating plant and fuel data from various other sources such as official catalogs of coal and biomass producers. Adaptation of the model to the real device requires careful setting of initial and boundary conditions as well as the source terms (e.g. for implementing controlled recirculation of solids). Detailed discussion is beyond the scope of this paper. Figure 2 shows the time progress of the cumulative (integral) quantities in a sample simulation and Fig. 3 contains the spatial profiles of the selected quantities at the end of the simulation. The model exhibits qualitative agreement with the expected behavior of the device, as can partly be observed in both figures. For example, the shape of gas velocity and temperature profiles as well as the ratio of total heat production to the heat transfer to the walls agree with the information in the technical documentation. Even though the validation of the model has not yet been completed, it already represents a cornerstone for further parameters fitting and step response measurement with potential application in the design of model-based predictive control.

# 6 Conclusion

Recently, there has been an extensive increase of interest in detailed non-stationary CFD simulations of fluidized bed combustion (see [16] for a review). There exist several comprehensive 3D simulation instruments based either on continuous, discrete, or combined (multiphase particle-in-cell [2]) approach. All of them generally require enormous computational resources.

On the other hand, we propose a model aimed for use as a vehicle for predictive control strategies development and testing. Its Eulerian–Eulerian multiphase flow model is based on the widely used equations of Gidaspow [10]. In addition, several results of independent modeling and measurements of the important phenomena of fluidized bed combustion are combined in the model, as described in Sect. 3. The choice of the variables and input parameters is in accordance with the purpose and possibilities of the control mechanisms [8, 9]. The quasi-1D approach provides a reasonable approximation of the underlying 3D geometry while maintaining low computational complexity. This results in faster than real time simulations on a single multi-core workstation.

# References

1. Anderson, J.D.: Computational Fluid Dynamics: The Basics with Applications. McGraw-Hill, New York (1995)
2. Andrews, M.J., O'Rourke, P.J.: The multiphase particle-in-cell (MP-PIC) method for dense particulate flows. Int. J. Multiphase Flow **22**(2), 379–402 (1996)
3. Basu, P.: Combustion and Gasification in Fluidized Beds. CRC Press, Boca Raton (2006)
4. Basu, P., Halder, P.K.: Combustion of single carbon particles in a fast fluidized bed of fine solids. Fuel **68**, 1056–1063 (1989)
5. Basu, P., Kefa, C., Jestin, L.: Boilers and Burners: Design and Theory. Springer, New York (1999)
6. Christiansen, J.: Numerical solution of ordinary simultaneous differential equations of the 1st order using a method for automatic step change. Numer. Math. **14**, 317–324 (1970)
7. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Ciarlet, P.G., Lions, J.L. (eds.) Handbook of Numerical Analysis, vol. 7, pp. 715–1022. Elsevier, New York (2000)
8. Findejs, J., Havlena, V., Jech, J., Pachner, D.: Model based control of the circulating fluidized bed boiler. In: Proceedings of the 9th IFAC Symposium on Power Plant and Power Systems, Control (2009)
9. Findejs, J., Havlena, V., Pachner, D.: Multivariable predictive circulating fluidized bed combustor control. ATP J. PLUS **2**, 12–15 (2008)
10. Gidaspow, D.: Multiphase Flow and Fluidization: Continuum and Kinetic Theory Description. Academic Press, Bostom (1994)
11. Holodniok, M., Klíč, A., Kubíček, M., Marek, M.: Methods of Analysis of Nonlinear Dynamical Models. Academia, Praha (1986)

12. Jiang, X., Zhou, L., Liu, J., Han, X.: A model on attrition of quartzite as a bed material in fluidized beds. Powder Technol. **195**, 44–49 (2009)
13. Modest, M.F.: Radiative Heat Transfer, 2nd edn. Academic Press, New York (2003)
14. Owoyemi, O., Lettieri, P.: A CFD study into the influence of the particle particle drag force on the dynamics of binary gas solid fluidized beds. In: The 12th International Conference on Fluidization—New Horizons in Fluidization, Engineering (2007)
15. Senkara, T.: Obliczenia cieplne pieców grzewczych w hutnictwie żelaza. Katowice: Wydaw. "Ślask" (1968). In Polish
16. Singh, R.I., Brink, A., Hupa, M.: Cfd modeling to study fluidized bed combustion and gasification. Appl. Therm. Eng. **52**, 585–614 (2013)
17. Snir, M., Otto, S., Huss-Ledermann, S., Walker, D., Dongarra, J.: The Complete MPI Reference. The MIT Press, Boston (1995)
18. Syamlal, M.: The particle-particle drag term in a multiparticle model of fluidization. Topical report DOE/MC/21353-2373, NTIS/DE87006500, National Technical Information Service, Springfield, VA (1987)
19. Tomeczek, J., Mocek, P.: Attrition of coal ash particles in a fluidized-bed reactor. AIChE J. **53**(5), 1159–1163 (2007)
20. Werther, J., Reppenhagen, J.: Catalyst attrition in fluidized-bed systems. AIChE J. **45**(9), 2001–2010 (1999)
21. Yang, W.C. (ed.): Handbook of Fluidization and Fluid-Particle Systems. Marcel Dekker, New York (2003)

# Simulation of Diluted Flow Regimes in Presence of Unsteady Boundaries

**Florian Bernard, Angelo Iollo and Gabriella Puppo**

**Abstract** The main feature of diluted flows is the presence of both continuum and kinetic regimes in the same field. The ES-BGK model is a kinetic model that preserves the asymptotic properties towards compressible Euler equations in the hydrodynamic regime, yet modeling momentum and kinetic energy diffusion for low Knudsen numbers. Here, this model is discretized by a finite-volume scheme on Cartesian meshes. The scheme is second order up to the possibly moving boundaries. To ensure a smooth transition between the hydrodynamic and the kinetic regime up to the walls, appropriate boundary conditions are devised. As an application, we present the simulation of an unsteady nozzle plume in a very low pressure environment.

## 1 Introduction

The Boltzmann equation models flow regimes where the mean free path $\lambda$ between particle shocks is larger than the characteristic length of the problem $L$ under consideration. The ratio $\lambda/L$ is called the Knudsen number and becomes large in the rarefied regime.

F. Bernard (✉)
Department of Mechanical and Aerospace Engineering, Politecnico di Torino,
Torino, Italy
e-mail: florian.bernard@math.u-bordeaux1.fr; florian.bernard@polito.it

A. Iollo
Univ. Bordeaux, IMB, UMR 5251, 33400 Talence, France

A. Iollo
Inria, 33400 Talence, France
e-mail: angelo.iollo@math.u-bordeaux1.fr

G. Puppo
Dip. di Scienza ed Alta Tecnologia, Università dell'Insubria, Como, Italy
e-mail: gabriella.puppo@uninsubria.it

Here, the ES-BGK model [1, 7], is considered for its capability to ensure a smooth transition between the rarefied and the continuum regime. With respect to the full Boltzmann equation, or to other models like DSMC [5] and BGK model [4], the computational cost is reasonable and it preserves the correct Prandtl number $Pr$.

In this work, a finite-volume method on Cartesian meshes is presented to discretized the ES-BGK model preserving the asymptotic properties. In this sense, a new boundary condition is proposed to avoid spurious effects due to the discretization [3]. In particular, this approach is validated with experimental data on a nozzle spreading jet in a low pressure environment.

## 2 Governing Equations

In the ES-BGK model, the collision term is approximated as a relaxation towards a Gaussian function:

$$\frac{\partial f}{\partial t}(x, \xi, t) + \xi \cdot \nabla_x f(x, \xi, t) = \frac{1}{\tau}(\mathscr{G}_f(x, \xi, t) - f(x, \xi, t)) \tag{1}$$

where $\tau$ is the relaxation time and $\mathscr{G}_f$ is the Gaussian distribution function, obtained as follows:

$$\mathscr{G}_f(x, \xi, t) = \frac{\rho(x, t)}{\sqrt{\det(2\pi \mathscr{T}(x, t))}} \exp\left(-\frac{(\xi - U(x, t))\mathscr{T}(x, t)^{-1}(\xi - U(x, t))^T}{2}\right) \tag{2}$$

where $R$ is the universal gas constant, $T(x, t)$, $U(x, t)$, $\rho(x, t)$ are the macroscopic values of temperature, velocity, density respectively. The tensor $\mathscr{T}$ is defined with the opposite stress tensor $\Theta(x, t)$ and the identity matrix $I$ as:

$$\mathscr{T}(x, t) = \frac{1}{Pr}T(x, t)I + (1 - \frac{1}{Pr})\Theta(x, t) \tag{3}$$

Macroscopic quantities are calculated from the moments of $f$ defined by:

$$\begin{pmatrix} \rho(x, t) \\ \rho(x, t)U(x, t) \\ E(x, t) \\ \rho(x, t)\Theta(x, t) \end{pmatrix} = \int_{\mathbb{R}^N} f(x, \xi, t)m(\xi)d\xi \quad \text{with} \quad m(\xi) = \begin{pmatrix} 1 \\ \xi \\ \frac{1}{2} |\xi|^2 \\ c \otimes c \end{pmatrix} \tag{4}$$

Here $c$ is the relative microscopic velocity $(\xi - U(x, t))$, $E$ is the total energy. We consider a mono-atomic gas for which the ratio of specific heats is $\gamma = 5/3$ and $N = 3$. Hence, this model does not take into account other energy degrees of freedom like in the case of a polyatomic gas.

The relaxation time for the ES-BGK model can be written in dimensionless form as:

$$\tau^{-1} = k\rho T^{1-\nu} \quad \text{with} \quad k = \frac{RT_0^\nu}{Pr\mu_0} = \frac{1}{PrKn_\infty} \tag{5}$$

where $\nu$ is the exponent of the viscosity law of the gas, $\mu_0$ is the reference viscosity of the gas at the reference temperature $T_0$ and $Kn_\infty$ the Knudsen number in reference conditions and $Pr$ is the Prandtl number.

## 3 Space and Time Discretization

The physical space $\Omega$ is discretized on a Cartesian grid with $n \times m$ cells:

$$\Omega = \bigcup_{\substack{i=1..n \\ j=1..m}} \Omega_{i,j} = \bigcup_{\substack{i=1..n \\ j=1..m}} [x_{i-1/2}, x_{i+1/2}] \times [y_{j-1/2}, y_{j+1/2}]$$

such that $(x_i, y_j)$ are the coordinates of the center of cell $(i, j)$ and $(x_{i+1/2}, y_j)$ are the coordinates of the center of the interface between cells $(i, j)$ and $(i + 1, j)$. On a space cell $\Omega_{i,j} = \left[x_i - \frac{\Delta x}{2}, x_i + \frac{\Delta x}{2}\right] \times \left[y_j - \frac{\Delta y}{2}, y_j + \frac{\Delta y}{2}\right]$, Eq. (1) is integrated with a finite volume method:

$$\frac{\partial f_{i,j}}{\partial t} + \xi \cdot \int_{\partial \Omega_{i,j}} f n_{\partial \Omega_{i,j}} d\sigma = \frac{1}{\tau_{i,j}}(\mathscr{G}_{f_{i,j}} - f_{i,j}) \tag{6}$$

where $f_{i,j} = \frac{1}{|\Omega_{i,j}|} \int_{\Omega_{i,j}} f dx dy$ and $\mathscr{G}_{f_{i,j}} = \frac{1}{|\Omega_{i,j}|} \int_{\Omega_{i,j}} \mathscr{G}_f dx dy$.

Since a uniform Cartesian grid is considered, the equation can be simply rewritten in terms of fluxes at each numerical interface (between two cells):

$$\frac{\partial f_{i,j}}{\partial t} + \frac{1}{\Delta x}(F_{i+\frac{1}{2},j} - F_{i-\frac{1}{2},j} + F_{i,j+\frac{1}{2}} - F_{i,j-\frac{1}{2}}) = \frac{1}{\tau}(\mathscr{G}_{f_{i,j}} - f_{i,j}) \tag{7}$$

with $F_{i+\frac{1}{2},j}$ the numerical flux between cell $\Omega_{i,j}$ and cell $\Omega_{i+1,j}$ (with a similar notation for the other fluxes) which is expressed as :

$$F_{i+\frac{1}{2},j} = \max(0, \xi_u) f_l + \min(0, \xi_u) f_r \tag{8}$$

with $f_r$ and $f_l$ the values of $f$ on the two sides of the interface and $\xi_u$ the first component of the microscopic velocity. The numerical expression of the distribution functions $f_l$, $f_r$ depends on the reconstruction used at the numerical interface. For a first order reconstruction, $f_l = f_{i,j}$ and $f_r = f_{i+1,j}$. For second order accuracy, a MUSCL reconstruction with slope limiters (MinMod for example) is employed.

In principle, the time discretization can be performed for all terms explicitly. But in this case, the time step will be determined by the space discretization ($\Delta x$), the maximum velocity of the velocity grid and the relaxation time $\tau$. For small Knudsen numbers, the relaxation part becomes very stiff ($\tau$ very small) and imposes a very strong restriction on the time step. Asher et al. [2] first presented IMEX schemes to cure this issue. Here, the IMEX scheme [6] is chosen. The relaxation term is treated implicitly while the convective part is non stiff but highly non linear which means that an explicit scheme is more efficient.

The time integration for a $\nu$-stages IMEX Runge-Kutta scheme reads as follows:

$$
\begin{aligned}
f_{i,j}^{n+1} &= f_{i,j}^n - \Delta t \sum_{k=1}^{\nu} \tilde{\omega}_k \xi \nabla_x f_{i,j}^{(k)} + \frac{\Delta t}{\tau} \sum_{k=1}^{\nu} \omega_k (\mathscr{G}_{f_{i,j}}^{(k)} - f_{i,j}^{(k)}) \\
f_{i,j}^{(k)} &= f_{i,j}^n - \Delta t \sum_{l=1}^{k-1} \tilde{A}_{k,l} \xi \nabla_x f_{i,j}^{(l)} + \frac{\Delta t}{\tau} \sum_{l=1}^{k} A_{k,l} (\mathscr{G}_{f_{i,j}}^{(l)} - f_{i,j}^{(l)}) \\
f_{i,j}^{(1)} &= f_{i,j}^n + \frac{\Delta t}{\tau} A_{1,1} (\mathscr{G}_{f_{i,j}}^{(1)} - f_{i,j}^{(1)})
\end{aligned}
\tag{9}
$$

where $A$ and $\tilde{A}$ are $\nu \times \nu$ matrices, with $\tilde{A}_{i,s} = 0$ if $s \geq i$ and $A_{i,s} = 0$ if $s > i$. These coefficients are derived from a double Butcher's tableaux. Here we take the second-order scheme of [11].

## 4 The Level Set Function

When an immersed solid is considered on a Cartesian grid, one needs to apply the wall boundary condition on a surface that is arbitrarily crossing the grid. To this end, the domain is decomposed in a fluid part and a solid part. In the solid the values of the physical variables are imposed in each cell since there is no calculation to perform. Such cells are called penalized cells. To decide whether or not a cell is penalized on a Cartesian mesh and to improve accuracy at the boundaries, we use the signed distance between a grid point and the immersed body. Introduced by Osher and Sethian [10], the level set function implicitly defines the solid interface $\Sigma$ in the computational domain by its zero isoline. It is defined by:

$$
\phi(x) = \begin{cases} dist_\Sigma(x) & \text{outside the solid} \\ -dist_\Sigma(x) & \text{inside the solid} \end{cases} \quad \text{and} \quad n(x) = \frac{\nabla \phi(x)}{|\nabla \phi(x)|} \tag{10}
$$

where $dist_\Sigma(x)$ is the minimum distance between the point considered (with coordinates $x$) and the solid interface $\Sigma$, $n(x)$ is the unit normal to the distance isoline pointing towards the fluid. In particular, for $\phi = 0$, $n_w(x)$ is the normal to the boundary pointing towards the fluid.

In the following test case of the nozzle plume, the contour of the jet is modeled as a moving boundary. The level set function defining such boundary is convected

with the imposed boundary velocity $u_\phi$:

$$\partial_t \phi + u_\phi \cdot \nabla \phi = 0 \tag{11}$$

This equation is solved with a WENO5 [8] discretization scheme in space and a standard Runge-Kutta 4 scheme for the integration in time.

Integrating (11) in time does not preserve the distance property of $\phi$. Therefore, a reinitialisation step is performed after each time integration step starting from the boundary ($\phi = 0$). In our case, this is done via a Fast Marching algorithm [12].

## 5 Wall Conditions

Two kinds of boundary conditions for kinetic models are usually found in the literature: the diffuse boundary condition and the specular reflection. In the following we consider only the specular reflection.

In the classical specular reflection, each particle hitting the wall is immediately reflected by the wall with the same tangential velocity and the opposite normal velocity: $\xi_{refl} = \xi - 2((\xi - U_w) \cdot n_w)n_w$, with $\xi_{refl}$ the particle velocity after reflection, $\xi$ the particle velocity before reflection, $U_w$ the wall velocity and $n_w$ the normal to the wall pointing towards the fluid. This holds true for each particle such that $\xi \cdot n_w > 0$. For $\xi \cdot n_w < 0$, the distribution function on the boundary is already known and equal to the one reconstructed at the boundary. The entire distribution function $f_s$ enforcing the boundary condition is then:

$$f_s = \begin{cases} f & \text{for } \xi \cdot n_w < 0 \\ f(\xi_{refl}) & \text{for } \xi \cdot n_w > 0 \end{cases} \tag{12}$$

However, because of the discretization of the velocity space, one needs to compute $f(\xi_{refl})$ where in general $\xi_{refl}$ does not correspond to a collocation point. Therefore $\xi_{refl}$ must be interpolated. This creates spurious mass and energy fluxes due to interpolation errors that can only be removed at significant cost (finer grid or higher-order interpolations), see [3].

To handle this problem, in the hydrodynamic limit, we have devised a new Euler-AP boundary condition. Let us assume that, in the fluid, the distribution function is a Gaussian ($Kn$ number close to 0). Then, imposing the impermeability condition at the wall corresponds to impose a Gaussian distribution function. However, in this case tangential velocity and temperature should be taken from the fluid and the velocity must have a zero wall-normal component. Therefore, tangential velocity and temperature are extrapolated from the fluid to the wall at the desired order. The density is computed invoking mass conservation. Finally, the same procedure as in (12) is applied to obtain the boundary condition $f_M$.

To build a fully asymptotic preserving boundary condition valid in more rarefied regimes, this model is added with the classical specular reflection with a coefficient

$\beta$: $f_b = \beta f_s + (1 - \beta) f_M$ with $\beta \in [0, 1]$ and such that it is close to zero in the inviscid limit for $Kn \to 0$. If $Kn$ is not close to zero, the classic specular reflection correctly takes over.

To set the value of $\beta$ we emphasize that $f_M$ corresponds to the specular reflection only when the distribution function in the fluid is close to the Gaussian. If it is not the case, the specular reflection is computed with $f_s$ (in particular in the rarefied regime). In our model, $\beta$ is set as follows:

$$\beta = \min\left(1, \frac{||f - \mathscr{G}_f||_{L^2}}{\max(f)tol}\right) \qquad (13)$$

with $tol$ a tolerance on the distance in $L^2$ norm between the closest interior domain distribution function $f$ and its corresponding Gaussian. Thus, if $\dfrac{||f - \mathscr{G}_f||_{L^2}}{\max(f)} <<$ $tol$, the specular reflection fully corresponds to the Euler-AP boundary condition. In the following, $tol = 10^{-3}$.
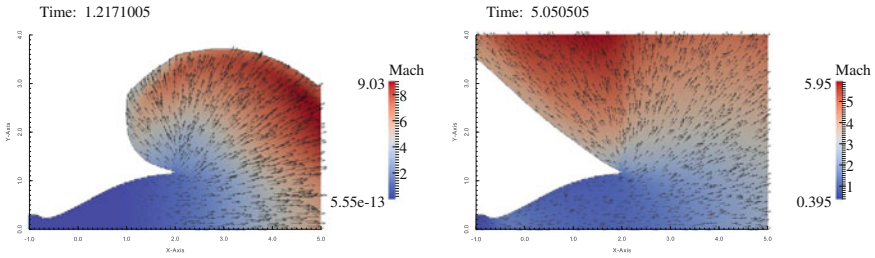
## 6 Numerical Example: Unsteady Nozzle Plume

We consider a qualitative validation of this numerical model against experimental data. A flow expands at the outlet of a nozzle in a low pressure atmosphere. Experiments by Latvala et al. in [9] were performed for different ambient pressures to determine the angle of the jet at the outlet of the nozzle. The area ratio between the throat and the outlet of the nozzle is 4.8. Here, we impose the total pressure ($P_{tot} = 1$) and the total temperature ($T_{tot} = 0.6$) at the inlet of the nozzle. The 1D isentropic flow formula gives M = 3.7763, T = 0.1738 and p = 0.0126 at the outlet. This pressure is called the adaptation pressure $P_c$.

We obtain the jet angle by tracking the contact discontinuity between the gas coming from the nozzle and the gas initially outside the nozzle with a level set function keeping fixed the point at the extremity of the nozzle ($x = 2$). At each time step, the level set function is transported according to the velocity of the fluid with the procedure described in Sect. 4 for moving boundaries. The velocity of the contact discontinuity is computed thanks to a Riemann problem where only the ambient pressure $P_{atm}$ is imposed.

At the initial state, the nozzle is filled with a gas at rest with $p = 1$ and $T = 0.6$. Outside of the nozzle, the gas is also at rest, with $p = P_{atm}$ and $T = 0.6$.

The solution for $P_c/P_{atm} = 2000$ is shown on Figs. 1 and 2 for different times. First, the flow goes out of the nozzle and turns back because of the abrupt expansion ($t = 1.2$ and $t = 5$). Then, when the flow stabilizes in the nozzle, a shock propagates from the inlet towards the outlet. On Fig. 2, at $t = 11.1$ the shock is at $x = 2.7$ and establishes the angle of the jet with the nozzle outlet.

Figure 3 shows the angle $\delta$ of the jet at the outlet for different pressure ratio $P_c/P_{atm}$, for ES-BGK and compressible Euler models. These results can be

**Fig. 1** Mach number and velocity vectors at t = 1.2 (*left*) and t = 5 (*right*) for $P_c/P_{atm} = 2000$



**Fig. 2** Mach number and velocity vectors at t = 11 (*left*) and steady state (*right*) for $P_c/P_{atm} = 2000$



**Fig. 3** Angle for different pressure ratios calculated with Euler and ES-BGK models

qualitatively compared with the experimental results obtained by Latvala et al. in [9] where it is shown, for $\gamma = 7/5$, that the evolution of the jet angle is linear with the logarithm of the pressure ratio. The same behaviour is observed in Fig. 3 for the ES-BGK model. The quantitative results cannot be directly compared to experiments since within the limit of our model, $\gamma = 5/3$. In order to consider $\gamma = 7/5$, one should include additional terms in the model as done in [1].

For small pressure ratios ($<$10) the ES-BGK and compressible Euler models give the same angle. When the ratio increases, the difference becomes larger and the kinetic model stays very close to a straight line. For this kind of pressure ratio, the local relaxation time increases outside the nozzle and becomes too large to consider the fluid at equilibrium. Thus, the continuum model tends to give a different solution. This emphasizes the necessity of using a kinetic model with an AP boundary condition since this problem cannot be solved with a continuum model for high pressure ratio. Also, a solution computed with a standard specular reflection wall condition in the nozzle would significantly pollute the simulation [3].

## 7 Conclusion

We have presented an integration scheme for the ES-BGK model discretized on Cartesian meshes. The scheme is second-order accurate in space and time. The wall condition is such that the hydrodynamic limit is preserved without spurious effects due to the discretization. The scheme is accurate and yet easily implemented in actual applications as shown by the nozzle plume case presented.

## References

1. Andries, P., Le Tallec, P., Perlat, J.P., Perthame, B.: The Gaussian-BGK model of Boltzmann equation with small Prandtl number. Eur. J. Mech. B. Fluids **19**(6), 813–830 (2000)
2. Ascher, U.M., Ruuth, S.J., Spiteri, R.J.: Implicit-explicit Runge-Kutta methods for time-dependent partial differential equations. IMACS J. Appl. Numer. Math. **25**(2–3), 151–167 (1997)
3. Bernard, F., Iollo, A., Puppo, G.: Accurate Asymptotic Preserving Boundary Conditions for Kinetic Equations on Cartesian Grids. Rapport de recherche RR-8471, INRIA (2014)
4. Bhatnagar, P.L., Gross, E.P., Krook, M.: A model for collision processes in gases. I. Small amplitude processes in charged and neutral one- component systems. Phys. Rev. **94**, 511–525 (1954)
5. Bird, G.A.: Molecular Gas Dynamics and the Direct Simulation of Gas Flows. Clarendon Press, Oxford Engineering Science Series, Oxford (1994)
6. Filbet, F., Jin, S.: An asymptotic preserving scheme for the ES-BGK model of the Boltzmann equation. J. Sci. Comput. **46**(2), 204–224 (2010)
7. Holway L.H. Jr.: Kinetic theory of Shock structure using an Ellipsoidal distribution Function. In: J.H. de Leeuw (ed.) Rarefied Gas Dynamics, vol. 1, p. 193. Academic, New York (1965)
8. Jiang, G.S., Shu, C.W.: Efficient implementation of weighted ENO schemes. J. Comput. Phys. **126**(1), 202–228 (1996)
9. Latvala, E.K., Anderson, T.P.: Experimental Determination of Jet Spreading from Supersonic Nozzles at High Altitudes. AEDC-TN-58-98 (ASTIA Document No. AD-208546) (1959)
10. Osher, S., Sethian, J.A.: Fronts propagating with curvature-dependent speed : algorithms based on hamilton jacobi formulations. J. Comput. Phys. **79**(12) (1988)
11. Pieraccini, S., Puppo, G.: Implicit-explicit schemes for BGK kinetic equations. J. Sci. Comput. **32**(1), 1–28 (2007)
12. Sethian, J.A.: Fast marching methods. SIAM Rev. **41**, 199–235 (1999)

# On the Use of the HLL-Scheme
# for the Simulation of the Multi-Species
# Euler Equations

**Phillip Berndt**

**Abstract** The HLL approximate Riemann solver is a reliable, fast and easy to implement tool for the under-resolved computation of inviscid flows. When applied to multi-species flows, it generates pressure oscillations at material interfaces. This is a well-known behaviour of conservative solvers and has been addressed as a problem by several authors before. We show that for this particular solver, the generation of pressure oscillations can be desired and is consistent with the underlying physics.

## 1 Introduction

The HLL solver, proposed by Harten, Lax and van Leer in [4] and later enhanced by Einfeldt with proper signal velocity estimates in [2], is based on the reduction of the Riemann problem to its two dominating waves. By demanding that the conservation law is fulfiled, the single intermediate state in between the waves is defined to be

$$q_\star = \frac{f(q_l) - f(q_r) - s_l q_l + s_r q_r}{s_r - s_l},\tag{1}$$

where $q_l$ and $q_r$ denote the left and right state, $s_1$ and $s_2$ are the dominating signal speeds and $f$ denotes the differential flux function of any one dimensional hyperbolic equation of the form $q_t + f(q)_x = 0$. A numerical flux can then be defined by either calculating a finite volume's average at time $\Delta t$, as in the classical Godunov scheme, or by another application of the conservation law. In both cases, the numerical flux is found to be the upwind flux if either $s_r < 0$ or $s_l > 0$, and

P. Berndt (✉)
Freie Universität Berlin, Department of Mathematics and Computer Science,
Arnimallee 6, 14195 Berlin, Germany
e-mail: pberndt@math.fu-berlin.de

$$F(q_l, q_r) = \frac{s_r f(q_l) - s_l f(q_r)}{s_r - s_l} + \frac{s_r s_l}{s_r - s_l} (q_r - q_l) \tag{2}$$

for $s_l < 0 < s_r$. We will for $f$ only consider the Euler equations, augmented with advection equations for the mass fractions ($\rho Y_i$) of the individual species and a thermally perfect[1] equation of state.

When these equations were first discussed, it was soon found that this extension of the Euler equations introduces some new numerical issues. Maintaining positivity, which had been a problem for density and pressure before, now also became relevant for the individual species' masses. This problem has been successfully addressed by Larrouturou [5], who suggested to define the species' mass density fluxes as the product of the density flux and the upwind mass fractions. A second problem can be readily observed when one tries to simulate two separated species in pressure equilibrium, moving with a common velocity. If the temperatures and gas constants differ, conservative simulations create pressure oscillations at the interface, which are not predicted by the equations. Several authors approached this problem, most notably R. Abgrall and S. Karni. Their paper [1] gives an overview on the different approaches in 2000. They all amount to the application of some non-conservative correction. To our knowledge, no satisfactory solution has been proposed since.

In the following, we will discuss both issues for the HLL solver and thermally perfect gas mixtures. We will show that a Larrouturou-type correction is not required to preserve positivity and that this correction would in fact increase the second problem. We will then argue that in the light of the HLL-scheme, it is an option not to apply any correction to stop the pressure oscillations. To conclude, we will elaborate on a test case demonstrating the downside of such corrections.

## 2 Positivity of the Mass Fractions

To show that the scheme preserves mass fraction positivity, it suffices to show that the intermediate states $q_\star$ preserve it. The remainder of the scheme boils down to updating every cell with a convex combination of its old time level's state and the surrounding $q_\star$'s.

Einfeldt defined the two signal speeds to be

$$s_1 = \min\{u_l - c_l, \bar{u} - \bar{c}\} \tag{3}$$

$$s_2 = \max\{u_r + c_r, \bar{u} + \bar{c}\}, \tag{4}$$

where $u$ denotes the velocity, $c$ the speed of sound and $\bar{u}, \bar{c}$ Roe-averaged velocity and sound speed. Note that $\bar{c}$ is always positive. (See e.g. [2] or [6].) For $s_1 > 0 \vee s_2 < 0$, the flux simplifies to upwinding, which obviously preserves mass fraction positivity

---

[1] i.e. $e(T, Y) = \sum_i \int^T Y_i c_{v,i}(T) \, d\tau$, $p = \rho R T$. Some authors call this "ideal".

if the CFL condition is fulfilled. In the other case, the mass fractions of $q_\star$ are, according to Eq. (1),

$$(\rho Y)_\star = \frac{u_l - s_l}{s_r - s_l}(\rho Y)_l + \frac{s_r - u_r}{s_r - s_l}(\rho Y)_r.$$ (5)

The denominator is always positive because of the constraint that $s_l < 0 < s_r$. We have to show that the numerators are also positive. For the first factor, Eq. (3) gives us two cases: If $u_l - c_l$ is smaller than $\bar{u} - \bar{c}$, the factor is trivially positive. In the other case,

$$u_l - s_l = u_l - (\bar{u} - \bar{c}) > (u_l - c_l) - (\bar{u} - \bar{c}) > 0.$$ (6)

The argument for $s_r - u_r$ is analogous. Since $(\rho Y)_\star$ is the sum of two positive mass fraction vectors, weighted by positive factors, it must remain positive. This concludes the proof.

We showed that Larrouturou's correction is not required to preserve mass fraction positivity in the HLL solver. Now we proceed to show that applying it would actually worsen the second problem. To this end, consider the following Riemann problem: Two species are initially separated and in pressure *and* temperature equilibrium. For now, let them be at rest. This problem is in equilibrium, so nothing should happen. The first term in the HLL flux from Eq. (2) vanishes, but the second one is a diffusion term, which will smear the discontinuity. The according update step amounts to a convex combination of conserved quantities. It is straightforward to show that the following statements hold for thermally perfect gas mixtures:

**Theorem 1** *Let $(\alpha_i)$ label the coefficients of a linear combination of states represented by the conserved quantities. Furthermore assume that the states have a common velocity.*

1. *The combined state preserves temperature equilibria*
2. *For initially equal pressures $p_0$, the combined state has a pressure of*

$$p = \sum_i \alpha_i \rho_i R_i T = p_0 \sum_i \alpha_i \frac{T}{T_i}.$$ (7)

*In particular, if the states are in temperature equilibrium and the combination is convex, the pressure is maintained.*

See below for a proof. By the theorem, the smearing of the discontinuity introduced by the HLL solver does not disturb the pressure and temperature equilibria. If one applies the Larrouturou correction to the flux, this changes: What has been a purely diffusive term $((\rho Y)_l - (\rho Y)_r)$ in the flux formula before, now becomes

$$\begin{cases} (\rho Y)_l & \text{if } \rho_l > \rho_r, \\ (\rho Y)_r & \text{else,} \end{cases}$$ (8)

so that the update no longer amounts to a linear combination and therefore introduces a pressure oscillation, which in the next time step also disturbs the temperature equilibrium. For non-vanishing velocities, the correction to the flux's advective terms introduces the same problem, so applying the correction solely to them would not resolve the problem.

While we only carried out this computation for the HLL scheme, other schemes which do not necessarily advect the species' mass fractions with the correct velocity might be affected as well. The two goals of maintaining positivity and preventing pressure oscillations can therefore not generally be treated independently.

*Proof (of Theorem 1)* We denote the different state variables going into the linear combination with subscript indexes and use superscript indexes for the different mass fractions $Y^j$. For the final state, we omit the subscript. First, assume that $T_i = T_0$ for all states $i$. For the internal energy

$$\rho e = \rho \int^T c_v(\tau) \, d\tau = \rho \sum_j Y^j \int^T c_v^j(\tau) \, d\tau = \sum_i \alpha_i \rho_i \int^T c_{v,i}(\tau) \, d\tau \quad (9)$$

holds. Note that the upper integral bound is the final temperature $T$ and not the states' temperature $T_i$. Since we do not consider different velocities, the internal energy is a conserved quantity, so

$$\sum_i \alpha_i \rho_i \int^T c_{v,i}(\tau) \, d\tau = \sum_i \alpha_i \rho_i e_i \quad (10)$$

must hold for any choice of $\alpha_i$. Consequently, the integrals must be equal to the corresponding energies. Finally, since energy is an injective function of $T$, $T = T_0$. For the pressure relations, assume $p_i = p_0$ for all $i$. By applying the ideal gas law, $p = \rho R T$,

$$p = \left( \sum_i \alpha_i \rho_i \right) \left( \hat{R} \sum_j \frac{1}{m^j} \frac{\sum_i \alpha_i \rho_i Y_i^j}{\sum_i \alpha_i \rho_i} \right) T = \left( \hat{R} \sum_j \frac{1}{m^j} \sum_i \alpha_i \rho_i Y_i^j \right) T$$

$$= \left( \hat{R} \sum_i \sum_j \frac{\alpha_i \rho_i X_i^j}{m_i} \right) T = \sum_i \alpha_i \rho_i R_i T = \sum_i \alpha_i p_i \frac{T}{T_i}. \quad (11)$$

In the last equation, $\hat{R}$ denotes the universal gas constant, $X^j$ the mole fractions, $m^j$ the species' molar masses, and $m_i$ the mean molar mass of the $i$th state. The expression is already in the form stated in the theorem. If $T_i = T$ for all $i$ the fraction cancels and we establish our final claim. $\qquad \square$

## 3 Pressure Oscillations at Species Interfaces

As sketched above, we prototypically investigate the simulation of two separated gases of different temperatures in pressure equilibrium, moving at a constant velocity. The Euler equations predict for this situation that the material interface moves at the given velocity. Godunov-type schemes fail to simulate them correctly for CFL numbers other than one, because the solution inside the finite volume containing the contact wave is projected to a single constant state. As has been stated in Theorem 1, this projected state does not generally have the same pressure the original states had. It is an interesting observation that this result is physically sound.

To illustrate this, assume in a gedankenexperiment that one could physically perform the Godunov method. The only requirement for this is that one must be able to instantaneously place and remove isolating walls in/from the domain of interest. By switching coordinate systems we can view the interface as being at rest and the position of the walls to change with each time step. If one removes those walls, the evolution step of the scheme takes place. After the walls are replaced, thermal and species diffusion equilibrate the volumes, corresponding to the scheme's projection step. The only difference to a simulation is that the dynamics is not governed by the Euler equations, but by real physics. The equilibrating is what concerns us here. From the solution of the Riemann problem, the content of a cell at the start of the equilibration is known: The material wave has progressed a fraction of $\alpha$ into the cell. On either side are the two gases, which we denote by subscripts 1 and 2. Thermodynamics allows us to split the equilibration into two processes: First, each gas isothermally expands to the whole finite volume. By Boyle's law, the final pressures are $p_1 = p_0 \alpha$ and $p_2 = p_0(1 - \alpha)$. The sum of the partial pressures equals the original pressure $p_0$, as one would expect from Dalton's law. In the second step, the temperatures isochorically equilibrate. By the ideal gas law we expect $\Delta p_i = \rho_i R_i \Delta T_i$. The pressure remains constant if and only if

$$\rho_1 R_1 \Delta T_1 = -\rho_2 R_2 \Delta T_2. \tag{12}$$

For the equilibration, on the other hand,

$$\rho_1 c_{v,1}(T_1)\mathrm{d}T_1 = -\rho_2 c_{v,2}(T_2)\mathrm{d}T_2 \tag{13}$$

holds, following the definition of a thermally perfect gas. The pressure is therefore generally not maintained, but will change due to thermal equilibration. A calculation using the convex combination of states from Theorem 1 leads both qualitatively and quantitatively to the same result. We showed that we are in the situation that the scheme exhibits an effect which is physically sound but does not comply with the model equations. In other words:

*A physical phenomenon which is not contained in the model equations was used to discretize said equations.*

**Table 1** Strength of pressure oscillations for different discontinuity widths

| Initial width of the interface (cells) | Max. pressure deviation (Pa) | Relative error |
|---|---|---|
| 1 | 4306.777902 | $4.306778 \cdot 10^{-02}$ |
| 10 | 628.360586 | $6.283606 \cdot 10^{-03}$ |
| 20 | 172.625298 | $1.726253 \cdot 10^{-03}$ |
| 50 | 20.347418 | $2.034742 \cdot 10^{-04}$ |
| 100 | 4.701430 | $4.701430 \cdot 10^{-05}$ |
| 150 | 1.673772 | $1.673772 \cdot 10^{-05}$ |

*Initial values* Riemann problem with Ar at 1200 K and $N_2$ at 300 K, both at 1 bar and at rest.

The obvious analogy to the role of diffusion for the traditional equations raises the question whether one should always counteract the pressure oscillations: With the traditional equations, it was found that the model's lack of diffusion is a source for well-posedness problems and that the qualitative insertion of the missing effect is required to ensure that the scheme converges to the physically relevant solution. (e.g. [7].) On the other hand, diffusion has to be kept as small as possible to comply with the model equations which do not predict it.

The HLL-solver can be interpreted as a scheme embracing one of the two possible extremes: It discards the whole jump discontinuity as a model phenomenon, prevented by diffusion from existing in real systems, and replaces it by a smeared state to introduce the maximal physically reasonable amount of diffusion.[2] If one takes this position, the pressure oscillations turn into an effect of a desired solver property and it becomes sensible not to counteract them at all.

We do not expect long-time contamination of the pressure field from the lack of a correction, because the effect is self-weakening: The scheme's diffusion is what leads to the oscillations, but it also smoothenes the contact discontinuity. As it smoothenes, adjacent gas constants become increasingly similar. By expanding Eq. (7) in $R$, we see that

$$p = p_0 + \mathcal{O}(\sum_i (R_i - R)). \tag{14}$$

To quantify this relation, see the exemplary calculation in Table 1. A consequence of this observation is that schemes which try to maintain the contact discontinuity (but do not track it) will introduce more pressure oscillations compared to those which quickly dissolve it. We would like to again emphasize that these thoughts do only make sense in the context of miscible species and where Dalton's law holds. Also, the statement that small pressure oscillations do not contaminate the field for long

---

[2] This thought led to the HLLE scheme, where Einfeldt reintroduced the contact wave as a linear growth rather than a sharp discontinuity into the Riemann fan. [3]

**Table 2** Riemann test problems for species diffusion

| Variable | Left IV | Right IV | Equilibrium (actual) | Equilibrium[a] (HLL-scheme) |
|---|---|---|---|---|
| $p$ | 1 bar | 1 bar | 0.9994 bar | 0.9994 bar |
| $T$ | 300 K | 800 K | 436.11 K | 436.11 K |
| $X$ | C | Ar | — | — |
| $p$ | 1 bar | 1 bar | 0.9227 bar | 0.9220 bar |
| $T$ | 300 K | 600 K | 369.08 K | 368.99 K |
| $X$ | $N_2$ | Ar | — | — |
| $p$ | 1 bar | 1 bar | 0.8423 bar | 0.8408 bar |
| $T$ | 300 K | 2000 K | 439.47 K | 439.30 K |
| $X$ | $N_2$ | Ar | — | — |
| $p$ | 1 bar | 1 bar | 0.8207 bar | 0.8191 bar |
| $T$ | 300 K | 4000 K | 458.09 K | 457.92 K |
| $X$ | $N_2$ | Ar | — | — |

[a] The simulations were carried out with 10 cells per species at cfl $= 0.5$ and stopped when the mass fractions were close to being homogeneous at $t = 0.298$ s.

times is only valid if a simulation is sufficiently under-resolved for highly active source terms (e.g. kinetics) to equilibrate before any incoming pressure waves can fasten the reaction.

## 4 Species Diffusion Test Case

The Riemann problem with which we began our investigation can be used as a test case: Two initially separated species at different temperatures are set up in a confined volume with impenetrable, isolating, reflecting boundaries. Due to species diffusion, this experiment will physically eventually reach the equilibrium state we derived above. We enforce numerical diffusion by overlaying the initial values with a small acoustic field at one of the domain's resonance frequencies, such that the field does not alter the average pressure. The mean pressure and temperature of any *physically* consistent scheme should reach said equilibrium values.

We have employed the HLL solver in a MUSCL-scheme with a second-order reconstruction in the primitive variables pressure, temperature, velocity and the mass fractions. The minmod-Limiter was used. For handling the multi-species EOS, an in-house chemical kinetics library using the GRI 3.0 mechanism was included with strang-splitting. The results for this scheme are given in Table 2.

The numerical equilibrium values are close to the predictions, which is the expected behaviour. A pressure-correcting scheme, in contrast, cannot be expected to reproduce the correct result, because it is designed to maintain the original pressure. Since it would still diffuse the species, its result would be neither physically correct nor in agreement with the Euler equations.

# 5 Conclusions

We have introduced a different view on the multi-species problems for the special case of thermally perfect gas mixtures and suggested that it might not always be desired to correct the pressure oscillations. We proved that the first-order HLL scheme preserves mass fraction positivity. This result is still incomplete: It remains to be shown whether the scheme is also monotonicity preserving and whether the second-order extension we have employed in the test section requires special treatment to preserve the positivity. The latter is not obvious for the forward-in-time reconstruction, since the mass fractions enter as $\mathrm{d}Y/\mathrm{d}t = -u\,\mathrm{d}Y/\mathrm{d}x$. We plan to accordingly extend the result in a future contribution.

# References

1. Abgrall, R., Karni, S.: Computations of compressible multifluids. J. Comput. Phys. **169**, 594–623 (2001)
2. Einfeldt, B.: On godunov-type methods for gas dynamics. SIAM J. Numer. Anal. **25**(2), 294–318 (1988)
3. Einfeldt, B.: A intuitionistic approach to the fluid continuum. http://discontinuous-flow.blogspot.de/2012/08/a-intuitionistic-approach-to-fluid.html (2012) Accessed 7 Feb 2014
4. Harten, A., Lax, P., van Leer, B.: On upstream differencing and godunov-type schemes for hyperbolic conservation laws. SIAM Rev. **25**(1), 35–61 (1983)
5. Larrouturou, B.: How to preserve the mass fractions positivity when computing compressible multi-component flows. J. Comput. Phys. **95**, 59–84 (1991)
6. Larrouturou, B., Fezoui, L.: On the equations of multi-component perfect or real gas inviscid flow. Nonlinear Hyperbolic Problems, Lecture Notes in Math. **1402**, 69–98 (1989)
7. Tadmor, E.: Numerical viscosity and the entropy condition for conservative difference schemes. Math. Comput. **43**(168), 369–381 (1984)

# A Conservative Well-Balanced Hybrid SPH Scheme for the Shallow-Water Model

**Christophe Berthon, Matthieu de Leffe and Victor Michel-Dansac**

**Abstract** A scheme defined by a hybridization between SPH method and finite volume method is considered. The aim of the present communication is to derive a suitable discretization of the source term to enforce the required well-balanced property. To address such an issue, we adopt a relevant reformulation of the flux function by involving the free surface instead of the water height. Such an approach gives a natural discretization of the topography source term in order to preserve the lake at rest. Moreover, we prove that the scheme is in conservative form, which is, in general, a very difficult task since we do not impose restrictive assumptions on the SPH method. Several 1D numerical experiments are performed to exhibit the properties of the scheme.

## 1 Introduction

The present work concerns the numerical approximation of the well-known shallow-water model. The model under consideration is given as follows:

$$
\begin{cases}
\partial_t h + \partial_x(hu) = 0, \\
\partial_t(hu) + \partial_x(hu^2 + \frac{g}{2}h^2) = -hg\,\partial_x Z,
\end{cases}
\tag{1}
$$

C. Berthon · V. Michel-Dansac (✉)
Laboratoire de Mathématiques Jean Leray, UMR 6629, 2 rue de la Houssinière,
BP 92208, 44322 Nantes Cedex 3, France
e-mail: victor.michel-dansac@univ-nantes.fr

C. Berthon
e-mail: christophe.berthon@univ-nantes.fr

M. de Leffe
HydrOcean, 1 rue de la Noë, CS 32122, 44321 Nantes Cedex 3, France
e-mail: matthieu.de-leffe@hydrocean.fr

where $h \geq 0$ denotes the water height, $u \in \mathbb{R}$ is the water velocity in the $x$ direction, and $g > 0$ stands for the gravity constant. The function $Z$ denotes the smooth topography. To shorten the notations, the system is rewritten in the following form:

$$\partial_t \Phi + \partial_x f(\Phi) = S, \quad \Phi = \begin{pmatrix} h \\ hu \end{pmatrix}, \quad f(\Phi) = \begin{pmatrix} hu \\ hu^2 + \frac{g}{2}h^2 \end{pmatrix}, \quad S = \begin{pmatrix} 0 \\ -hg\partial_x Z \end{pmatrix}. \tag{2}$$

By adopting a finite volume method to approximate the weak solutions, a usual property to be satisfied concerns the lake at rest preservation. Indeed, the stationary solution given by $u = 0$ and $h + Z = $ cst, must be exactly preserved by the numerical method (for instance, see [2, 3, 5] and references therein).

Here, we do not consider a classic finite volume scheme, but we adopt a hybrid method deriving from the SPH techniques. More precisely, the SPH method (issuing from the particle methods) involves a like interface numerical flux function. According to [7, 14], this like interface flux function is substituted by a like finite volume flux function derived from approximate Riemann solvers [5, 9, 12].

In the present paper, we exhibit a source term discretization to make well-balanced this hybrid numerical technique. The paper is organized as follows. In the next section, we briefly recall the gradient evaluation derived from the SPH technique, and the hybrid version by considering approximate Riemann solvers. Next, in Sect. 3, after [4], we adopt a relevant reformulation of the model to introduce a suitable well-balanced discretization of the topography source term. The full discrete scheme is proved to preserve the required lake at rest, and it is in conservation form without any additional assumptions. In Sect. 4, numerical experiments are performed in order to illustrate the relevance of the scheme. A short conclusion is given in the last section.

## 2 Introduction to the SPH Method and Finite Volume Hybridization

The Smoothed Particle Hydrodynamics (SPH) method was introduced to perform astrophysical simulations. Recent works (for instance see [13] and references therein) extend the SPH method in the field of CFD. Now, we present the derivation of the SPH scheme to approximate the weak solutions of (2).

First, it is worth noticing that, for all real functions $f : \mathbb{R} \to \mathbb{R}$, we have the following relation: $f(x) = (f * \delta)(x) = \int_{\mathbb{R}} f(y)\delta(x - y)dy$, with $\delta$ the Dirac measure. The particle approximation relies on a suitable regularization of this Dirac measure.

To address such as issue, after [10, 11], a kernel $W \in C_0^1(\mathbb{R}) \cap L^1(\mathbb{R})$ is introduced, which is usually some bell-shaped function, depending on both center $x$ and smoothing length $h$. It must satisfy the consistency conditions [11] given by $\int_{\mathbb{R}} W(x, h)dx = 1$ and $\int_{\mathbb{R}} W'(x, h)dx = 0$.

Now, after [10, 11], the particle approximation of $f$, given by $(f * W)(x) = \int_{\mathbb{R}} f(y)W(x - y, h)dy$, is nothing but a second-order accurate approximation of the function $f$. Using the Green formula, we easily deduce an approximation of $f'$, given by $\int_{\mathbb{R}} f(y)W'(x - y, h)dy$.

Unfortunately, this particle approximation involves integrals which cannot be exactly evaluated. As a consequence, a quadrature formula is adopted as follows: $\int_{\mathbb{R}} f(x)dx \simeq \sum_{j\in\mathscr{P}} \omega_j f_j$, where $x_j$ are the quadrature points, $f_j$ denotes the evaluation of $f$ at point $x_j$, and $\omega_j$ stands for the associated weight. Within the SPH method, the quadrature points are made of particles $x_i$ with volume $\omega_i$, and $\mathscr{P}$ denotes the set of interacting particles $x_j$ close enough to the particle $x_i$. We then get the following approximation:

$$\Pi^h(f)_i = \sum_{j\in\mathscr{P}} \omega_j f_j W_{ij}, \quad \Pi^h(f')_i = \sum_{j\in\mathscr{P}} \omega_j f_j W'_{ij}, \quad W_{ij} = W(x_i - x_j, h).$$
(3)

From now on, let us underline that a natural property to be satisfied by this particle approximation is $\Pi^h(1) = 1$ and $\Pi^h(1') = 0$, which reads

$$\sum_{j\in\mathscr{P}} \omega_j W_{ij} = 1 \quad \text{and} \quad \sum_{j\in\mathscr{P}} \omega_j W'_{ij} = 0.$$
(4)

Such relations are not always satisfied by usual choices for the kernel $W$ (see [13]).

By adopting the derivative discretization formula (3), the SPH scheme for a general set of equations (2) is given by

$$\frac{1}{\Delta t}\omega_i\left(\Phi_i^{n+1} - \Phi_i^n\right) + \sum_{j\in\mathscr{P}} \omega_i\omega_j(f_i^n + f_j^n)W'_{ij} = \omega_i S_i^n,$$

with $f_i^n = f(\Phi_i^n)$, $\Delta t$ the time step, and $\omega_i\Phi_i^n$ the vector of conserved variables for the particle $x_i$.

Concerning the source term discretization $S_i^n$, one may adopt the particle approximation (3). However, in order to satisfy the required well-balanced property, we will introduce a specific approximation of the topography in the next section.

To conclude this brief presentation of the SPH scheme, we now show the finite volume approximate Riemann solver hybridization as introduced in [14]. Indeed, the derivative discrete operator involves an interface flux approximation given by $\frac{1}{2}(f_i^n + f_j^n)$. In [14], this flux approximation is substituted by the numerical flux function coming from usual Godunov-type scheme (for instance Godunov, HLL, HLLC, Roe scheme [8]). As a consequence, we consider the following modified hybrid SPH scheme:

$$
\begin{cases}
\dfrac{1}{\Delta t}\omega_i \left(h_i^{n+1} - h_i^n\right) + \displaystyle\sum_{j\in\mathscr{P}} 2\omega_i\omega_j \, (hu)_{ij} \, W'_{ij} = 0, \\[2ex]
\dfrac{1}{\Delta t}\omega_i \left(h_i^{n+1}u_i^{n+1} - h_i^n u_i^n\right) + \displaystyle\sum_{j\in\mathscr{P}} 2\omega_i\omega_j \left(hu^2 + \dfrac{g}{2}h^2\right)_{ij} W'_{ij} = \omega_i S_i,
\end{cases}
\tag{5}
$$

where $f^{\Delta x}(\Phi_i^n, \Phi_j^n) = ((hu)_{ij}, (hu^2 + \frac{g}{2}h^2)_{ij})$ stands for the numerical flux function issuing from a usual finite volume scheme.

## 3 A Well-Balanced Scheme

Neither usual SPH techniques (for instance Monaghan SPH formulation [11], Vila formulation [13]) nor the here presented hybrid schemes combining SPH and Riemann solvers [14] are able to preserve the lake at rest steady state.

   In order to derive a lake at rest preserving scheme, we adopt a recent equivalent reformulation of the PDE. After [7, 14], the flux function, which is in the center of the hybridization, is reformulated by considering the free surface $H = h + Z$ and the velocity. Indeed, within the required lake at rest, these two quantities stay constant, which is of prime importance in the numerical flux definition.

   Here, we assume $H > 0$ and we introduce $X = h/H$ a water height like fraction. To shorten the notations, we set $V = {}^t(H, Hu)$. In the following statement, by considering $V$, we reformulate the system (1) (see [3, 4]).

**Lemma 1** *The weak solutions of* (1) *satisfy the following system:*

$$
\begin{cases}
\partial_t h + \partial_x (X(Hu)) = 0, \\[1ex]
\partial_t (hu) + \partial_x \left(X(Hu^2 + \dfrac{g}{2}H^2)\right) = \dfrac{g}{2}\partial_x(hZ) - gh\partial_x Z.
\end{cases}
\tag{6}
$$

   Let us emphasize that these reformulations involve the flux function but for the new variables $V$. As a consequence, as soon as a lake at rest is considered, this flux function only involves constant states (see (10) later on). This turns out to be the main ingredient to get the required well-balanced property.

   Now, we suggest to adopt the hybrid scheme (5) but for the equivalent formulation (6). As a consequence, the hybrid SPH scheme under consideration now reads

$$
\begin{cases}
\dfrac{1}{\Delta t}\omega_i \left(h_i^{n+1} - h_i^n\right) + \displaystyle\sum_j 2\omega_i\omega_j X_{ij} \, (Hu)_{ij} \, W'_{ij} = 0, \\[2ex]
\dfrac{1}{\Delta t}\omega_i \left(h_i^{n+1}u_i^{n+1} - h_i^n u_i^n\right) + \displaystyle\sum_j 2\omega_i\omega_j X_{ij} \left(Hu^2 + g\dfrac{H^2}{2}\right)_{ij} W'_{ij} \\[2ex]
\qquad\qquad\qquad\qquad\qquad = \omega_i \left(\dfrac{g}{2}\partial_x(hZ) - gh\partial_x Z\right)_i.
\end{cases}
\tag{7}
$$

Concerning the here involved numerical flux function $((Hu)_{ij}, (Hu^2 + gH^2/2)_{ij})$, we directly adopt $f^{\Delta x}(V_i, V_j)$. To complete the scheme, we characterize the new formulation of the source term. Let us first notice the following easy relation:

$$\frac{g}{2}\partial_x(hZ) - gh\partial_x Z = \frac{g}{2}\partial_x\left(H^2 X(1 - X)\right) - gHX\partial_x\left(H(1 - X)\right).$$

In fact, a straightforward application of the SPH discretization is not relevant and we need to consider an additional correction term. We thus adopt the following SPH like discretization of the source term:

$$\omega_i\left(\frac{g}{2}\partial_x(hZ) - gh\partial_x Z\right)_i = \frac{g}{2}\sum_j 2\omega_i\omega_j\left(X_{ij}H_{ij} - 2\bar{H}_i\bar{X}_i\right)\left(H_{ij}\left(1 - X_{ij}\right)\right)W'_{ij}$$

$$+ g\sum_j 2\omega_i\omega_j\bar{H}_i^2\bar{X}_i(1 - \tilde{X}_i)W'_{ij}, \tag{8}$$

where $\bar{H}_i$, $\bar{X}_i$ and $\tilde{X}_i$ are averages to be defined. Let us remark that the correction term $g\sum_j 2\omega_i\omega_j\bar{H}_i^2\bar{X}_i(1 - \tilde{X}_i)W'_{ij}$ vanishes as soon as the kernel function satisfies the consistency conditions (4). This correction term is, in fact, a representation of zero.

Equipped with the hybrid SPH scheme (7)–(8), we now exhibit a suitable definition for $\bar{X}_i$ to enforce the expected well-balanced property.

**Theorem 1** *Assume both free surface averages to satisfy:*

$$H_{ij} = \bar{H}_i = H, \qquad as\ soon\ as\ H_i = H_j = H.$$

*Assume $\bar{X}_i$ is defined by*

$$\bar{X}_i = \frac{1}{2}\frac{\sum_j \omega_j X_{ij}^2 W'_{ij}}{\sum_j \omega_j\left(X_{ij} - 1\right)W'_{ij} + (\tilde{X}_i - 1)\sum_j \omega_j W'_{ij}}. \tag{9}$$

*Then the scheme (7)–(8) preserves the lake at rest.*

*Proof* At time $t^n$, we assume the approximate solution $\Phi_i^n$ be given by the lake at rest. Then, for all $i$ in $\mathbb{Z}$, we have $h_i^n + Z_i = H$ a positive constant and $u_i^n = 0$. The proof consists in establishing $\Phi_i^{n+1} = \Phi_i^n$. Since the numerical flux function $f^{\Delta x}$ is consistent, it preserves the constant states. Hence, we have the following sequence of equalities:

$$f^{\Delta x}(\Phi_i^n, \Phi_{i+1}^n) = f^{\Delta x}\left(\begin{pmatrix} H \\ 0 \end{pmatrix}, \begin{pmatrix} H \\ 0 \end{pmatrix}\right) = f\begin{pmatrix} H \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ g\frac{H^2}{2} \end{pmatrix}. \tag{10}$$

From the water height evolution issuing from (7), we immediately get $h_i^{n+1} = h_i^n$. Next, concerning the discharge evolution, because of the consistency properties of

the involved average functions, $H_{ij}$ and $\bar{H}_i$, the source term discretization (8) now reads

$$\omega_i \left( \frac{g}{2} \partial_x(hZ) - gh\partial_x Z \right)_i = \frac{g}{2} H^2 \sum_j 2\omega_i \omega_j \left( (X_{ij} - 2\bar{X}_i)(1 - X_{ij}) + \bar{X}_i(1 - \tilde{X}_i) \right) W'_{ij}.$$

Finally, by definition of $\bar{X}_i$, given by (9), a straightforward computation gives $\omega_i \left( \frac{g}{2} \partial_x(hZ) - gh\partial_x Z \right)_i = \frac{g}{2} H^2 \sum_j 2\omega_i \omega_j X_{ij} W'_{ij}$. As a consequence, the updated discharge, given by (7), gives $u_i^{n+1} = 0$, and the proof is achieved.                                                                      □

Let us underline that the formula (9), to define $\bar{X}_i$, is consistent with an evaluation of $X$ at particle $x_i$. Indeed, from the SPH space derivative approximation (3), we notice that $\frac{1}{2} \sum_j \omega_j X_{ij}^2 W'_{ij}$ is consistent with $\frac{1}{2} \partial_x X^2$ while $\sum_j \omega_j (X_{ij} - 1) W'_{ij}$ is consistent with $\partial_x(X - 1)$. Since $\sum_j \omega_j W'_{ij}$ is consistent with zero, then $\bar{X}_i$ is consistent with $\frac{1}{2} \partial_x X^2 / \partial_x(X - 1) = X$.
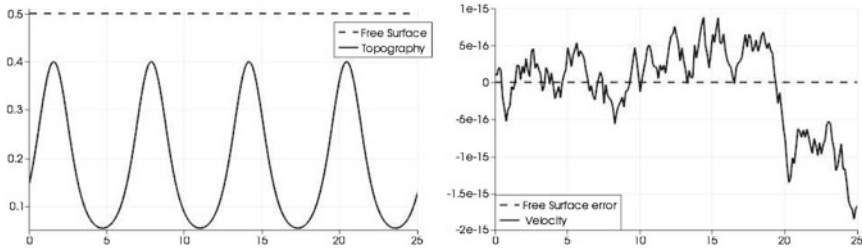
To conclude this section, we remark that the required well-balanced property is established independently of the definitions of $X_{ij}$ and $\bar{H}_i$. Here, we adopt the averages introduced in [3, 4]: $\bar{H}_i = h_i^n + Z_i$ and $X_{ij} = X_i$ if $(Hu)_{ij} > 0$, $X_j$ otherwise.

In fact, at this level, we notice that the proposed scheme satisfies an additional stronger property. Indeed, when adopting SPH type scheme to approximate the solution of homogeneous hyperbolic systems (i.e. with vanishing source term), in general it is not possible to preserve the constant solutions. By considering an initial data made of a uniform constant state, the SPH approach makes some particles move and the constant initial data is no longer preserved. Since the derived scheme is well-balanced, it obviously preserves such constant solutions as soon as the topography is flat, i.e. $Z = \text{cst}$. Moreover, we can exhibit a precise definition of the average functions to preserve the conservation form of the scheme: $\sum_{i \in \mathbb{Z}} \omega_i h_i^{n+1} = \sum_{i \in \mathbb{Z}} \omega_i h_i^n$ and $\sum_{i \in \mathbb{Z}} \omega_i h_i^{n+1} u_i^{n+1} = \sum_{i \in \mathbb{Z}} \omega_i h_i^n u_i^n$. The conservation of the water height is directly deduced from the evolution law for $h_i^{n+1}$ given by (7). Next, considering the updated formula for the discharge, we easily obtain
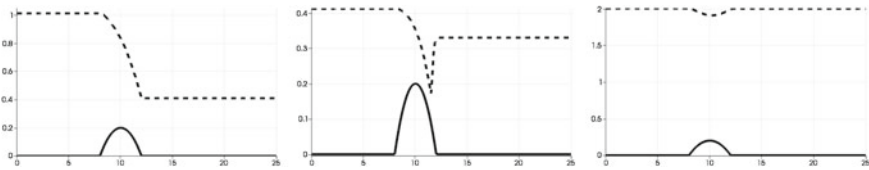
$$\sum_{i \in \mathbb{Z}} \omega_i h_i^{n+1} u_i^{n+1} = \sum_{i \in \mathbb{Z}} \omega_i h_i^n u_i^n - g \sum_{i \in \mathbb{Z}} \omega_i \bar{H}_i \bar{X}_i \sum_j 2\omega_j H_{ij}(1 - X_{ij}) W'_{ij}$$
$$+ g \sum_{i \in \mathbb{Z}} \omega_i \bar{H}_i^2 \bar{X}_i (1 - \tilde{X}_i) \sum_j 2\omega_j W'_{ij}$$

Now, we have to define the average functions ($H_{ij}$, $\bar{H}_i$, $\bar{X}_i$ and $\tilde{X}_i$) such that the discharge conservation is recovered as soon as the topography function is a given constant $Z$. Of course, providing that the consistency conditions (4) holds true, we have just to consider average functions such that $H_{ij}(1 - X_{ij}) = Z$.

If (4) is not satisfied, we enter the delicate problem of the *inconsistency* of the SPH technique. Let us assume that the average functions satisfy the following condition as soon as the topography is a given constant $Z$: $H_{ij}(1 - X_{ij}) = Z$ and

**Fig. 1** *Left* free surface profile for the lake at rest with the above defined topography. *Right* velocity and free surface errors for this lake at rest. Both graphs show the solutions at time $t = 600\,\text{s}$



**Fig. 2** From *left* to *Right*, the three test cases defined in [6], respectively G1, G2, G3. The *dashed line* represents the free surface and the *full line* is the topography. The graphs show the solutions at time $t = 600\,\text{s}$

$\bar{H}_i(1 - \tilde{X}_i) = Z$. Then we immediately recover the expected conservation of the discharge. Such average functions can be easily obtained. For instance, let us set $\bar{H}_i = h_i^n + Z_i$, $\tilde{X}_i = h_i^n/(h_i^n + Z_i)$ and $H_{ij} = H_i$ if $(Hu)_{ij} > 0$, $H_j$ otherwise.

## 4 Numerical Experiments

We now illustrate the relevance of the proposed SPH scheme. For all the tests, the computational domain is [0, 25], 200 particles are used, and the gravity constant is equal to 9.81.

To test the well-balanced property, we consider a topography defined by $Z(x) = 0.4e^{(sin(x)-1)}$. The initial conditions are $h(x, 0) + Z(x) = 0.5$, and $u(x, 0) = 0$.

Figure 1 shows that the free surface is unperturbed with an oscillating topography. The velocity is, as expected, close to 0, up to $10^{-15}$. The perturbations appearing the in the velocity are of the order of magnitude of the machine precision, which is confirmed by a simulation in quadruple precision, where the perturbations are close to 0, up to $10^{-33}$.

The next three test cases come from [6]. The topography is flat with a bump for $x \in [8, 12]$, as follows: $Z(x) = 0.2 - 0.05(x - 10)^2$. The transcritical flow without shock (G1), with shock (G2) and subcritical flow (G3) test cases are performed according to the initial and boundary conditions given by [6].

Figure 2 shows good agreement with the exact results (see [3, 6] for instance).

**Table 1** Discharge errors for the three test cases described above. Comparisons between three schemes: the modified SPH scheme as well as the ones introduced in [1, 3]

| Test case | Hydrostatic reconstruction | | Hydrostatic upwind | | SPH scheme | |
|---|---|---|---|---|---|---|
| | $L^2$ error | $L^\infty$ error | $L^2$ error | $L^\infty$ error | $L^2$ error | $L^\infty$ error |
| G1 | 4.35E-2 | 1.92E-2 | 5.98E-2 | 1.87E-2 | 5.67E-2 | 1.85E-2 |
| G2 | 4.88E-2 | 3.31E-2 | 4.68E-2 | 2.85E-2 | 5.50E-2 | 4.02E-2 |
| G3 | 9.62E-2 | 3.07E-2 | 9.78E-2 | 2.70E-2 | 9.83E-2 | 2.74E-2 |

In Table 1, the discharge errors turn out to be similar to other methods like hydrostatic reconstruction.

## 5 Conclusion

By adopting a suitable reformulation of the shallow-water model, we have derived a relevant discretization of the topography source term to enforce a hybrid SPH scheme (introduced in [7, 14]) to be well-balanced. Numerical simulations have been performed to illustrate the interest of such a topography source term discretization. Indeed, usual and hybrid SPH schemes are known to not preserve the constant state because of the kernel function which does not satisfy the consistency conditions (4). The proposed technique corrects such a failure.

## References

1. Audusse, E., Bouchut, F., Bristeau, M.O., Klein, R., Perthame, B.: A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water ows siam. J. Sci. Comp. **25**, 2050–2065 (2004)
2. Bermudez, A., Vazquez-Cendon, M.E.: Upwind methods for hyperbolic conservation laws with source terms. Comput. Fluids. **23**, 1049–1071 (1994)
3. Berthon, C., Foucher, F.: Efficient well-balanced hydrostatic upwind schemes for shallow-water equations. J. Comput. Phys. **231**(15), 4993–5015 (2012)
4. Berthon, C., Foucher, F.: Hydrostatic upwind schemes for shallow-water equations. Finite Volumes for Complex Appl. VI Springer Proc. Math. **4**, 97–106 (2011)
5. Bouchut F.: Nonlinear stability of finite volume methods for hyperbolic conservation laws and well-balanced schemes for sources. Frontiers in Mathematics. Birkhuser Verlag, Basel (2004)
6. Goutal, N., Maurel, F.: Proceedings of the 2nd workshop on dam-break wave simulation. Tech. rep. EDF-DER, HE-43/97/016/B (1997)
7. Guilcher, P.M.: Contribution au développement d'une méthode SPH pour la simulation numérique des interactions houle-structure. Thèse de doctorat, École Centrale de Nantes (2008)

8. Harten, A., Lax, P., van Leer, B.: On upstream differencing and godunov-type schemes for hyperbolic conservation laws. SIAM Rev. **25**(1), 35–61 (1983)
9. Kröner D.: Finite volume schemes in multidimensions. Numerical analysis 1997 (Dundee), Pitman Res. Notes Math. Ser. 380, 179–192. Longman, Harlow (1998)
10. Mas-Gallic, S., Raviart, P.A.: A particle method for first-order symmetric systems. Numerische Mathematik **51**(3), 323–352 (1987)
11. Monaghan, J.J.: Smoothed particles hydrodynamics. Ann. Rev. Astron. Astrophys. **30**, 543–574 (1992)
12. Toro, E.F.: Riemann solvers and numerical methods for fluid dynamics. 3rd edn. A practical introduction. Springer, Berlin (2009)
13. Vila, J.P.: On particle weighted methods and smooth particle hydrodynamics. Math. Models Methods Appl. Sci. **09**, 161 (1999)
14. Vila, J.P.: Sph renormalized hybrid methods for conservation laws: applications to free surface flows. Meshfree Methods for Partial Differential Equations II **43**, 207–229 (2005)

# Asymptotic-Preserving Scheme Based on a Finite Volume/Particle-In-Cell Coupling for Boltzmann-BGK-Like Equations in the Diffusion Scaling

**Anaïs Crestetto, Nicolas Crouseilles and Mohammed Lemou**

**Abstract**  This work is devoted to the numerical simulation of the collisional Vlasov equation in the diffusion limit using particles. To that purpose, we extend the Finite Volumes/Particles hybrid scheme developed in [5], based on a micro-macro decomposition technique introduced in [1] or [13]. Whereas a uniform grid was used to approximate both the micro and the macro part of the full distribution function in [13], we use here a particle approximation for the kinetic (micro) part, the fluid (macro) part being always discretized by standard finite volume schemes. There are many advantages in doing so: (i) the so-obtained scheme presents a much less level of noise compared to the standard particle method; (ii) the computational cost of the micro-macro model is reduced in the diffusion limit since a small number of particles is needed for the micro part; (iii) the scheme is asymptotic preserving in the sense that it is consistent with the kinetic equation in the rarefied regime and it degenerates into a uniformly (with respect to the Knudsen number) consistent (and deterministic) approximation of the limiting equation in the diffusion regime.

A. Crestetto (✉)
Laboratoire de Mathématiques Jean Leray, Université de Nantes, 2 rue de la Houssinière, 44322 Nantes, France
e-mail: anais.crestetto@univ-nantes.fr

N. Crouseilles
INRIA Rennes Bretagne-Atlantique, Projet IPSO, Institut de Recherche Mathématique de Rennes, Université de Rennes 1, 263 avenue du Général Leclerc, 35042 Rennes, France
e-mail: nicolas.crouseilles@inria.fr

M. Lemou
Institut de Recherche Mathématique de Rennes—CNRS, Université de Rennes 1, 263 avenue du Général Leclerc, 35042 Rennes, France
e-mail: mohammed.lemou@univ-rennes1.fr

# 1 Introduction

Particle systems appearing in several physical applications like plasma or radiative transfer can be studied at different scales. A kinetic description is necessary when the system is far from thermodynamical equilibrium. It is based on the representation of the set of particles by a distribution function $f$ depending on time $t$, space $x$ and velocity $v$, $f$ verifying a partial differential equation of Vlasov-type. When the system stays near equilibrium, the problem can be reduced using a macroscopic description, only depending on $t$ and $x$. Several strategies can be used to solved multiscale problems (see for example [8, 9, 11] or [2]), among them, the micro-macro decomposition introduced in [1] leads to a coupling of two equations: a macroscopic one for the mean part of $f$ (in velocity) and a microscopic one for the remainder part (called perturbation).

This work is devoted to the design of an Asymptotic-Preserving (AP) scheme (see [10]) for the following kinetic equation in the diffusion scaling

$$\partial_t f + \frac{1}{\varepsilon} v \partial_x f + \frac{1}{\varepsilon} E \partial_v f = \frac{1}{\varepsilon^2} (\rho M - f), \tag{1}$$

where $x \in [0, L_x]$, $\rho = \int f \, dv$ is the charge density, $E$ is the electric field given by the Poisson equation $\partial_x E = \rho - 1$, $M$ is either the absolute Maxwellian (in the BGK-case $v \in \mathbb{R}$) or equal to 1 (in the radiative transport equation (RTE)-case $v \in [-1, 1]$) and $\varepsilon$ is the Knudsen number, parameter of the frequency of collisions between particles, that may be of order one or tend to zero in the diffusion limit.

The strategy will be the use of the micro-macro decomposition. However, following [5], we want to use particles to discretize the micro part so that in the limit regime, the numerical cost is reduced since a few number of particles will be necessary to sample the (small) non equilibrium part. The main difficulty compared to phase space grid approaches [6, 13] remains in the fact that the use of particles requires a splitting between transport and source terms whereas in [6, 13], the stiffest (source) term is used to stabilize the stiff transport term.

The outline of the paper is the following. We derive the micro-macro model in Sect. 2 and its numerical discretization in Sect. 3. Some numerical results are given in Sect. 4. Section 5 is devoted to the conclusion and some perspectives.

# 2 Derivation of the Micro-Macro Equations

This section is devoted to the derivation of the micro-macro model. Let us first introduce a velocity-set $V = \mathbb{R}$ in the Vlasov-Poisson-BGK-case and $V = [-1, 1]$ in the RTE-case and define the null space of the linear collisional BGK-operator $Q(f) = \rho M - f$ by $\mathcal{N} = \text{Span}\{M\} = \{f = \rho M, \text{ with } \rho = \langle f \rangle\}$, where $\langle h \rangle := \int h \, dv$ and $M(v) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{v^2}{2}\right)$ is the absolute Maxwellian in the

Vlasov-Poisson-case or $M = 1$ in the RTE-case. We now define the orthogonal pro-
jector $\Pi$ in $L^2\left(M^{-1}dv\right)$ onto $\mathcal{N}$ as $\Pi h := \langle h \rangle M$. Following [6, 13], we decompose
$f$ as $f = \rho M + g$, where $\rho := \langle f \rangle$ and $g := f - \rho M$, and rewrite the kinetic
equation (1) into the equivalent micro-macro model

$$
\begin{cases}
\partial_t \rho + \dfrac{1}{\varepsilon} \partial_x \langle vg \rangle = 0, \\[2mm]
\partial_t g + \dfrac{1}{\varepsilon}(I - \Pi)[v\partial_x(\rho M + g) + E\partial_v(\rho M + g)] = -\dfrac{1}{\varepsilon^2}g.
\end{cases} \tag{2}
$$

The micro equation on $g$ makes appear stiff terms that need a particular treatment
in order to get an AP scheme. The strategy of [12] is used and recalled here. We
rewrite the flux term of the micro equation

$$
(I - \Pi)(v\partial_x(\rho M + g)) = vM\partial_x\rho + v\partial_x g - \partial_x \Pi(vg),
$$
$$
(I - \Pi)(E\partial_v(\rho M + g)) = E\partial_v(\rho M + g) = -vME\rho + E\partial_v g,
$$

so that the micro equation becomes

$$
\partial_t g + \dfrac{1}{\varepsilon}[vM\partial_x\rho + v\partial_x g - \partial_x \langle vg \rangle M - vME\rho + E\partial_v g] = -\dfrac{1}{\varepsilon^2}g. \tag{3}
$$

Starting from (3), we rewrite it as

$$
\partial_t (e^{t/\varepsilon^2} g) = -\dfrac{e^{t/\varepsilon^2}}{\varepsilon}\left[vM\partial_x\rho + v\partial_x g - \partial_x \langle vg^n \rangle M - vME\rho + E\partial_v g\right].
$$

Integrating in time between $t^n$ and $t^{n+1}$ leads to

$$
e^{t^{n+1}/\varepsilon^2}g^{n+1} = e^{t^n/\varepsilon^2}g^n - \dfrac{1}{\varepsilon}\int_{t^n}^{t^{n+1}} e^{t/\varepsilon^2}dt\left[vM\partial_x\rho^n\right.
$$
$$
\left. + v\partial_x g^n - \partial_x \langle vg^n \rangle M - vME^n\rho^n + E^n\partial_v g^n\right],
$$

and multiplying by $e^{-t^{n+1}/\varepsilon^2}$ gives

$$
g^{n+1} = e^{-\Delta t/\varepsilon^2}g^n - \varepsilon(1 - e^{-\Delta t/\varepsilon^2})\left[vM\partial_x\rho^n\right.
$$
$$
\left. + v\partial_x g^n - \partial_x \langle vg^n \rangle M - vME^n\rho^n + E^n\partial_v g^n\right].
$$

By using the discrete time derivative, we finally get

$$
\dfrac{g^{n+1} - g^n}{\Delta t} = \dfrac{(e^{-\Delta t/\varepsilon^2} - 1)}{\Delta t}g^n - \varepsilon\dfrac{(1 - e^{-\Delta t/\varepsilon^2})}{\Delta t}\left[vM\partial_x\rho^n + v\partial_x g^n\right.
$$
$$
\left. - \partial_x \Pi(vg^n)M - vME^n\rho^n + E^n\partial_v g^n\right],
$$

which we approximate, up to terms of order $\mathcal{O}(\Delta t^2)$, by

$$\partial_t g = \frac{e^{-\Delta t/\varepsilon^2} - 1}{\Delta t} g$$
$$- \varepsilon \frac{1 - e^{-\Delta t/\varepsilon^2}}{\Delta t} [vM\partial_x\rho + v\partial_x g - \partial_x \Pi(vg)M - vME\rho + E\partial_v g]. \quad (4)$$

Let us remark that this equation does not contain any stiff term. Moreover, the two following properties are verified:

- consistency: $\forall \varepsilon > 0$ fixed, as $\Delta t \to 0$, we recover the initial micro equation (3),
- AP property: $\forall \varepsilon > 0$ fixed, as $\varepsilon \to 0$, we get $g = -\varepsilon(vM\partial_x\rho - vME\rho)$, which injected in the macro equation provides the right limit model given by $\partial_t \rho - \partial_x (\partial_x \rho - E\rho) = 0$ (see [6]).

The Sect. 3 is devoted to the numerical discretization of the modified micro-macro model

$$\partial_t \rho + \frac{1}{\varepsilon} \partial_x \langle vg \rangle = 0,$$
$$\partial_t g = \frac{e^{-\Delta t/\varepsilon^2} - 1}{\Delta t} g - \varepsilon \frac{1 - e^{-\Delta t/\varepsilon^2}}{\Delta t} [vM\partial_x\rho + v\partial_x g - \partial_x \Pi(vg)M - vME\rho + E\partial_v g]$$
$$(5)$$

by an hybrid scheme, that couples finite volumes for the macro part $\rho$ to a particle method for the micro part $g$. As in [5], where the hydrodynamic limit was studied, we expect a reduction of computational time when $\varepsilon \to 0$, related to the few number of particles needed to represent $g$ at the limit.

## 3 Finite Volumes/Particles Discretization

We present in this section the Finite Volumes/Particle-In-Cell (PIC) coupling developed for solving (5). Such a coupling is explained in more details in [5] for the hydrodynamic limit.

Let us consider a classical uniform discretization of the spatial domain $x \in [0, L_x]$ denoted by $(x_i)_{0 \le i \le N_x}$ and the following approximations: $\rho_i^n \approx \rho(t^n, x_i)$ and $E_i^n \approx E(t^n, x_i)$. The Poisson equation $\partial_x E = \rho - 1$ for $E$ is solved thanks to finite volumes without difficulty. We now focus on the two other Eq. (5).

### 3.1 Particle Approximation for g

Our goal is to extend the particle discretization developed in [5] to the diffusion scaling. To that purpose, we exploit the reformulation (4). As already said in [5],

we have to use a splitting procedure between the transport part and the source part. Then, the algorithm is the following

- solve $\partial_t g + \varepsilon \frac{(1-e^{-\Delta t/\varepsilon^2})}{\Delta t} v \partial_x g + \varepsilon \frac{(1-e^{-\Delta t/\varepsilon^2})}{\Delta t} E \partial_v g = 0$,
- solve $\partial_t g = \frac{(e^{-\Delta t/\varepsilon^2}-1)}{\Delta t} g - \varepsilon \frac{(1-e^{-\Delta t/\varepsilon^2})}{\Delta t} [vM\partial_x\rho + \partial_x\langle vg\rangle M - vME\rho]$.

In the PIC method (described for example in [3]), the distribution function $g$ is represented by a set of $N$ particles of position $x_k$, velocity $v_k$ and weight $\omega_k$ and approximated by $g(t,x,v) = \sum_{k=1}^{N} \omega_k(t)\delta(x - x_k(t))\delta(v - v_k(t))$. Then, the transport part is solved with the (non stiff) characteristics

$$\dot{x}_k = \varepsilon \frac{(1 - e^{-\Delta t/\varepsilon^2})}{\Delta t} v_k, \quad \dot{v}_k = \varepsilon \frac{(1 - e^{-\Delta t/\varepsilon^2})}{\Delta t} E(t, x_k), \qquad (6)$$

$E(t, x_k)$ being computed by a deposition step knowing $E_i^n$ on the mesh. The source part is solved using the equation satisfied by the weights

$$\begin{aligned}\dot{\omega}_k = \frac{(e^{-\Delta t/\varepsilon^2}-1)}{\Delta t}\omega_k - \varepsilon\frac{(1-e^{-\Delta t/\varepsilon^2})}{\Delta t}[v_k M(v_k)(\partial_x\rho(t, x_k) \\ -E(t, x_k)\rho(t, x_k)) + \partial_x\langle vg\rangle(t, x_k)M(v_k)].\end{aligned} \qquad (7)$$

In more details, from an initial repartition of the $N$ particles $(x_k^0, v_k^0)$ in the phase-space domain of size $L_x \times L_v$, with $\omega_k^0 = g(t = 0, x_k, v_k)L_x L_v/N$, (6) is approximated by

$$x_k^{n+1} = x_k^n + \varepsilon(1 - e^{-\Delta t/\varepsilon^2})v_k^n, \quad v_k^{n+1} = v_k^n + \varepsilon(1 - e^{-\Delta t/\varepsilon^2})E^n(x_k^n). \quad (8)$$

Then, we compute the momentum $\langle vg^{n+1/2}\rangle$ of $g^{n+1/2}$ using this new position:

$$\langle vg^{n+1/2}\rangle|_{x=x_i} \approx \sum_{k=1}^{N} \omega_k^n B_\ell(x_i - x_k^{n+1})v_k^{n+1}, \qquad (9)$$

$B_\ell \geq 0$ is a B-spline function of order $\ell$:

$$B_\ell(x) = (B_0 * B_{\ell-1})(x), \quad \text{with} \quad B_0(x) = \begin{cases} \frac{1}{\Delta x} & \text{if } |x| < \Delta x/2, \\ 0 & \text{else.} \end{cases} \qquad (10)$$

We rewrite the weight equation as

$$\omega_k^{n+1} = \omega_k^n + (e^{-\Delta t/\varepsilon^2} - 1)\omega_k^n - \varepsilon(1 - e^{-\Delta t/\varepsilon^2})[\alpha_k^n + \beta_k^n], \qquad (11)$$

with $\alpha_k^n = v_k M(v_k)[\partial_x\rho^n(x_k) - E^n(x_k)\rho^n(x_k)]\frac{L_x L_v}{N}$ and $\beta_k^n = \partial_x\langle vg^n\rangle(x_k)M(v_k)\frac{L_x L_v}{N}$.

To compute $\alpha_k^n$ (resp. $\beta_k^n$), since $\rho^n$ (resp. $\langle vg^n \rangle$) is known on the spatial grid, we approximate $\partial_x \rho^n$ (resp. $\partial_x \langle vg^n \rangle$) by centered finite differences: $(\partial_x \rho^n)_i \approx \frac{\rho_{i+1}^n - \rho_{i-1}^n}{2\Delta x}$ (resp. $(\partial_x \langle vg^n \rangle)_i \approx \frac{\langle vg^n \rangle_{i+1} - \langle vg^n \rangle_{i-1}}{2\Delta x}$) and evaluate at $x = x_k$ using an interpolation.

**Remark** We have now a new approximation of $g^{n+1}$ given by its particle discretization. We have to ensure that the micro-macro structure $f = \rho M + g$ with $\rho = \int f \, dv$ is preserved numerically. To do that, we correct the weights $\omega_k^{n+1}$, adapting an idea of [7]. We do not detail this procedure here but refer the reader to [5].

**Chapman-Enskog expansion** When $\varepsilon$ goes to zero, we immediately observe that $\omega_k^{n+1} = -\varepsilon \alpha_k^n + \mathcal{O}(\varepsilon^2)$ (since $\omega_k^n = \mathcal{O}(\varepsilon) \, \forall n \geq 1$). Computing the momentum of $g^{n+1}$ means that we use (9) with $g^{n+1}$, or in the limit regime

$$
\begin{aligned}
\langle vg^{n+1} \rangle|_{x=x_i} &\approx -\varepsilon \sum_{k=1}^{N} \alpha_k^n B_\ell(x_i - x_k) v_k + \mathcal{O}(\varepsilon^2), \\
&\approx -\varepsilon \left[ \langle v^2 M \rangle (\partial_x \rho^n - E^n \rho^n) \right]|_{x=x_i} + \mathcal{O}(\varepsilon^2) \\
&\approx -\varepsilon (\partial_x \rho^n - E^n \rho^n)(x_i) + \mathcal{O}(\varepsilon^2).
\end{aligned}
$$

Injecting in the macro equation then leads to a discretization of $\partial_t \rho - \partial_x (\partial_x \rho - E\rho) = 0$, which corresponds to the right asymptotic model (see [6]).
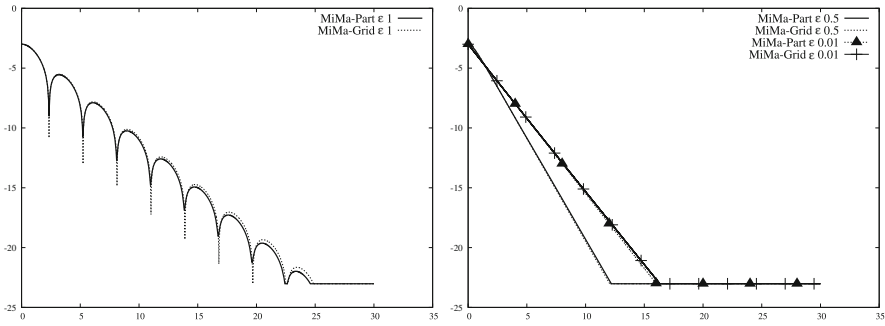
### 3.2 Coupling Strategy

After the computation of $g^{n+1}$ by the PIC method, we compute $\rho^{n+1}$ thanks to a standard finite volume method. We use for example the following scheme:

$$
\rho_i^{n+1} = \rho_i^n - \frac{1}{\varepsilon} \frac{\Delta t}{\Delta x} \left( \langle vg^{n+1} \rangle_{i+\frac{1}{2}} - \langle vg^{n+1} \rangle_{i-\frac{1}{2}} \right), \tag{12}
$$

where $\langle vg^{n+1} \rangle_{i+\frac{1}{2}}$ is computed with (9).

Finally, the algorithm reduces to:

- Initialization of $(x_k, v_k)$ and $\omega_k$.
- (1) Advance micro part:

    – advance the characteristics with (8),
    – compute $\langle vg^n \rangle$ with (9),
    – advance the weights equation with (11).

- (2) Correction step for preserving the micro-macro structure as in [5].
- (3) Advance macro part:

    – compute $\langle vg^{n+1} \rangle$ with (9),
    – compute $\rho^{n+1}$ with (12).

**Fig. 1** Landau damping test case. MiMa-Part compared to MiMa-Grid. Electric energy as a function of time $t$ for $\varepsilon = 1$ on the *left* and $\varepsilon = 0.5$ and $10^{-2}$ on the *right*
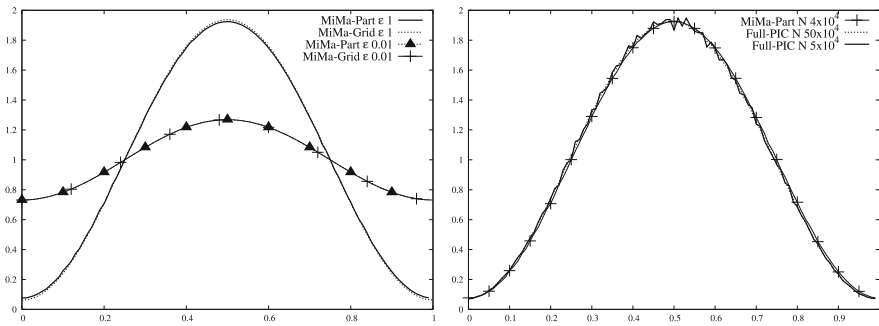
## 4 Numerical Results

We validate our model, denoted by MiMa-Part, on two classical test cases and compare it to a full particle method ($f$ is discretized by particles and not only $g$, see [3]) denoted by Full-PIC and to a micro-macro scheme using a Eulerian discretization of phase space, denoted by MiMa-Grid (which corresponds to the scheme developed in [6]).

We first consider the linear Landau damping case, where $f$ is initially given by $f(0, x, v) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{v^2}{2}\right)(1 + \alpha \cos(kx))$, $x \in \left[0, \frac{2\pi}{k}\right]$ with periodic conditions in $x$ and $v \in \mathbb{R}$ (cut to $[-10, 10]$, assuming that the number of gas particles having a larger absolute velocity is negligible). We take here $k = 0.5$ and $\alpha = 10^{-2}$. The hybrid MiMa-Part scheme is compared to MiMa-Grid for different values of $\varepsilon$, from $10^{-2}$ to 1, in Fig. 1. We look at the time evolution of $\log \|E(t)\|_{L^2}$ which is known to decrease linearly in time. The kinetic regime ($\varepsilon = 1$—on the left) is well described but the number of needed particles is quite big: $N = 5 \times 10^5$. For $\varepsilon = 0.5$ (on the right), we note that the boundary layer is captured by both methods. For small values of $\varepsilon$ (for example $10^{-2}$ but decreasing $\varepsilon$ does no more change the curves - on the right), MiMa-Part leads to the diffusion limit, as well as MiMa-Grid. But here, 200 particles are sufficient to represent in a good way the perturbation $g$ and to capture the limit. The cost of MiMa-Part then reduces as $\varepsilon \to 0$, whereas MiMa-Grid keeps the same complexity.

We then consider the RTE-testcase given by $f(0, x, v) = 1 + \cos\left(2\pi\left(x + \frac{1}{2}\right)\right)$, $E = 0 \,\forall t$, $x \in [0, 1]$ with periodic conditions in $x$ and $v \in [-1, 1]$. Results obtained at $t = 0.1$ are presented in Fig. 2. On the left, MiMa-Part is compared to MiMa-Grid for $\varepsilon = 1$ and $\varepsilon = 10^{-2}$. In both regimes, our hybrid scheme gives a good representation of the density $\rho(x)$. These results can also be compared to those of [12]. On the right, we compare MiMa-Part to a full PIC method when $\varepsilon = 1$. From the PIC point of view, our hybrid scheme can be seen as a $\delta f$ method (see [4] for example). We thus take the same advantages: the noise due to the probabilistic

**Fig. 2** RTE-testcase. Density $\rho$ as a function of $x$ at time $t = 0.1$. MiMa-Part compared to MiMa-Grid for $\varepsilon = 1$ and $10^{-2}$ on the *left* and compared to a full PIC method for $\varepsilon = 1$ on the *right*

character of the particles discretization is reduced since it affects only the perturbation $g$, and not the whole function $f$. This noise appears on the representation of $\rho$ when $N$ is too small, and for example in the black curve obtained with Full-PIC and $N = 5 \times 10^4$. With the same order of $N$, the black line labeled with crosses corresponding to MiMa-Part and $N = 4 \times 10^4$ is not affected by this noise. Finally, for obtaining a smooth curve with Full-PIC, we have to take $N = 5 \times 10^5$ (see the dashed line). The cost of the model is directly linked to $N$. To obtained the two smooth curves for $\varepsilon = 1$, the computational time is $0.12$ s for MiMa-Part and $0.47$ s for Full-PIC.

## 5 Conclusion and Perspectives

A first extension of the AP hybrid method developed in [5] is presented in this paper, concerning the diffusion scaling. Same conclusions are observed: the scheme is AP and the number of needed particles to represent $g$ in a good way decreases as $\varepsilon \to 0$. The cost of the hybrid method reduces then at the diffusion limit, whereas it does not depend on $\varepsilon$ in standard phase-space grid methods.

Other possible extensions may be considered and will be the subject of future works. First, it would be interesting to deal with Dirichlet boundary conditions (instead of periodic ones) for enlarging the application field. For the same reason, more general collision operators should be considered, combining this approach with relaxation techniques as in [12]. Extension to higher dimensions of the phase-space is also possible and a comparison with semi-Lagrangian schemes would be interesting.

# References

1. Bennoune, M., Lemou, M., Mieussens, L.: Uniformly stable numerical schemes for the Boltzmann equation preserving the compressible Navier-Stokes asymptotics. J. Comput. Phys. **227**, 3781–3803 (2008)
2. Berthon, C., Turpault, R.: A numerical correction of the M1-model in the diffusive limit. Discrete Continuous Dynamical Systems Series S **5**, 245–255 (2012)
3. Birdsall, C.K., Langdon, A.B.: Plasma Physics via Computer Simulation. McGraw-Hill, New York (1985)
4. Brunner, S., Valeo, E., Krommes, J.A.: Collisional delta-f scheme with evolving background for transport time scale simulations. Phys. Plasmas **6**, 4504–4521 (1999)
5. Crestetto, A., Crouseilles, N., Lemou, M.: Kinetic/fluid micro-macro numerical schemes for Vlasov-Poisson-BGK equation using particles. Kinet. Relat. Models **5**, 787–816 (2012)
6. Crouseilles, N., Lemou, M.: An asymptotic preserving scheme based on a micro-macro decomposition for collisional Vlasov equations: diffusion and high-field scaling limits. Kinet. Relat. Models **4**, 441–477 (2011)
7. Degond, P., Dimarco, G., Pareschi, L.: The moment guided Monte Carlo method. Int. J. Numer. Meth. Fluids **67**, 189–213 (2011)
8. Dimarco, G., Pareschi, L.: Asymptotic preserving implicit-explicit Runge-Kutta methods for non linear kinetic equations. SIAM J. Numer. Anal. **51**, 1064–1087 (2013)
9. Golse, F., Jin, S., Levermore, C.D.: A domain decomposition analysis for a two-scale linear transport problem. Math. Model. Numer. Anal. **37**, 869–892 (2003)
10. Jin, S.: Efficient asymptotic-preserving (AP) schemes for some multiscale kinetic equations. SIAM J. Sci. Comput. **21**, 441–454 (1999)
11. Klar, A.: An asymptotic-induced scheme for nonstationary transport equations in the diffusive limit. SIAM J. Numer. Anal. **35**, 1073–1094 (1998)
12. Lemou, M.: Relaxed micro-macro schemes for kinetic equations. C. R. Math. **348**, 455–460 (2010)
13. Lemou, M., Mieussens, L.: A new asymptotic preserving scheme based on micro-macro formulation for linear kinetic equations in the diffusion limit. SIAM J. Sci. Comput. **31**, 334–368 (2008)

# Some Applications of a Two-Fluid Model

**Fabien Crouzet, Frédéric Daude, Pascal Galon, Jean-Marc Hérard, Olivier Hurisse and Yujie Liu**

**Abstract** We present in this paper some comparisons of numerical results and experimental data in some two-phase flows involving rather high pressure ratios. A two-fluid two-phase flow model has been used herein, but we also report a few results obtained with some simpler single-fluid two-phase flow models.

## 1 Introduction

The correct modelling of two-phase flows still requires a further investigation of models and methods, but also demands more details and a thorough comparison with available experimental data. For most of the water-vapour applications arising within the framework of nuclear power plants, the vapour phase is dilute; however, the mean flow may sometimes contain a much larger amount of vapour (this may

F. Crouzet · F. Daude · Y. Liu
EDF R&D, AMA, LAMSID, UMR EDF/CNRS/CEA 2832, 1 avenue du Général de Gaulle, 92141 Clamart, France
e-mail: fabien.crouzet@edf.fr

F. Daude
e-mail: frederic.daude@edf.fr

Y. Liu
e-mail: yujie.liu@edf.fr

P. Galon
CEA Saclay, Gif sur Yvette, France
e-mail: pascal.galon@cea.fr

J.-M. Hérard (✉) · O. Hurisse
EDF R&D, MFEE, 6 quai Watier, 78400 Chatou, France
e-mail: jean-marc.herard@edf.fr

O. Hurisse
e-mail: olivier.hurisse@edf.fr

occur in the upper part of steam generators, or more likely in some severe accident configurations following the boiling crisis, or in water-hammer situations), and thus relative velocities may become large. This, among other reasons, has motivated the focus on a class of two-fluid models for which the numerical simulation of highly unsteady flows is relevant. Actually, when restricting to the statistical averaging formalism, we know that standard tools may be used in order to derive meaningful models, in order to tackle unsteady and inhomogeneous two-phase flow patterns.

The two-fluid two-phase flow model discussed herein belongs to a wider class that has been investigated in [3, 4, 8, 10, 11, 16, 17] among other references. It requires the computation of seven unknowns (statistical void fraction of the vapour, mean densities, mean velocities and mean pressures). As recalled in [7, 15] for instance, partial differential equations may be derived for statistical void fractions, and partial mass, momentum and total energy within each phase ; equations of state which provide the mean internal energy within each phase must be prescribed, and some other closure laws for cross-correlations and interfacial transfer terms are also necessary.

We recall in Sect. 2 the governing equations and their main properties ; afterwards we briefly describe the basics of the Finite Volume scheme that is used for numerical simulations. Then we focus on the main part, which consists in reporting some numerical results that have been obtained in [18], thus including a comparison with experimental data [19, 21], but also with other numerical results.

## 2 Governing Equations

Classical notations are used, hence $\alpha_k(x, t)$ will denote the statistical void fraction of phase $k = l, v$, and will comply with the constraint $\alpha_l(x, t) + \alpha_v(x, t) = 1$. Variables $\rho_k$, $U_k$, $P_k$ respectively denote the mean density, the mean velocity, the mean pressure within phase $k$, and we define partial masses $m_k = \alpha_k \rho_k$. The total energy $E_k$ within phase $k = l, v$ is defined by: $E_k = \rho_k e_k(P_k, \rho_k) + \rho_k(U_k^2)/2$, where $e_k(P_k, \rho_k)$ stands for the internal energy. The state variable $W$ will be noted:

$$W^t = (\alpha_v, m_l, m_v, m_l U_l, m_v U_v, \alpha_l E_l, \alpha_v E_v)$$

Thus, when neglecting the contribution of viscous effects and turbulence, the form of the governing equations of mean quantities in the two-fluid model is, for $k = l, v$:

$$\partial_t (\alpha_v) + V_{int}(W)\partial_x (\alpha_v) = \phi_v(W)$$
$$\partial_t (m_k) + \partial_x (m_k U_k) = \Gamma_k(W)$$
$$\partial_t (m_k U_k) + \partial_x \left(m_k U_k^2\right) + \partial_x (\alpha_k P_k) - \Pi_{int}(W)\partial_x (\alpha_k) = D_k(W) + \Gamma_k(W)\overline{U}_{int}$$
$$\partial_t (\alpha_k E_k) + \partial_x (\alpha_k U_k(E_k + P_k)) + \Pi_{int}(W)\partial_t (\alpha_k) = \psi_k(W) + \overline{U}_{int} D_k(W) + \Gamma_k(W)\overline{H}_{int}$$
$$(1)$$

Contributions $\Gamma_k(W)$, $D_k(W)$ and $\psi_k(W)$ take interfacial mass transfer, drag effects and interfacial heat transfer into account. Besides, the term $\phi_k(W)$ arising

in the governing equation of the statistical void fraction $\alpha_k$ is due to the statistical averaging [7, 15] of the topological equation. The following constraints also hold:

$$\sum_{k=l,v} \Gamma_k(W) = 0 \;\; ; \;\; \sum_{k=l,v} \psi_k(W) = 0 \;\; ; \;\; \sum_{k=l,v} D_k(W) = 0 \;\; ; \;\; \sum_{k=l,v} \phi_k(W) = 0. \tag{2}$$

and we define: $\overline{U}_{int} = (U_l + U_v)/2$ and: $\overline{H}_{int} = U_l U_v / 2$. Furthermore, we define $V_{int}(W)$ as:

$$V_{int}(W) = \xi(W)U_l + (1 - \xi(W))U_v . \tag{3}$$

where $\xi(W)$ lies in $[0, 1]$. Physically relevant functions $\xi(W)$ have been proposed in [8], and will be recalled at the end of this section. We also introduce the specific entropy $S_k(P_k, \rho_k)$ in each phase, which complies with:

$$c_k^2 \partial_{P_k} (S_k) + \partial_{\rho_k} (S_k) = 0 \tag{4}$$

-noting $c_k(W)$ the speed of acoustic waves within phase $k$- and temperatures: $1/T_k = \partial_{P_k} (S_k) / \partial_{P_k} (e_k)$ ; we also set: $\mu_k = e_k + P_k/\rho_k - T_k S_k$ . Besides, source terms $\Gamma_l(W), \phi_l(W), \psi_l(W), D_l(W)$ are defined as (see property 1):

$$\Gamma_l(W) = K_\Gamma(W)(\mu_v(W)/T_v - \mu_l(W)/T_l) \qquad ; \quad D_l(W) = \quad K_U(W)(U_v - U_l) ;$$
$$\psi_l(W) = K_T(W)(T_v - T_l) \qquad\qquad\qquad ; \quad \phi_l(W) = \quad K_P(W)(P_l - P_v)$$

The first three closure laws are in agreement with classical formulations (see [7, 16]), and the last one for $\phi_l(W)$ is physically relevant: it simply means that the statistical void fraction of the liquid phase locally increases when $P_l > P_v$. The -positive-scalar functions in the drag contribution and in the heat transfer closure law may be chosen as:

$$K_U(W) = m_l m_v((m_l + m_v)\tau_U(W))^{-1},$$

$$K_T(W) = m_l m_v C_{l-v}((m_l + m_v)\tau_T(W))^{-1},$$

$$K_P(W) = \alpha_l \alpha_v((P_l + P_v)\tau_P(W))^{-1}.$$

Here, $\tau_{U,P,T}(W)$ denote velocity-pressure-temperature relaxation time scales, and we also set : $K_\Gamma(W) = K'_\Gamma(W)/\tau_\Gamma(W)$. Closure laws for $\tau_{U,P,T,\Gamma}(W)$ can be found in the literature (see [9] for a review concerning $\tau_P$). Eventually, we assume that $\Pi_{int}(W)$ is a convex combination of both pressures, thus:

$$\Pi_{int}(W) = \chi(W)P_l + (1 - \chi(W))P_v \tag{5}$$

with:

$$\chi(W) = \frac{(1 - \xi(W))/T_l}{(1 - \xi(W))/T_l + \xi(W)/T_v} \tag{6}$$

**Property 1:**
*For smooth solutions W of* (1) *with closure laws* (3), (5), (6), *the governing equation of the entropy of the two-fluid model* $\eta(W) = \sum_{k=l,v} m_k S_k$ *is:*

$$\partial_t (\eta(W)) + \partial_x \left( \sum_{k=l,v} m_k U_k S_k \right) = \Gamma_l(W)(\mu_v(W)/T_v - \mu_l(W)/T_l)$$

$$+ D_l(W)(U_v - U_l)(1/(2T_v) + 1/(2T_l))$$
$$+ \psi_l(W)(T_v - T_l)/(T_v T_l)$$
$$+ \phi_l(W)(P_l - P_v)((1 - \chi(W))/T_v + \chi(W)/T_l)$$

Obviously, when $\xi(W) = 0$ (or $\xi(W) = 1$), one retrieves the standard Baer-Nunziato model [3], where the interface velocity $V_{int}(W)$ corresponds to the mean velocity of the vanishing phase [3, 4, 10, 17]. We finally recall two basic properties:

**Property 2:** *The set of equations associated with the left-hand side of* (1) *has seven real eigenvalues which read:*

$$\lambda_1 = V_{int}(W) \tag{7}$$
$$\lambda_2 = U_v, \quad \lambda_3 = U_v - c_v(W), \quad \lambda_4 = U_v + c_v(W), \tag{8}$$
$$\lambda_5 = U_l, \quad \lambda_6 = U_l - c_l(W), \quad \lambda_7 = U_l + c_l(W) \tag{9}$$

*Associated righteigenvectors span the whole space* $\mathcal{R}^7$, *if:* $|U_k - V_{int}(W)|/c_k \neq 1$.

**Property 3:** *Fields associated with eigenvalues* $\lambda_{2,5}$ *are linearly degenerate. Other fields associated with eigenvalues* $\lambda_{3,4,6,7}$ *are non linear. The 1-field is linearly degenerate if:* $\xi(W)(1 - \xi(W)) = 0$, *or if:* $\xi(W) = m_l/(m_l + m_v)$.

If the 1-field is linearly degenerate, unique jump conditions can be written within each single field. Thus, for schemes that provide convergent approximations when the mesh is refined, we expect that approximations converge towards the unique shock solution. Other properties can be found in [5].

## 3 Finite Volume Scheme

The basic algorithm that is used to compute approximations of the whole system relies on an entropy-consistent fractional step method including an evolution step and a relaxation step. Details on schemes can be found in references [8, 12–14].

- *Evolution step*
  This step computes approximate solutions of the homogeneous system:

$$\begin{cases} \partial_t\,(\alpha_v) + V_{int}(W)\partial_x\,(\alpha_v) = 0 \\ \partial_t\,(m_k) + \partial_x\,(m_k U_k) = 0 \\ \partial_t\,(m_k U_k) + \partial_x\,\left(m_k U_k^2\right) + \partial_x\,(\alpha_k P_k) - \Pi_{int}(W)\partial_x\,(\alpha_k) = 0 \\ \partial_t\,(\alpha_k E_k) + \partial_x\,(\alpha_k U_k(E_k + P_k)) + \Pi_{int}(W)\partial_t\,(\alpha_k) = 0 \end{cases} \quad (10)$$

through the time interval $[t^n, t^n + \Delta t]$, with given initial values $W^n$. The Finite Volume solver that is used to compute interface fluxes either relies on a non-conservative version of the Rusanov scheme, on the approximate VFRoe-ncv Godunov scheme (see [8]), or on the relaxation scheme introduced in [20] (see [1, 2] too). An explicit CFL condition enforces the time step. This provides a set of approximations $\tilde{W}$. An extensive verification of convective schemes can be found in [6], with focus on solutions on one-dimensional Riemann problems.
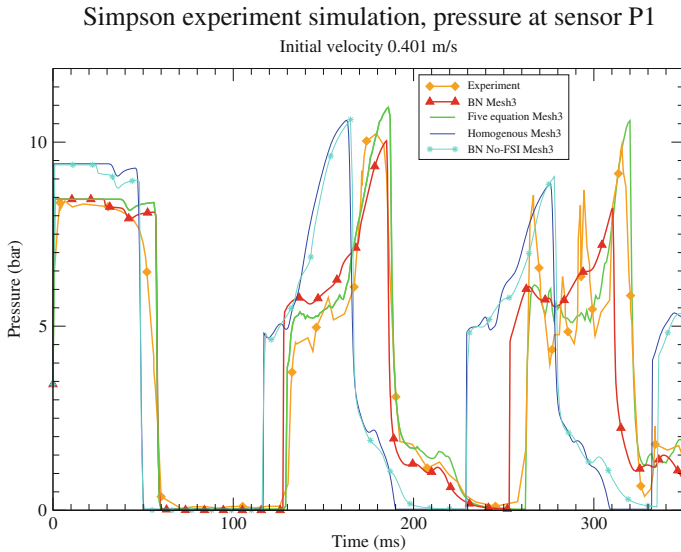
• *Relaxation step*

Given discrete cell values of $\tilde{W}$, we compute approximations of the coupled set of ODEs corresponding to relaxation terms, that is:

$$\begin{cases} \partial_t\,(\alpha_v) = \phi_v(W) \\ \partial_t\,(m_k) = \Gamma_k(W) \\ \partial_t\,(m_k U_k) = D_k(W) + \Gamma_k(W)\overline{U}_{int} \\ \partial_t\,(\alpha_k E_k) + \Pi_{int}(W)\partial_t\,(\alpha_k) = \psi_k(W) + \overline{U}_{int} D_k(W) + \Gamma_k(W)\overline{H}_{int} \end{cases} \quad (11)$$

The most difficult task in the building of the Finite Volume solver is due to the mass transfer term and to the contribution $\phi_k$. In particular, difficulties arise when enforcing the conservative form for the mixture, and meanwhile requesting that void fractions and pressures should remain in their physical range. Many details on this part of the algorithm can be found in [12–14].

## 4 A Comparison of Computational Results with Experimental Data

We provide numerical results and a comparison with experimental data for two distinct cases characterized by high pressure variations. A stiffened gas equation of state (EOS) has been used in the liquid phase, whereas a perfect gas EOS is retained for the vapour phase. A non-conservative version of the Rusanov scheme has been used for all numerical experiments presented in the sequel. As mentionned before, other stable and accurate schemes proposed in the literature may be used, but we emphasize that one of the most difficult tasks also consists in building efficient schemes in order to account for mass transfer and pressure relaxation.
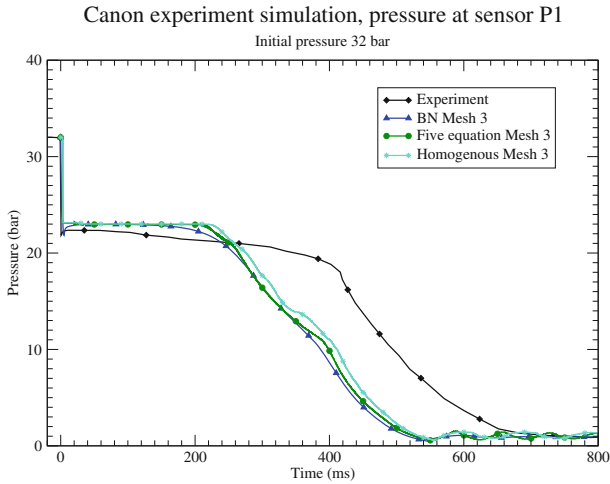
**Fig. 1** Time evolution of the pressure $P = \alpha_v P_v + (1 - \alpha_v) P_l$ in Simpson experiment (*orange squares*). Numerical results: *Red triangles*—two-fluid model with fluid-structure interaction, *Green line*—five-equation homogeneous model with fluid-structure interaction, *Light blue line*—two-fluid model without fluid-structure interaction, *Dark blue line*—three-equation homogeneous model without fluid-structure interaction
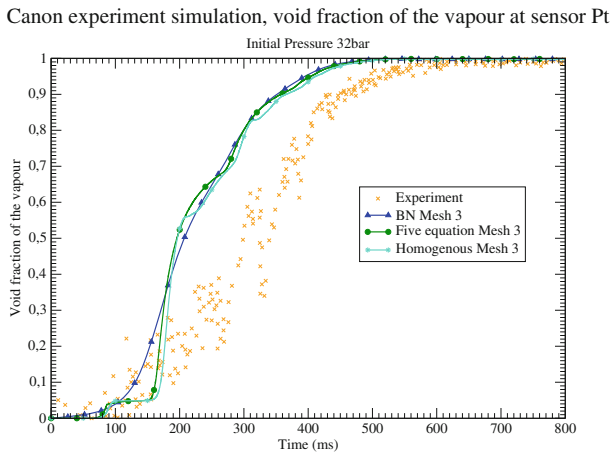
## 4.1 Simpson experiment

This experiment is described in [21]. A big tank is filled with water that flows in a small pipe, the diameter of which is 19 mm; at the very beginning of the recording, the velocity of the fluid is equal to 0.4 m/s, the pressure in the tank is $3.419 \times 10^5$ Pa, the temperature is $T = 296$ K. The pipe of 36 m length is suddenly closed at its right end; thus it results in a violent water-hammer. A shock wave is created and propagates to the left towards the tank. Three pressure captors have been inserted along the pipe, and focus is given here on the one that is close to the right closed exit. The one-dimensional mesh in the pipe contains 12000 regular cells, and the CFL has been set to 1/2. Numerical results obtained with a finer mesh with 36000 cells hardly differ from the latter -absolute differences are less than 1 %. In Fig. 1, the time evolution of the mean pressure $P = \alpha_v P_v + (1 - \alpha_v) P_l$ for this captor has been displayed (orange squares), and a comparison with numerical results obtained with the two-fluid approach on the fine mesh can be done (light blue line). The red triangles refer to the two-fluid approach when one accounts for the elasticity of the pipe (see [18]). Obviously, the prediction of maximum and mimimum values in the transient, as well as occurences of sudden increases and decreases, highly depends on the elasticity of the pipe, and whether it has been accounted for or not in the whole model. This pattern is even emphasized in some other experiments (for instance in

Canon experiment simulation, pressure at sensor P1



**Fig. 2** Time evolution of the pressure $P = \alpha_v P_v + (1 - \alpha_v) P_l$ in Canon experiment (*black squares*). Numerical results: *Dark blue*—two-fluid model/*Green*—five-equation homogeneous model, *Light blue*—three-equation homogeneous model

Canon experiment simulation, void fraction of the vapour at sensor Pt



**Fig. 3** Time evolution of the vapour statistical fraction in Canon experiment (*orange crosses* for different runs). Numerical results: *Dark blue triangles*—two-fluid model, *Green dots*—five-equation homogeneous model, *Light blue*—three-equation homogeneous model

Romander experiment, where a wave propagates in a pipe including a rigid section and an elasto-plastic section, see [18]).

## *4.2 Canon experiment*

In this second experiment [19], a closed rigid pipe initially filled with pressurized water is suddenly opened at its right end. This results in a sudden vaporisation of the fluid, and a left-going rarefaction wave is propagating in the liquid region. The initial pressure $P = \alpha_v P_v + (1 - \alpha_v) P_l$ in the pipe is $32 \times 10^5 \, Pa$, the initial uniform temperature is $T_v = T_l = 493$ K, and the fluid ($\alpha_v = 10^{-3}$) is at rest: $U_v = U_l = 0$. The same EOS have been used within the liquid and vapour phases for this second experiment, and the time step is still chosen in agreement with the constraint: $CFL = 1/2$. The mesh for which results are displayed contains 8000 cells along the pipe axis. Several data have been collected, and results presented in Fig. 2 (black squares) correspond to the time evolution of the pressure close to the right end. A sudden decrease can be oberved first, followed by an almost contant state corresponding to the saturation pressure; afterwards a second smooth decrease occurs, together with an intense vaporization (see Fig. 3), until the atmospheric pressure is reached. Vapour statistical fractions have been recorded at the same place as time goes on, for different experimental runs (orange crosses in Fig. 3). Numerical results obtained with the two-fluid model on a fine mesh have been plotted on both Figs. 2 and 3, together with approximations provided by two different homogeneous models (a five-equation model and a three-equation model). Obviously the vaporization occurs sooner in the simulation than in the experiment.

## References

1. Ambroso, A., Chalons, C., Coquel, F., Galié, T.: Relaxation and numerical approximation of a two-fluid two-pressure diphasic model. ESAIM: M2AN, 43(6), 1063–1097 (2009)
2. Ambroso, A., Chalons, C., Raviart, P.A.: A Godunov-type method for the seven-equation model of compressible two-phase flow. Comput. Fluids **54**, 67–91 (2012)
3. Baer, M.R., Nunziato, J.W.: A two-phase mixture theory for the deflagration to detonation transition (DDT) in reactive granular materials. IJMF **12**(6), 861–889 (1986)
4. Bdzil, J.B., Menikoff, R., Son, S.F., Kapila, A.K., Stewart, D.S.: Two phase modelling of deflagration to detonation transition in granular materials: a critical examination of modelling issues. Phys. Fluids **11**, 378–402 (1999)
5. Coquel, F., Hérard, J.M., Saleh, K., Seguin, N.: Two properties of two-velocity two-pressure models of two-phase flows. Commun. Math. Sci. **12**(3), 593–600 (2014)
6. Crouzet, F., Daude, F., Galon, P., Helluy, P., Hérard, J.M., Hurisse, O., Liu, Y.: Approximate solutions of the Baer Nunziato model. ESAIM Proc. **40**, 63–82 (2013)
7. Drew, D.A., Passman, S.L.: Theory of Multi-Component Fluids, Applied Mathematical Sciences, vol. 135. Springer, New York (1999)
8. Gallouët, T., Hérard, J.M., Seguin, N.: Numerical modelling of two phase flows using the two-fluid two-pressure approach. Math. Mod. Meth. Appl. Sci. **14**(5), 663–700 (2004)
9. Gavrilyuk, S.: The structure of pressure relaxation terms: one-velocity case (2013). (preprint)

10. Gavrilyuk, S., Saurel, R.: Mathematical and numerical modelling of two-phase compressible flows with micro-inertia. J. Comput. Phys. **175**, 326–360 (2002)
11. Glimm, J., Saltz, D., Sharp, D.H.: Two-phase flow modelling of a fluid mixing layer. J. Fluid Mech. **378**, 119–143 (1999)
12. Hérard, J.M., Hurisse, O.: Schémas d'intégration du terme de relaxation des pressions phasiques pour un modèle bifluide hyperbolique, EDF report H-I81-2009-01514-FR, (2009)
13. Hérard, J.M., Hurisse, O.: Computing two-fluid models of compressible two-phase flows with mass transfer, AIAA paper 2012–2959 (2012)
14. Hérard, J.M., Hurisse, O.: A fractional step method to compute a class of compressible gas liquid flows. Comput. Fluids **55**, 57–69 (2012)
15. Hérard, J.M., Liu, Y.: Une approche bifluide statistique de modélisation des écoulements diphasiques à phases compressibles. EDF report H-I81-2013-01162-FR (2013)
16. Ishii, M., Hibiki, T.: Thermofluid Dynamics of Two-Phase Flow. Springer, Berlin (2006)
17. Kapila, A.K., Son, S.F., Bdzil, J.B., Menikoff, R., Stewart, D.S.: Two phase modeling of a DDT: structure of the velocity relaxation zone. Phys. Fluids **9**(12), 3885–3897 (1997)
18. Liu, Y.: Contribution à la vérification et á la validation d'un modèle diphasique bifluide instationnaire. PhD thesis, Université Aix-Marseille, Marseille, France, 11/09/2013. http://tel.archives-ouvertes.fr/tel-00864567
19. Riegel, B.: Contribution à l'étude de la décompression d'une capacité en régime diphasique. PhD thesis, Université de grenoble (1978)
20. Saleh, K.: Analyse et Simulation Numérique par Relaxation d'Ecoulements Diphasiques Compressibles. Contribution au Traitement des Phases Evanescentes. PhD thesis, Université Pierre et Marie Curie, Paris, France, 26/11/2012. http://tel.archives-ouvertes.fr/tel-00761099
21. Simpson, A.R.: Large water-hammer pressures due to column separation in sloping pipes (transient cavitation). PhD thesis, University of Michigan (1986)

# Numerical Simulation of Flow in a Meridional Plane of Multistage Turbine

**Jiří Fürst, Jaroslav Fořt, Jan Halama, Jiří Holman, Jan Karel, Vladimír Prokop and David Trdlička**

**Abstract**   The paper presents a numerical method, which simulates the circumferentially averaged steady flow of a compressible fluid in a multistage turbine. The method is considered in the analytic mode with known geometry. It is intended as a fast tool to turbine designers, which provides the distribution of the flow parameters in the meridional plane, gives the information about mass flow and estimates the efficiency of turbine. The method is based on the solution of the circumferentially averaged three-dimensional Euler equations complemented by the source terms related to the turbine geometry and to the loss prediction model. The meridional plane is discretized by a structured grid. Equations are solved by a finite volume method with the AUSM type numerical flux. Examples including the transonic flow in a turbine stator and in a stage are presented.

## 1 Introduction

The design of a multi-stage turbine is a very complex problem. Designers at a certain step propose setups from typical turbine components to meet the given operating conditions. This step usually brings necessity to simulate the flow inside a multi-stage turbine for different geometries and flow parameters. Methods based on the streamline curvature and stream functions were widely used in the past. They are unfortunately not able to handle transonic flow and they cannot guarantee the conservation of transported quantities. Fully three-dimensional simulations of the turbulent flow are still inapplicable in this initial step of design, mainly due to the excessive CPU time. The desirable method must have the low CPU time consumption and it should deliver

J. Fürst · J. Fořt · J. Halama (✉) · J. Holman · J. Karel · V. Prokop · D. Trdlička
FME CTU Prague, Karlovo nám. 13, 121-35  Prague, Czech Republic
e-mail: Jan.Halama@fs.cvut.cz

J. Fürst
e-mail: Jiri.Furst@fs.cvut.cz

results 'close' to the three-dimensional simulations. There is a variety of simplified approaches ranging from a quasi 1D solvers, e.g. [4], to circumferentially averaged Euler solvers, e.g. [3] or to circumferentially averaged Navier-Stokes solvers, e.g. [6]. Neglected phenomena (averaging, viscous effects, ...) are included in the form of source terms, e.g. [2] or [7]. The presented method is based on the idea of [3]. It solves the circumferentially averaged three-dimensional Euler equations coupled with different loss prediction models (dissipation phenomena). It is able to simulate the flow field in the meridional plane. The choice of equations permits the use of rather coarse grid (no need for the grid refinement along the walls). This is important, since the low CPU time consumption is one of the key requirements for the presented method. The given blade geometry defines the shape of the midplane between the pressure and the suction sides of blades. The shape of the midplane controls the direction of the flow. The loss prediction model includes also an incidence and a deviation corrections.

## 2 Model of the Circumferentially Averaged Flow

Consider the Euler equations in the frame defined by the cylindrical coordinates attached to the respective blade row (the relative frame of reference)

$$\frac{\partial W}{\partial t} + \frac{1}{r}\frac{\partial (rF)}{\partial r} + \frac{1}{r}\frac{\partial G}{\partial \varphi} + \frac{\partial H}{\partial z} = B, \tag{1}$$

$$
W = \begin{bmatrix} \rho \\ \rho v_r \\ \rho v_\varphi \\ \rho v_z \\ e \end{bmatrix}, \quad
F = \begin{bmatrix} \rho v_r \\ \rho v_r^2 + p \\ \rho v_r v_\varphi \\ \rho v_r v_z \\ v_r(e+p) \end{bmatrix}, \quad
G = \begin{bmatrix} \rho v_\varphi \\ \rho v_\varphi v_r \\ \rho v_\varphi^2 + p \\ \rho v_\varphi v_z \\ v_\varphi(e+p) \end{bmatrix}, \quad
H = \begin{bmatrix} \rho v_z \\ \rho v_z v_r \\ \rho v_z v_\varphi \\ \rho v_z^2 + p \\ v_z(e+p) \end{bmatrix},
$$

$$
B = \begin{bmatrix} 0; & \rho\frac{v_\varphi^2}{r} + \frac{p}{r} + \rho r\omega^2 + 2\rho\omega v_\varphi; & -\rho\frac{v_r v_\varphi}{r} - 2\rho\omega v_r; & 0; & 0 \end{bmatrix}^T,
$$

where $r$, $\varphi$, $z$ denote the cylindrical coordinates, $t$ denotes the time, $\rho$ is used for the density, $v_r$, $v_\varphi$ and $v_z$ are the velocity components, $e$ is the total energy per unit volume, $p$ is the pressure and $\omega$ is the angular velocity (it is equal to zero for the stator cascade). Assume one blade passage as the solution domain $D = \{[r, \varphi, z] \in \mathbb{R}^3;$ $[r, z] \in D_{rz}, \ \varphi_1(r, z) < \varphi < \varphi_2(r, z)\}$, where $D_{rz} = \{[r, z] \in \mathbb{R}^2; \ r_1(z) < r < r_2(z), z_1(r) < z < z_2(r)\}$ is the projection of $D$ onto the meridional plane ($zr$-plane). The discretization of the domain $D$ is based on the idea of having single cell in the circumferential direction and a common finite volume discretization of the domain $D_{rz}$. Let us denote the projection of an arbitrary finite volume from $D$

into $D_{rz}$ by $K$. Then one can integrate the Eq. (1) with limits $[r, z] \in K \subset D_{rz}$ and $\varphi_1(r, z) < \varphi < \varphi_2(r, z)$

$$\iint\limits_{K} \int\limits_{\varphi_1(r,z)}^{\varphi_2(r,z)} \left( \frac{\partial W}{\partial t} + \frac{1}{r} \frac{\partial (r F)}{\partial r} + \frac{1}{r} \frac{\partial G}{\partial \varphi} + \frac{\partial H}{\partial z} \right) r d\varphi \, dr dz = \iint\limits_{K} \int\limits_{\varphi_1(r,z)}^{\varphi_2(r,z)} B r d\varphi \, dr dz.$$

(2)

Integration with respect to $\varphi$ yields to

$$\iint\limits_{K} \left( \frac{\partial (br W)}{\partial t} + \frac{\partial (br F)}{\partial r} + \frac{\partial (br H)}{\partial z} \right) dr dz$$

$$= \iint\limits_{K} (br B - (F_1, G_1, H_1)\mathbf{n}_1 + (F_2, G_2, H_2)\mathbf{n}_2) \, dr dz,$$

(3)

where $b = \varphi_2 - \varphi_1$ and $\mathbf{n}_i = (\partial \varphi_i / \partial r, -1/r, \partial \varphi_i / \partial z)$ is the normal vector on the boundary $\partial D_{\varphi,i} = \{[r, \varphi, z] \in \mathbb{R}^3; [r, z] \in D_{rz}, \varphi = \varphi_i(r, z)\}$ for $i = \{1, 2\}$. Consider the finite volume $K$ downstream the leading and upstream the trailing edges, then the non-permeability condition applied on $\partial D_{\varphi,i}$ yields $(F_i, G_i, H_i)\mathbf{n}_i = [0, p_i \mathbf{n}_i, 0]^T$, i.e. the right hand side of the Eq. (3) can be written as

$$\iint\limits_{K} \left( br B + p_2 \begin{bmatrix} 0 \\ r\partial \varphi_2 /\partial r \\ -1 \\ r\partial \varphi_2 /\partial z \\ 0 \end{bmatrix} - p_1 \begin{bmatrix} 0 \\ r\partial \varphi_1 /\partial r \\ -1 \\ r\partial \varphi_1 /\partial z \\ 0 \end{bmatrix} \right) dr dz.$$

The above form of the right hand side of the Eq. (3) is also valid for the part of $\partial D_{\varphi,i}$, where the periodicity conditions $(F_1, G_1, H_1) = (F_2, G_2, H_2)$, $\mathbf{n}_1 = \mathbf{n}_2$ and $p_1 = p_2$ are considered. The resulting form of the governing equations is

$$\iint\limits_{K} \frac{\partial (br W)}{\partial t} dz dr + \iint\limits_{K} div(br H, br F) dz dr = \iint\limits_{K} br Q dz dr,$$

(4)

$$Q = B + \frac{p}{b} \left[ 0; \quad \frac{\partial b}{\partial r}; \quad 0; \quad \frac{\partial b}{\partial z}; \quad 0 \right]^T.$$
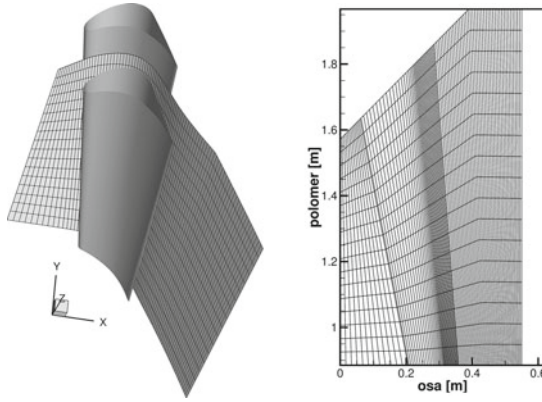
Considering a common explicit finite volume method, one obtains

$$(br W)_K^{n+1} = (br W)_K^n - \frac{\Delta t}{\mu(K)} \sum_l [(br H)_l^n \Delta r_l - (br F)_l^n \Delta z_l] + (br Q)_K^n + S_K,$$
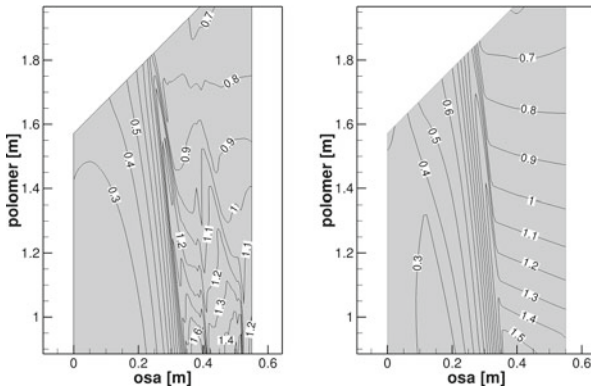
(5)

where the subscript $\cdot_l$ denotes the edges of volume $K$, $n$ is the time level and $S_K = [0, \rho \mathbf{f}^{tot}, 0]^T$ is the additional source term, where $\mathbf{f}^{tot} = \mathbf{f}^b + \mathbf{f}^d$. The external force $\mathbf{f}^b$ is applied between blades and it is perpendicular to the middle plane defined as $\varphi(r, z) = (\varphi_1(r, z) + \varphi_2(r, z))/2 + \varphi_{cor}$. The role of $\mathbf{f}^b$ is to force the fluid to flow along the middle plane, which is given by the blade geometry and it is modified by the incidence and the deviation corrections $\varphi_{cor}$ coming from the loss prediction model. The external force $\mathbf{f}^d$ has the direction of flow and its magnitude is given by the loss prediction model, i.e. it decelerates the flow (dissipation of kinetic energy).

## 3 Numerical Method

The numerical solution is based on a finite volume method for the system (5) coupled with some empirical loss prediction model and the definition of the geometry of a channel. The computational domain (subset of meridional plane) is discretized by the structured quadrilateral grid. The finite volume method uses the AUSM type flux. The coupling between the finite volume method and the loss prediction model has the following steps. Consider the solution $W_{i,j}^n$ at the point $[z_i, r_j, t^n]$ is known. The computational grid has the uniform spacing in the radial direction, therefore grid lines $j = const$ can be roughly considered as streamlines. The solution $W_j^n$ for each particular $j = const$ line is used as an input data for the loss prediction model, which returns the value of the total pressure loss, which is further expressed as the entropy rise $\Delta s_j^{loss}$. We further compute the entropy rise $\Delta s_j^n$ related to the solution $W_j^n$. The solution should satisfy $\lim \Delta s_j^n = \Delta s_j^{loss}$ for $n \to \infty$, i.e. to have solution with the prescribed losses. The correction of external force $\Delta \mathbf{f}_j^d$ is related to the difference $\Delta s_j^{loss} - \Delta s_j^n$. Once the entropy of the solution and from the loss prediction model are equilibrated, the force $\mathbf{f}^d$ becomes constant. The force $\mathbf{f}_j^d$ is appropriately distributed along $j = const$ line between leading and trailing edges. The loss model also gives the correction of the flow direction $\varphi_{cor}$. The component of the velocity, which is normal to the midplane $\varphi(r, z) = (\varphi_1(r, z) + \varphi_2(r, z))/2 + \varphi_{cor}$, defines the correction of force $\Delta \mathbf{f}_{i,j}^b$, which eliminates this normal component. If the normal component is equal to zero, which is the desired state, the force $\mathbf{f}^b$ becomes constant. The force $\mathbf{f}^b$ is applied between leading and trailing edges to mimic the guidance of the flow by the blades. Currently, constant loss model without $\varphi_{cor}$ correction and AMDC-KO [2] loss model are implemented. Next step of numerical algorithm is the evaluation of values in ghost cells (implementation of boundary conditions). Finally $W_{i,j}^{n+1}$ is computed using explicit two stage Runge–Kutta time integration. To enhance the robustness of the proposed method, the addition of external forces is relaxed.

**Fig. 1** The middle section from the 3D Euler simulation (*left*), the meridional plane for the circumferentially averaged simulation (*right*)
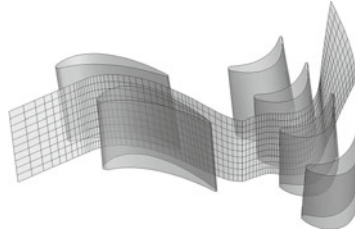


**Fig. 2** The Mach number contours in the meridional plane. Full 3D Euler simulation (*left*), the circumferentially averaged simulation (*right*)
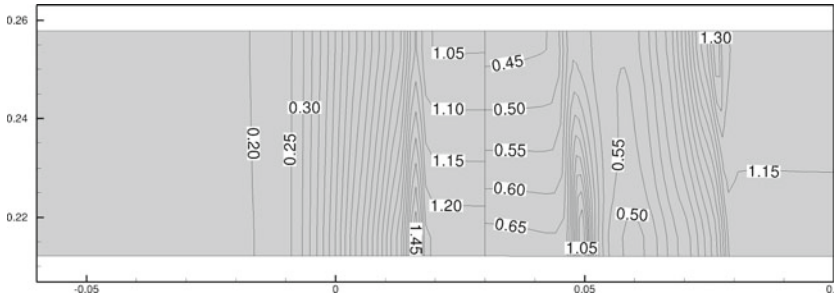
## 4 Results of Simulations

The first example represents the flow in the stator cascade from the low pressure part of a steam turbine. The Fig. 1 shows the the shape of the 'middle' plane between blades and the discretization of domain for the presented method (projection of one blade passage into the meridional plane). The Fig. 2 compares the numerical results achieved by the solution of the full 3D Euler equations and by the circumferentially averaged 3D Euler equations (the presented method without any loss model). The qualitative comparison of the Mach number contours in the meridional plane for the 3D Euler simulation and for the presented method ('single cell' in circumferential direction) gives an idea, how well can the presented method simulate the transonic flow for a complex shape of the meridional section.
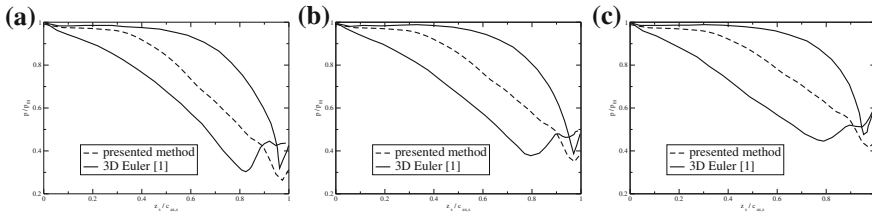
**Fig. 3** The shape of the *middle* section for the high pressure core stage of a gas turbine
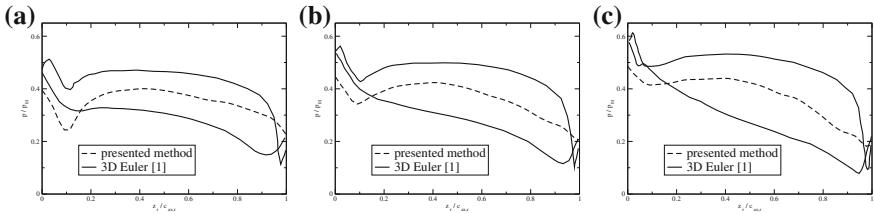


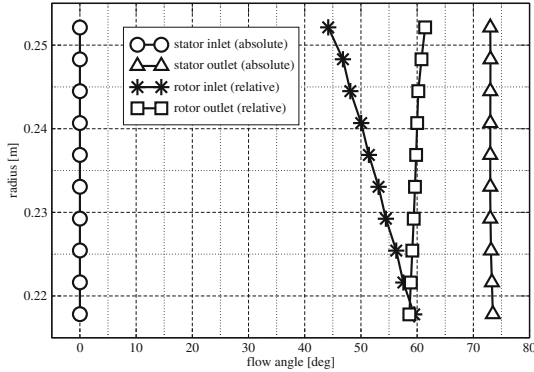**Fig. 4** The Mach number contours in the meridional plane for the high pressure core stage



**Fig. 5** The distribution of the pressure along the stator blade profiles at three radial locations. **a** Hub. **b** Mid. **c** Tip

The second example is the flow in the high pressure core stage of a gas turbine, see the Fig. 3. The stage geometry and experimental data are available in [5]. The presented results correspond to the cold air test described in [5]. Some results are also compared with the results of the 3D Euler simulation from [1]. The considered inlet total pressure is the atmospheric pressure 101.3 kPa and the inlet total temperature is 288.2 K. The flow at the stage inlet has axial direction. The rotor rotates with 8081 rev/min. The ratio of the static pressure behind the stage to the inlet total pressure is 0.225. The total pressure loss has been set 5.5 % for each cascade, it corresponds to the design efficiency considered in [5]. The Fig. 4 shows the contours of the relative Mach number in the meridional plane. The stator blades are prismatic, therefore there is no significant gradient of the solution in the radial direction. The maximum of the Mach number at the hub downstream the stator cascade corresponds to the design value. The Figs. 5 and 6 show the pressure distribution along the blade profiles for

**Fig. 6** The distribution of the pressure along the rotor blade profiles at three radial locations. **a** Hub. **b** Mid. **c** Tip



**Fig. 7** The radial distributions of the relative flow angle in three axial locations (*upstream* the stator cascade, in the gap between stator and rotor and *downstream* the rotor)

the stator and the rotor at three different radial locations. The full line is used for the results of the 3D Euler simulation [1], which gives usually the realistic values of the pressure for the design conditions (no flow separation). The dashed line is the result of the presented method with the single cell in circumferential direction. Due to the single cell, there is only one dashed line, which can be understood as a certain average of the pressure between the pressure and the suction sides—the dashed line should be somewhere in the middle between the both full lines. This may not be true in the vicinity of the leading and trailing edges. It is important, that the presented method is able to approximate well the expansion through the blade channel. This is the main advantage with respect to former methods, which did not take into account the blade geometry. The radial distributions of the flow angles (measured from the axial direction) are plotted in the Fig. 7. The absolute and the relative flow angles in the gap between the stator and the rotor ($\alpha_1$ and $\beta_1$) and the relative flow angle downstream the rotor ($\beta_2$) at the hub, middle section and tip are summarized in the Table 1 together with the value of the total mass flow through the stage. Results achieved by the presented method are compared with the design values and the measured values presented in [5] and with the results of three-dimensional simulation based on the solution of the Euler equations [1]. The agreement between

**Table 1** Comparison of the mass flow and the flow angles

| Case | $\dot{m}$ (kg/s) | $\alpha_1^H$ (deg) | $\alpha_1^M$ (deg) | $\alpha_1^T$ (deg) | $\beta_1^H$ (deg) | $\beta_1^M$ (deg) | $\beta_1^T$ (deg) | $\beta_2^H$ (deg) | $\beta_2^M$ (deg) | $\beta_2^T$ (deg) |
|---|---|---|---|---|---|---|---|---|---|---|
| The design value [5] | 3.708 | 74.4 | 73.0 | 71.8 | 59.3 | 50.8 | 39.3 | 58.7 | 59.6 | 60.7 |
| The experiment [5] | 3.856 | | 72.2 | | | 48.3 | | | 56.4 | |
| 3D Euler [1] | 3.950 | | 72.7 | | | 48.0 | | | 55.0 | |
| The presented method | 3.768 | 73.5 | 73.0 | 73.0 | 59.5 | 52.2 | 44.2 | 58.8 | 59.7 | 60.6 |

The absolute and the relative flow angles downstream the stator are denoted by $\alpha_1$ and $\beta_1$ respectively. The relative flow angle downstream the rotor is denoted by $\beta_2$. The superscripts $\cdot^H$, $\cdot^M$ and $\cdot^T$ refer to hub, mid and tip locations

all results is good. Certain differences can be found for inlet relative flow angle to the rotor, where the presented method gives a slightly higher value at the tip. Nevertheless one has to remember, that even in the original paper [5] there is a difference in stator outlet angle between the design and the real geometry used in experimental setup.

## 5 Conclusions

The presented method based on the iterative coupling of the circumferentially averaged Euler equations with a loss model is able to simulate the transonic flow in the meridional plane of a multistage turbine. Since all viscous effects are modeled by the loss prediction model, the computational grid has less cells compared to methods based on the averaged Navier–Stokes equations. It allows to obtain results for the flow in a multi-stage turbine in a relatively short time (several minutes on today PC's). The value of a loss given by the prediction model is included in the form of external force $\mathbf{f}^d$, which has the same direction as the flow, i.e. it decelerates the flow. Entropy produced by a numerical method (numerical diffusion) is compensated using the same source term. The developed method has a modular character, user can choose from several thermodynamic models (the perfect gas, the steam according to IAPWS IF-97) and from several loss prediction models. First tests have shown, that the presented method is able to provide a reliable information about the mass flow and the radial profiles of the flow angles, of the pressure and of the velocity.

## References

1. Arts, T.: Calculation of the three-dimensional, steady, inviscid flow in a transonic axial turbine stage. J. Eng. Gas Turbines Power. 107(2), 286-292 (1985). doi:10.1115/1.3239713
2. Aungier, R.H.: Turbine Aerodynamics: Axial-Flow and Radial-Inflow Turbine Design and Analysis. ASME Press, New York (2006)
3. Fahua, G., Anderson, M.R.: Cfd-based throughflow solver in turbomachinery design systems, pp. 1259–1267, Paper No. GT2007-7389. doi: 10.1115/GT2007-27389

4. Leonard, O., Adam, O.: A quasi-one-dimensional cfd model for multistage turbomachines. Int. J. Therm. Sci. **17**(1), 7–20 (2008)
5. Moffit, T.P., Szanca, E.M., Whitney, W.J., Behning, F.P.: Design and cold-air test of single uncooled core turbine with high work output. Technical Report NASA Technical paper 1680, Lewis Research Center, Cleveland (1980)
6. Simon, J.F., Leonard, O.: Modeling of 3-d losses and deviations in a throughflow analysis tool. Int. J. Therm. Sci. **16**(3), 208–214 (2007)
7. Zhu, J., Sjolander, S.A.: Improved profile loss and deviation correlations for axial-turbine blade rows, pp. 783–792, Paper No. GT2005-69077. doi:10.1115/GT2005-69077

# Application of a Two-Fluid Model to Simulate the Heating of Two-Phase Flows

**Jean-Marc Hérard, Olivier Hurisse, Antoine Morente and Khaled Saleh**

**Abstract**   This paper is dedicated to the simulation of two-phase flows on the basis of a two-fluid model that allows to account for the disequilibrium of velocities, pressures, temperatures and chemical potentials (mass transfer). The numerical simulations are performed using a fractional step method treating separately the convective part of the model and the source terms. The scheme dealing with the convective part of the model follows a Finite Volume approach and is based on a relaxation scheme. In the sequel, a special focus is put on the discretization of the terms that rule the mass transfer. The scheme proposed is a first order implicit scheme and can be verified using an analytical solution. Eventually, a test case of the heating of a mixture of steam and water is presented, which is representative of a steam generator device.

## 1 Introduction

Most of the industrial processes used for generating electricity require the use of fluids, and especially water. The water is used either as a coolant fluid or to ensure the production of mechanical work through the turbines which are motionned by steam. If we focus on a nuclear power plant based on a Pressurized Water Reactor (PWR), the water is used as liquid or vapour depending on the circuit under consideration.

J.-M. Hérard · O. Hurisse (✉) · A. Morente
EDF R&D, 6 Quai Watier, 78400 Chatou, France
e-mail: olivier.hurisse@edf.fr

J.-M. Hérard
e-mail: jean-marc.herard@edf.fr

A. Morente
e-mail: antoine.morente@edf.fr

K. Saleh
IRSN, BP-3, 13115 Saint-Paul-lez-Durance Cedex, France
e-mail: khaled.saleh@irsn.fr

In particular, the secondary circuit of a nuclear power plant contains steam and liquid water. Moreover, vaporization and condensation phenomena take place in different parts of that circuit. In this industrial context, the two fluid approach is often retained to perform fine 3D simulations in complex geometries.

For instance, the well-known standard two-fluid model [9] is widely used in industrial numerical codes. This model allows to deal with the velocity and temperature disequilibrium, and to take into account the mass transfer between the phases by a source term measuring the distance to the saturation (most of the time in terms of enthalpy or temperaure). In this model, the pressure is assumed to be the same for the two phases at every point and every time. This pressure equilibrium is based on the mechanical assumption of large interfaces between the two phases [9] and it neglects the thermodynamical aspect of the pressure equilibrium. Indeed, the classical thermodynamics theory states that two phases of the same fluid are in equilibrium if and only if: the pressures, the temperatures and the chemical potentials are equal for the two phases. In our opinion, it is crucial to recover this equilibrium condition in a model used to perform numerical simulations of two-phase flows, mainly if mass transfer is an important feature of the problem. We thus choose a model that also takes explicitly into account the pressure disequilibrium between the phases.

The two-fluid model used in the sequel is related to the so-called Baer-Nunziato model [1, 10]. Its formal derivation has been performed following a statistical approach in [8]. In one space dimension, the corresponding system possesses seven independent variables: the statistical fraction of liquid, the statistical mean temperatures, the statistical mean pressures and the statistical mean velocities. The space-time evolution of these variables is described by a set a PDEs whose convective part is hyperbolic and whose source terms are chosen to comply with the entropy inequality, based on the physical mixture entropy. Non-conservative products are present in the equations but some specific closures [4] allow to define discontinuous solutions in a unique manner.

The whole numerical scheme proposed here is based on a operator-splitting method [15]. We first account for the convective part of the system thanks to the explicit relaxation scheme proposed in [2, 13]. The source terms are then successively discretized by four implicit ODE schemes. Very good agreement with experiments has been found in [11] (using a Rusanov scheme for the convective part) focusing on situations where the mass transfer occurs due to a pressure drop. We propose here a one-dimensional test case close to the OECD test case [14]: the mass transfer is due to the heating of saturated water which flows in a pipe.

## 2 The Two-Fluid Model

The system of PDEs governing the time-space evolution of the variables is:

$$\begin{cases} \partial_t(\alpha_v) + V_i(W)\partial_x(\alpha_v) = \phi_v(W), \\ \partial_t(m_l) + \partial_x(m_l U_l) = -\Gamma_v(W), \\ \partial_t(m_l U_l) + \partial_x(m_l U_l^2 + \alpha_l P_l) - P_i(W)\partial_x(\alpha_l) = -D_v(W) - \Gamma_v(W)\overline{U}_{int}, \\ \partial_t(\alpha_l E_l) + \partial_x(\alpha_l U_l(E_l + P_l)) + P_i(W)\partial_t(\alpha_l) = -\psi_v(W) - \overline{V}_{int}D_v(W) - \Gamma_v(W)\overline{H}_{int}, \\ \partial_t(m_v) + \partial_x(m_v U_v) = \Gamma_v(W), \\ \partial_t(m_v U_v) + \partial_x(m_v U_v^2 + \alpha_v P_v) - P_i(W)\partial_x(\alpha_v) = D_v(W) + \Gamma_v(W)\overline{U}_{int}, \\ \partial_t(\alpha_v E_v) + \partial_x(\alpha_v U_v(E_v + P_v)) + P_i(W)\partial_t(\alpha_v) = \psi_v(W) + \overline{V}_{int}D_v(W) + \Gamma_v(W)\overline{H}_{int}, \end{cases} \tag{1}$$

where $\alpha_k$ denote the statistical fractions and satisfy $\alpha_l + \alpha_v = 1$, $\rho_k$ denote the densities, $m_k = \alpha_k \rho_k$ are the partial masses, $U_k$ the velocities, $P_k$ the pressures and $E_k$ the total energies which read $E_k = \rho_k(e_k + U_k^2/2)$. The specific internal energies $e_k$ are obtained through an EOS defined with respect to the pressures and densities: $e_k = e_k(\rho_k, P_k)$. Closure laws have to be provided for the velocities $V_i(W)$, $\overline{V}_{int}(W)$, $\overline{U}_{int}(W)$, for the pressure $P_i(W)$ and for the energy $\overline{H}_{int}(W)$, where $W = (\alpha_l, m_l, m_l U_l, \alpha_l E_l, m_v, m_v U_v, \alpha_v E_v)$. We follow the choice proposed in [4, 6, 7]: $V_i(W) = U_v$, $P_i(W) = P_l$, $\overline{U}_{int} = \overline{V}_{int} = (U_l + U_v)/2$ and $\overline{H}_{int} = U_l U_v/2$. We also define the total mass $m = m_l + m_v$, the mean velocity $U$ with $mU = m_l U_l + m_v U_v$, and the total energy of the mixture $E = \alpha_l E_l + \alpha_v E_v$.

The source terms for the pressure relaxation $\phi_v(W)$, for the mass transfer $\Gamma_v(W)$, for the drag force $D_v(W)$ and for the heat exchange $\psi_v(W)$ are then chosen according to the entropy inequality for the mixture $s = m_l s_l(\rho_l, P_l) + m_v s_v(\rho_v, P_v)$ and the associated entropy-flux $\eta_s = m_l U_l s_l(\rho_l, P_l) + m_v U_v s_v(\rho_v, P_v)$, where $s_k$ are the physical phasic specific entropies. The source terms can then be chosen as:

$$\begin{aligned} \Gamma_v(W) &= \frac{1}{\tau_g(W)} \frac{m_l m_v}{(m_l+m_v)(|\mu_v|/T_v+|\mu_l|/T_l)}(\mu_l/T_l - \mu_v/T_v), \\ D_v(W) &= \frac{1}{\tau_u(W)} \frac{m_l m_v}{m_l+m_v}(U_l - U_v), \\ \psi_v(W) &= \frac{1}{\tau_t(W)} \frac{m_l C_{V,l} m_v C_{V,v}}{m_l C_{V,l}+m_v C_{V,v}}(T_l - T_v), \\ \phi_v(W) &= \frac{\alpha_l \alpha_v}{K_p(W)}(P_v - P_l), \end{aligned} \tag{2}$$

with the positive characteristic time scales $\tau_g$, $\tau_u$, $\tau_t$, and the positive parameter $K_p$ which has the dimension of a kinematic vicosity [5]. The chemical potentials are denoted by $\mu_k = e_k + P_k/\rho_k - T_k s_k$, $T_k = T_k(\rho_k, P_k)$ stand for the temperatures and $C_{V,k}$ are the specific heat capacities.

Model (1) with the closures proposed above is defined for a statistical liquid fraction in ]0, 1[. Otherwise, if for instance $\alpha_l = 0$, the quantities $\rho_l$, $U_l$ and $e_l$ are not defined in a unique manner. It is important to note that due to the choice of the closures for (1), $\alpha_l$ remains in ]0, 1[ if the initial condition for $\alpha_l$ belongs to ]0, 1[ everywhere on the spacial domain and if $\alpha_l$ is in ]0, 1[ on the boundary of the domain (especially at the inlets). Other properties of this model can be found in [3, 4, 8].

With this model, the thermodynamical equilibrium is reached if and only if the temperatures, the pressures and the chemical potentials are equal. In the pressure-

temperature plane, the set of couples $(P, T)$ which are solutions of the system:

$$T_l = T = T_v, \quad P_l = P = P_v, \quad \frac{\mu_l(T_l, P_l)}{T_l} = \frac{\mu_v(T_v, P_v)}{T_v} \Leftrightarrow \mu_l(T, P) = \mu_v(T, P), \tag{3}$$

represents the so-called saturation curves for which the two phases co-exist in a stable manner. For any couple $(P, T)$ which is not solution of (3), only one of the two phases is stable (i.e. the other one tends to vanish). When considering Stiffened Gas EOS in the pressure-temperature plane, the chemical potential reads:

$$\begin{aligned}
\frac{\mu_k(T_k, P_k)}{T_k} &= \gamma_k C_{V,k} - s_k(T_k, P_k), \\
s_k(T_k, P_k) &= s_{k,0} + \gamma_k C_{V,k} \ln(C_{V,k} T_k) - (\gamma_k - 1) C_{V,k} \ln\left(\frac{P_k + P_{inf,k}}{\gamma_k - 1}\right),
\end{aligned} \tag{4}$$

where $\gamma_k > 1$, $C_{V,k}$ and $P_{inf,k}$ are constant. We can exhibit explicitly the saturation curve for the temperature with respect to the pressure. It is defined only if $\gamma_v C_{V,v} \neq \gamma_l C_{V,l}$ and reads:

$$T_{sat}(P) = e^{\left(\frac{\beta_l - \beta_v + \gamma_v C_{V,v} - \gamma_l C_{V,l}}{\gamma_v C_{V,v} - \gamma_l C_{V,l}}\right)} \left(\frac{(P+P_{inf,v})^{(C_{V,v}(\gamma_v - 1))}}{(P+P_{inf,l})^{(C_{V,l}(\gamma_l - 1))}}\right)^{\left(\frac{1}{\gamma_v C_{V,v} - \gamma_l C_{V,l}}\right)}, \tag{5}$$

where $\beta_k = s_{k,0} + \gamma_k C_{V,k} \ln(C_{V,k}) + (\gamma_k - 1) C_{V,k} \ln(\gamma_k - 1)$ are the constant parts of the entropies $s_k(P_k, T_k)$. The saturation curve for the pressure with respect to the temperature can not be written explicitly.

## 3 Discretization Scheme

The overall scheme is based on a fractional step method [15]. We first account for the convection terms, which corresponds to system (1) with $\Gamma_v(W) = D_v(W) = \psi_v(W) = \phi_v(W) = 0$. In the sequel, this step is achieved using the relaxation scheme described in [2]. It is not recalled here and the convergence curves obtained for analytical test cases can be found in [12]. This scheme has proven to be accurate and has shown good capability to treat small values of $\alpha_k$, which are very important features for industrial simulations.

In the second step of the algorithm, source terms $\Gamma_v(W)$, $D_v(W)$, $\psi_v(W)$ and $\phi_v(W)$ are accounted for successively through the corresponding ODE system with the time step $\Delta t$ fixed by the convection scheme. The corresponding schemes for $D_v(W)$, $\psi_v(W)$ and $\phi_v(W)$ are implicit and are described in [7, 11, 12]. For each source term, analytical solutions can also be found in these references. We focus here on the scheme that handles the mass transfer term $\Gamma_v(W)$. Paying attention to the properties of mass, momentum and total energy conservation for the mixture, the ODE system for the mass transfer obtained from system (1) is:

$$
\begin{cases}
\partial_t \left( \alpha_v \rho_v \right) = \Gamma_v, \\
\partial_t \left( \alpha_v \rho_v U_v \right) = \overline{U}_{int} \Gamma_v, \\
\partial_t \left( \alpha_v E_v \right) = \overline{H}_{int} \Gamma_v, \\
\partial_t \alpha_v = \partial_t (m) = \partial_t (mU) = \partial_t (E) = 0.
\end{cases}
\tag{6}
$$

Starting from an initial value $W^n$ of $W$ at time $t^n$, we describe now how the value $W^{n+1}$ is computed at time $t^{n+1} = t^n + \Delta t$.

We first approximate system (6) by taking $\tau_g = \tau_g(W(t = 0))$. The solutions for the statistical fractions are obvious: $\alpha_k(t) = \alpha_k(t = 0)$, which enables to write the source term $\Gamma_v$ as a function of the densities and the specific internal energies: $\Gamma_v = \tilde{\Gamma}_v(\rho_l, e_l, \rho_v, e_v)$. Moreover, thanks to the closures for $\overline{U}_{int}$ and $\overline{H}_{int}$, the internal energies remain constant:

$$
\partial_t (m_v e_v) = 0 \quad \text{and} \quad \partial_t (m_l e_l) = 0.
\tag{7}
$$

If we now use the fact that the mass of the mixture is conserved, $\Gamma_v$ can be written as a function of $m_l(t)$ (or $m_v(t)$) and the initial conditions:

$$
\begin{aligned}
\Gamma_v &= \tilde{\Gamma}_v(\rho_l, e_l, \rho_v, e_v) \\
&= \tilde{\Gamma}_v \left( \frac{m_l}{\alpha_l(t=0)}, \frac{(m_l e_l)(t=0)}{m_l}, \frac{(m_l + m_v)(t=0) - m_l}{\alpha_v(t=0)}, \frac{(m_v e_v)(t=0)}{(m_l + m_v)(t=0) - m_l} \right) \\
&= \bar{\Gamma}_v(m_l).
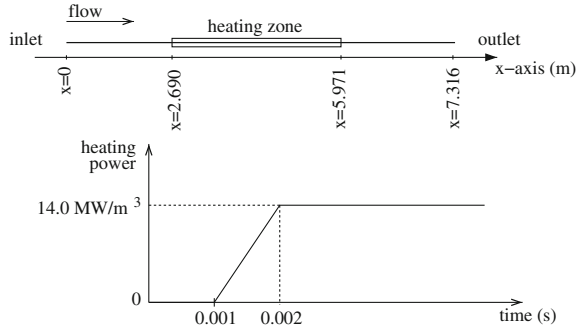\end{aligned}
\tag{8}
$$

A straightforward consequence of this property is that the mass equation (i.e. the first equation of (6)) can be solved independently of the other equations. In general, the source term $\Gamma_v$ can not be explicitely integrated. We thus solve the mass equation using one time-step $\Delta t$ of the Euler implicit scheme:

$$
m_l^{n+1} = m_l^n - \Delta t \, \bar{\Gamma}_v(m_l^{n+1}), \quad \text{with} \quad m_l^{n+1} \in [0, m_l^n + m_v^n].
\tag{9}
$$

The solution $m_l^{n+1}$ at the end of the time step may be computed by a dichotomy algorithm. The function $Y \to \bar{\Gamma}_v(Y)$ is non-linear and might be non smooth. We can state the following result setting $F(Y) = m_l^n - Y - \Delta t \, \bar{\Gamma}_v(Y)$. Since $\Gamma_v$ vanishes for $m_l = 0$ or $m_v = 0$ we obviously have $F(0) = m_l^n$ and $F(m_l^n + m_v^n) = -m_v^n$. If we assume that $Y \to \bar{\Gamma}_v(Y)$ is continuous, its form ensures that if the masses $m_k^n$ are positive, then the masses $m_k^{n+1}$ are also positive. Finally, if $F$ is continuous and strictly monotone on $[0, m_l^n + m_v^n]$, there exists a unique solution to (9) in $]0, m_l^n + m_v^n[$. Once the mass $m_l^{n+1}$ has been computed, the term $\bar{\Gamma}_v(m_l^{n+1})$ and the remaining equations can be updated using one step of the implicit Euler scheme:

$$
\begin{cases}
\alpha_l^{n+1} = \alpha_l^n, \quad m_l^{n+1} = m_l^n - \Delta t \, \bar{\Gamma}_v(m_l^{n+1}), \quad m_v^{n+1} = m^n - m_l^{n+1}, \\
(m_l U_l)^{n+1} = (m_l U_l)^n - \Delta t \, \overline{U}_{int}^{n+1} \bar{\Gamma}_v(m_l^{n+1}), \quad (m_v U_v)^{n+1} = mU_l^n - (m_l U_l)^{n+1}, \\
(\alpha_l E_l)^{n+1} = (\alpha_l E_l)^n - \Delta t \, \overline{H}_{int}^{n+1} \bar{\Gamma}_v(m_l^{n+1}), \\
(\alpha_v E_v)^{n+1} = (\alpha_v E_v)^n + (\alpha_l E_l)^n - (\alpha_l E_l)^{n+1}.
\end{cases}
\tag{10}
$$

**Fig. 1** Sketch of the test case:
geometrical domain and time-
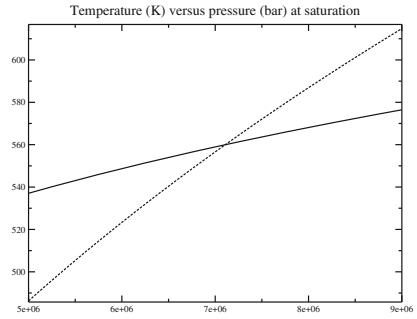schedule of the heating



In fact, the two momentum equations form a $2 \times 2$ linear system whose determinant $\Delta_{G,u}$ is always positive if and only if the partial masses are positive, since: $\Delta_{G,u} = (m_l^n m_v^{n+1} + m_l^{n+1} m_v^n)/2$. Once the velocities $U_k^{n+1}$ are known, the update of the total energies $E_l^{n+1}$ and $E_v^{n+1}$ is straightforward. This scheme is a first-order scheme which ensures the conservation of the total mass, the total momentum and the total energy of the mixture. The positivity of the fractions and the partial masses is ensured.
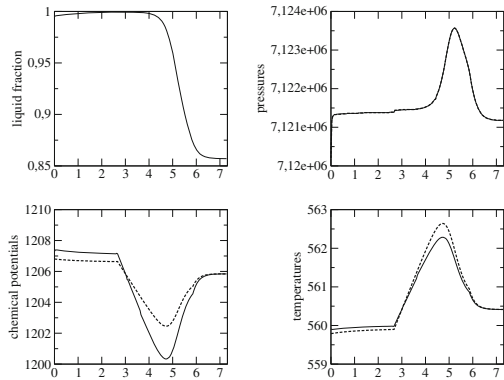
## 4 Heated Saturated Water in a Pipe

The test case is derivated from the OECD/CSNI benchmark problem [14]. It consists in heating saturated water flowing in a one-dimensional pipe. The increase of heat of the fluid leads to vaporization of the water which is advected. The sketch of the case is depicted in Fig. 1. Since we do not account for the head loss in the pipe - as proposed in the OECD/CSNI benchmark problem - we do not need to wait for a stationnary state to be established in the pipe before beginning to heat the fluid. In fact the initial conditions given below already represent a stationnary state. Hence the time schedule of the present case is slightly different.

The initial conditions are chosen at a pressure of $P = 71.0$ bars and a temperature close to the saturation temperature $T = 559.75$ K. They are: $\alpha_l = 0.99$, $\rho_l = 739.8$ kg/m$^3$, $\rho_v = 37.1$ kg/m$^3$, $U_l = U_v = 1.468$ m/s, $P_l = P_v = 71.0$ bars. The EOS parameters are chosen to get these values and to recover the values of the phasic celerities and a temperature saturation-curve (5) close to the real one in the vicinity of the pressure $P = 71.0$ bars and the temperature $T = 559.75$ K. It yields: $C_{V,v} = 1329.45$ J/kg/K, $\gamma_v = 1.257$, $P_{inf,v} = 0$, $s_v^0 = -16274.14$ J/kg/K, $C_{V,l} = 285.14$ J/kg/K, $\gamma_l = 3$, $P_{inf,l} = 2.29 \times 10^8$ bars, $s_l^0 = 0$. The saturation curve is shown on Fig. 2 together with a tabulated saturation curve. The difference is not negligible. Actually, due to the higher slope of the stiffened gas saturation curve, we may underestimate the vapour production. For the inlet boundary-condition the values are the same as the initial condition values. These values provide an equilibrium state since velocities, pressures, temperatures and chemical potentials are equal.
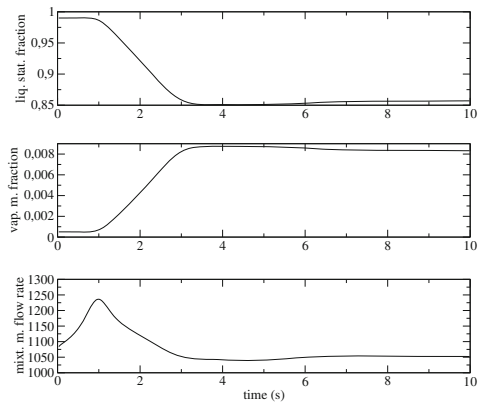
**Fig. 2** Saturation curves for
the temperature with respect
to the pressure on the interval
[ 50, 90 bars]. The *plain
line* represents a reference
saturation curve, whereas the
*dashed line* represents the
saturation curve obtained with
our stiffened gas EOS

**Fig. 3** Thermodynamical
variables along the x-axis at
time $t = 10$ s. The *plain lines*
represent the liquid variables
and the *dashed lines* represent
the vapour variables

**Fig. 4** Liquid statistical
fraction, vapour mass fraction
and mixture mass flow rate at
the outlet of the domain along
the simulation

We are interested in the stationnary state that is reached after 10 s of physical
time. The results obtained with the code presented in the previous sections are given
on Figs. 3 and 4. They correspond to an industrial mesh with 200 uniform cells.
The CFL condition $1/2$ applied to the convection scheme leads to a time step of
$1.5 \times 10^{-5}$ s. The latter is smaller than the time scales which are: $\tau_g = 2.0 \times 10^{-4}$ s
and $\tau_t = 1.0 \times 10^{-4}$ s. The parameter for pressure relaxation is chosen in accordance

with [5]: $K_p = 1.226\,10^{-4}$. Figure 3 represents the thermodynamical variables along the x-axis at time $t = 10\,\text{s}$ and Fig. 4 gives the vapour mass fraction and the mixture mass flow rate at the outlet of the domain for the whole simulation time. It can be noticed that the heating mainly results in the increase of the temperature and that the pressures do not vary a lot. At the outlet of the domain, the liquid fraction starts to evolve at time $t = 1\,\text{s}$, which corresponds to the time necessary for the vapour generated to reach the outlet (the vapour travels at almost 1.4 m/s and there is almost 1.4 m between the downstream edge of the heating zone and the outlet).

# References

1. Baer, M., Nunziato, J.W.: A two-phase mixture theory for the deflagration-to-detonation transition in reactive granular materials. Int. J. Multiph. Flows **12**, 861–889 (1986)
2. Coquel, F., Hérard, J.M., Saleh, K., Seguin, N.: A robust entropy-satisfying finite volume scheme for the isentropic Baer Nunziato model. Math. Model. Numer. Anal. **48**(1), 165–206 (2014)
3. Coquel, F., Hérard, J.M., Saleh, K., Seguin, N.: Two properties of two-velocity two-pressure models for two-phase flows. Comm. Math. Sci. **12**(3), 593–600 (2014)
4. Gallouët, T., Hérard, J.M., Seguin, N.: Numerical modeling of two-phase flows using the two-fluid two-pressure approach. Math. Models Meth. Appl. Sci. **14**(5), 663–700 (2004)
5. Gavrilyuk, S.: The structure of pressure relaxation terms: one velocity case. EDF R&D report H-I83-2014-00276-EN (2014)
6. Hérard, J.M., Hurisse, O.: Computing two-fluid models of compressible water-vapour flows with mass transfer. AIAA paper 2012–2959. http://www.aiaa.org (2012)
7. Hérard, J.M., Hurisse, O.: A fractional step method to compute a class of compressible gas-liquid flows. Comput. Fluids **55**, 57–69 (2012)
8. Hérard, J.M., Liu, Y.: Une approche bifluide statistique de modélisation des écoulements diphasiques à phases compressibles. EDF R&D report H-I81-2013-01162-FR (2013)
9. Ishii, M., Hibiki, T.: Thermo-fluid dynamics of two-phase flow. Springer, New York (2006)
10. Kapila, A., Bdzil, J., Menikoff, R., Son, S., Stewart, D.: Two-phase modelling of ddt in granular materials: reduced equations. Phys. Fluids **13**, 3002–3024 (2001)
11. Liu, Y.: Contribution à la vérification et à la validation d'un modèle diphasique bifluide instationnaire. Ph.D. thesis, Aix Marseille University, Marseille, France. http://tel.archives-ouvertes.fr/tel-00864567. Accessed 11 Sept 2013
12. Morente, A., Hurisse, O., Hérard, J.M.: Vérification d'un code pour les écoulements diphasiques. EDF R&D report H-I83-2013-03283-FR (2013)
13. Saleh, K.: Analyse et simulation par relaxation d'écoulements diphasiques compressibles. Ph.D. thesis, Pierre et Marie Curie University, Paris, France. http://tel.archives-ouvertes.fr/tel-00761099. Accessed 26 Nov 2012
14. Werner, W.: First CSNI numerical benchmark problem. Comparison report. CSNI report 47 (1980)
15. Yanenko, N.N.: Méthodes à pas fractionnaires. Armand Colin (1968)

# Modeling Phase Transition and Metastable Phases

**François James and Hélène Mathis**

**Abstract** We propose a model that describes phase transition including metastable phases present in the van der Waals Equation of State (EoS). We introduce a dynamical system that is able to depict the mass transfer between two phases, for which equilibrium states are both metastable and stable states, including mixtures. The dynamical system is then used as a relaxation source term in a isothermal two-phase model. We use a Finite Volume scheme (FV) that treats the convective part and the source term in a fractional step way. Numerical results illustrate the ability of the model to capture phase transition and metastable states.

## 1 Introduction

Metastable vapor is a gaseous state where the pressure is higher than the saturation pressure. Such states are very unstable and a very small perturbation brings out a droplet of liquid inside the gas. Such a phenomenon can appear at saturated pressure (or at saturated temperature for metastable liquid) for instance inside a nozzle such as a fuel injector or in a cooling circuit of pressurized water reactor. In the last decades considerable research has been devoted to the modeling of two-phase flows with phase transition. However the exact expressions of the mass transfer term are usually unknown (see [2]). In particular, to our knowledge, there is very few literature about the transfer term able to depict metastable states. In [7] and [8] the authors consider a 6-equation model where relaxation to equilibrium is achieved by chemical and pressure relaxation terms whose kinetics are considered infinitely fast.

F. James
MAPMO, Univ. Orléans and CNRS, UMR CNRS 7349, 45067 Orléans Cedex 2, France
e-mail: francois.james@univ-orleans.fr

H. Mathis (✉)
LMJL, Univ. Nantes, 2 rue de la Houssinière, BP 92208, 44322 Nantes Cedex 3, France
e-mail: helene.mathis@univ-nantes.fr

We intend here to provide a new model able to depict phase transition and metastable states with non-infinite relaxation speed. It is based on the use of the van der Waals EoS, that is well-known to depict stable and metastable states below the critical temperature. However this EoS is not valid in the so-called spinodal zone where the pressure is a decreasing function of the density. This leads to instabilities and computational failure and most commonly the pressure is corrected using the Maxwell equal area rule construction to recover a constant pressure. However such a correction removes the metastable regions. We propose transfer terms obtained through an optimization problem of the Helmholtz free energy of the two-phase system. For sake of simplicity we assume the system to be isothermal. We obtain a dynamical system that is able to depict mass transfer including metastable states and that dissipates the total Helmholtz free energy. The equilibria of the dynamical system are either stable states or metastable states or a mixture state that satisfies the pressures and chemical potentials equalities. This dynamical system is used as a transfer term in a isothermal two-phase model in the spirit of [6] and [1]. We use a classical FV scheme that treats the convective and the source terms in a splitting approach.

Section 2 is devoted to the thermodynamics of a binary mixture and presents the major properties of the van der Waals EoS. Section 3 is devoted to the construction of the dynamical system based on results of the previous section. In particular we show that metastable states are attractors of the dynamical system. In Sect. 4 we briefly present the splitting FV scheme we use and give numerical results where some metastable vapor appears.

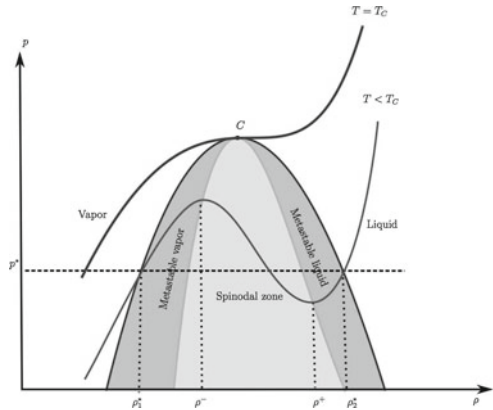## 2 Thermodynamics and van der Waals Equation of State

In this Section we first recall the thermodynamics theory for a single isothermal fluid and introduce the different potentials of the van der Waals EoS, then we state the mathematical framework for the thermodynamics of immiscible binary mixtures.

### 2.1 Thermodynamics of a Single Phase

Consider a single fluid of mass $M > 0$ occupying a volume $V > 0$. At constant temperature if the fluid is homogeneous and at rest, its behavior is entirely described by the Helmholtz free energy function $E(M, V)$ which belongs to $C^2(\mathbb{R}_+ \times \mathbb{R}_+)$ and is positively homogeneous of degree 1 (PH1). Thus, at fixed volume $V$, one can introduce the specific Helmholtz free energy $f$ and the specific energy $e$ that are functions of the density $\rho = M/V$

$$f(\rho) = E(\rho, 1), \qquad \rho e(\rho) = E(\rho, 1). \tag{1}$$

**Fig. 1** Phase diagram for the van der Waals EoS in the $(p, \rho)$ plan. Below the critical point $C$, the *isotherm curve* decreases in the so-called spinodal zone delimited by the densities $\rho^- < \rho^+$. In that area the isotherm is commonly replaced by an *horizontal segment* that coincides with the isobaric line at constant pressure $p^*$. Such a construction defines the two densities $\rho_1^*$ and $\rho_2^*$



We also introduce the pressure $p$ and the chemical potential $\mu$ that are partial derivatives of the free energy $E$, respectively with respect to $V$ and $\rho$. By homogeneity these can be written as functions of $\rho$ solely:

$$p(\rho) = -\partial_V E(\rho, 1), \qquad \mu(\rho) = \partial_M E(\rho, 1). \tag{2}$$

Again thanks to the homogeneity of the energy function, one has

$$f(\rho) = \rho\mu(\rho) - p(\rho), \qquad f'(\rho) = \mu(\rho). \tag{3}$$

Stable pure phases are characterized by a convex energy function, which leads to a nondecreasing pressure law. We consider a classical example of a fluid that may experience phase transitions, namely the van der Waals monoatomic fluid. At fixed temperature $T$ its Helmholtz free energy is given by

$$E(M, V) = -\frac{aM^2}{V} + RT\left(M \log \frac{M}{V - Mb} - M\right), \tag{4}$$

where $R$ stands for the perfect gas constant and $a$ and $b$ are positive constants, $a$ accounts for binary interactions and $b$ is a specific covolume. Below a critical temperature $T_C$ the associated pressure is not monotone with respect to (wrt) the density (see Fig. 1): in a region called the spinodal zone, delimited for a given temperature by the densities $\rho^- < \rho^+$, the pressure decreases wrt the density, thus leading to unstable states. In that region the isotherm is commonly replaced by the Maxwell area rule in order to recover that phase transition happens at constant pressure and constant chemical potential. However this construction removes admissible regions where the pressure law is still nondecreasing. Such regions are called the metastable regions (see Fig. 1).

## *2.2 Equilibrium of a Two-Phase Mixture*

We consider now two immiscible phases of a same pure fluid of total mass $M$ and volume $V$ at a fixed subcritical temperature. Each phase $i = 1, 2$, is depicted by its mass $M_i \geq 0$ and its volume $V_i \geq 0$. We assume that both phases are characterized by the same van der Waals extensive Helmholtz free energy $E$ function of $M_i$ and $V_i$, given by (4). By the conservation of mass, the mass of the binary system is $M = M_1 + M_2$ and immiscibility implies $V = V_1 + V_2$.

According to the second principle of thermodynamics (see [4]), for fixed mass $M$ and volume $V$ the stable equilibrium states of the system are the solutions to the constrained optimization problem

$$\inf\{E(M_1, V_1) + E(M_2, V_2)|\ V_1 + V_2 = V,\ M_1 + M_2 = M\},$$

which can be rewritten using (1) in term of the specific Helmholtz free energy at fixed density $\rho$:

$$\inf\{\alpha_1 f(\rho_1) + \alpha_2 f(\rho_2)|\ \alpha_1 + \alpha_2 = 1,\ \alpha_1 \rho_1 + \alpha_2 \rho_2 = \rho\}, \tag{5}$$

where $\alpha_i = V_i/V \in [0, 1]$ denotes the volume fraction and $\rho_i = M_i/V_i$ is the density of the phase $i = 1, 2$. In the sequel the fractions $\alpha_i$ are written as functions of $\rho, \rho_1$ and $\rho_2$ such that $\alpha_1(\rho, \rho_1, \rho_2) = (\rho - \rho_2)/(\rho_1 - \rho_2)$ and $\alpha_2(\rho, \rho_1, \rho_2) = 1 - \alpha_1(\rho, \rho_1, \rho_2)$.

Note that $\alpha_1$ and $\alpha_2$ are simultaneously non zero if and only if $\rho_1 \neq \rho_2$. In that case we shall always assume without loss of generality that $\rho_1 < \rho_2$ and $\rho \in [\rho_1, \rho_2]$. The total Helmholtz free energy $F : \mathbb{R}^3_+ \to \mathbb{R}$ of the binary system is given by

$$F(\rho, \rho_1, \rho_2) = \alpha_1(\rho, \rho_1, \rho_2)f(\rho_1) + \alpha_2(\rho, \rho_1, \rho_2)f(\rho_2). \tag{6}$$

Depending on the saturation of the volume fractions, one can characterize the equilibria of the optimization problem (5).

**Proposition 1**

1. *Pure states*: *if $\alpha_1 = 0$ (resp. $\alpha_2 = 0$) then only the phase 2 (resp. 1) is stable.*
2. *Mixture*: *if $\alpha_1 \alpha_2 \neq 0$, then the equilibrium state is characterized by one of the following equivalent properties*

   a. *equality of the chemical potentials and the pressures*

   $$\mu(\rho_1) = \mu(\rho_2) = \mu^*, \quad p(\rho_1) = p(\rho_2) = p^*, \tag{7}$$

   b. *Maxwell area rule on the chemical potential*

   $$\int_0^1 \mu(\rho_2 + t(\rho_1 - \rho_2))dt = \mu(\rho_1) = \mu(\rho_2) = \mu^*. \tag{8}$$

*The densities such that (7) or (8) holds are denoted $\rho_1^*$ and $\rho_2^*$, see Fig. 1.*

The most important consequence of this result is that in the metastable zones there are two possible equilibrium states corresponding to a pure metastable state and a stable mixture state. Hence the EoS at equilibrium is not single-valued. The difference between stable and metastable states lies in their dynamical behaviour with respect to perturbations, see [5].

## 3 Dynamical System and Phase Transition

We turn now to the study of dynamical stability of equilibrium states. First we address the homogenous case, introducing a dynamical system for which the equilibria are both stable and metastable states as well as states in the spinodal area such that (7) or (8) is satisfied. Next the dynamical system is plugged as a relaxation source terms in a isothermal two-fluid model. Some properties of the full model are given: hyperbolicity, existence of an energy function that decreases in time.

### 3.1 Dynamical System

Assuming that $\rho$, $\rho_1$ and $\rho_2$ are only time-dependent, we introduce the following dynamical system, which derives from the optimality conditions of Proposition 1:

$$
\begin{aligned}
\dot{\rho} &= 0, \\
\dot{\rho}_1 &= -(\rho - \rho_1)(\rho - \rho_2)\left(\rho_2(\mu(\rho_2) - \mu(\rho_1)) + p(\rho_1) - p(\rho_2)\right), \\
\dot{\rho}_2 &= (\rho - \rho_1)(\rho - \rho_2)\left(\rho_1(\mu(\rho_1) - \mu(\rho_2)) - p(\rho_1) + p(\rho_2)\right).
\end{aligned}
\tag{9}
$$

Straightforward computations show that the total Helmholtz free energy $F$ defined by (6) decreases in time along the solutions of this system. It can also be proved that if $\rho_1(0) < \rho_2(0)$ then $\rho_1(t) < \rho_2(t)$ for all $t > 0$. Hence one can assume without loss of generality that

$$
\rho_1(t) < \rho_2(t), \qquad \rho_1(t) \leq \rho \leq \rho_2(t).
\tag{10}
$$

We focus now on the equilibria which can be reached by the model.

**Theorem 1** *Under assumption* (10)*, the equilibrium states of system* (9) *are*

1. *Pure states: $\alpha_2 = 0$, $\rho = \rho_1$, any $\rho_2 \neq \rho$ (resp. $\alpha_1 = 0$, $\rho = \rho_2$, any $\rho_1 \neq \rho$).*
   *More precisely,*
   a. *if $\rho \notin [\rho^-, \rho^+]$, then the equilibrium is an attractor and corresponds to monophasic stable or metastable states (see Fig. 1),*
   b. *if $\rho \in ]\rho^-, \rho^+[$, then the equilibrium is a repeller and corresponds to non admissible states belonging to the spinodal zone,*

2. *Mixture states: the unique state such that $0 < \alpha_1 < 1$ and relation (7) or (8) is satisfied.*

A remarkable feature of this system is that a perturbation of a pure metastable state involving the other phase leads to a mixture equilibrium state, corresponding to the definition of metastable state [5].

## *3.2 The Isothermal Model*

The previous dynamical system (9) is now coupled with a modified version of the isothermal two-phase model proposed in [1] (see also [6]). Now the unknowns $\rho, \rho_1, \rho_2$ are functions of time $t$ and space $x$. The model admits a mixture pressure $\alpha_1 p(\rho_1) + \alpha_2 p(\rho_2)$ and one velocity $u$ for both phases. It reads

$$\partial_t \rho + \partial_x (\rho u) = \frac{1}{\varepsilon} \dot{\rho} = 0,$$

$$\partial_t \rho_i + \partial_x (\rho_i u) = \frac{1}{\varepsilon} \dot{\rho}_i, \quad i = 1, 2 \qquad (11)$$

$$\partial_t (\rho u) + \partial_x (\rho u^2 + \alpha_1 p(\rho_1) + \alpha_2 p(\rho_2)) = 0,$$

where the source terms are given by the dynamical system (9) and account for mass and mechanical transfer. The parameter $\varepsilon > 0$ is a relaxation parameter that represents the relaxation time to reach thermodynamical equilibrium.

If $\rho, \rho_1, \rho_2 \notin [\rho^-, \rho^+]$ then the convective part of the model (11) is hyperbolic with the eigenvalues

$$\lambda_1 = u - c, \quad , \lambda_2 = \lambda_3 = u, \quad \lambda_4 = u + c, \qquad (12)$$

where the sound velocity is $c = \sqrt{\frac{1}{\rho} (\alpha_1 \rho_1 p'(\rho_1) + \alpha_2 \rho_2 p'(\rho_2))}$.

**Proposition 2** *The function $\mathscr{E}(\rho, \rho_1, \rho_2, u) = \dfrac{\rho u^2}{2} + \alpha_1 f(\rho_1) + \alpha_2 f(\rho_2)$, satisfies the following equation*

$$\partial_t(\mathscr{E}) + \partial_x(u(\mathscr{E} + \alpha_1 p(\rho_1) + \alpha_2 p(\rho_2)) = (\partial_{\rho_1} F)\dot{\rho}_1 + (\partial_{\rho_2} F)\dot{\rho}_2 \leq 0. \qquad (13)$$

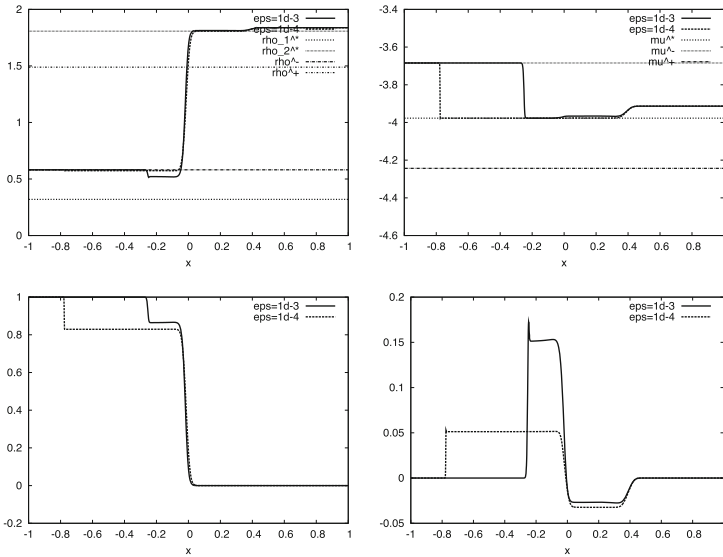Note that $\mathscr{E}$ is not an entropy of the system since $f$ is a non-convex function of the density.

**Fig. 2** First line: density $\rho$ (*left*) chemical potential $\mu$ (*right*). Second line: fraction $\alpha_1$ (*left*), velocity $u$ (*left*)

## 4 Numerical Illustration

We present here numerical results that assess the ability of the model to capture phase transition including metastable states. We use a standard FV method to approximate the Cauchy problem

$$\partial_t W + \partial_x F(W) = S(W), \qquad W(0, x) = W_0(x), \ x \in \mathbb{R}, \tag{14}$$

where $W = (\rho, \rho_1, \rho_2, \rho u)^T$, $F(W) = (\rho u, \rho_1 u, \rho_2 u, \rho u^2 + \alpha_1 p(\rho_1) + \alpha_2 p(\rho_2))^T$, and $S(W) = (0, \frac{1}{\varepsilon}\dot{\rho}_1, \frac{1}{\varepsilon}\dot{\rho}_2, 0)^T$. Neumann boundary conditions are implemented. We use a fractional step approach. We denote $\Delta t$ the time step and $\Delta x$ the length of the cell $(x_{i-1/2}, x_{i+1/2})$ on the regular 1D-mesh. Let $W^n$ be the FV approximation at time $t^n = n\Delta t$, $n \in \mathbb{N}$. The first step corresponds to the approximation of the convective part which provides the solution $W^{n,-}$ at time $t^{n,-}$. It is treated by a classical Rusanov scheme. The second step is the approximation of the source terms (relaxation), at this stage we merely use an explicit Euler method.

We consider the van der Waals equation (4) in its reduced form, see [5], with $R = 8/3$, $a = 3$ and $b = 1/3$ at constant subcritical temperature $T = 0.85$. The extrema of the isotherm curve are $\rho^- = 0.581079$ and $\rho^+ = 1.488804$. The Maxwell construction on the chemical potential defines the densities $\rho_1^* = 0.319729$ and $\rho_2^* = 1.807140$ such that $\mu(\rho_1^*) = \mu(\rho_2^*) = -3.977178$ and $p(\rho_1^*) = p(\rho_2^*) = 0.504492$. If the Riemann problem consists in an initial constant pressure and constant chemical potential state, the numerical scheme preserves this state exactly as it is expected.

Another test case consists in an initial constant pressure state which is subjected to a disequilibrium in chemical potential. The initial data are $\rho_L = \rho_{1,L} = \rho^-$, $\rho_{2,L} = 1.6$, $\rho_R = \rho_{2,R} = 1.837840$, $\rho_{1,R} = 0.2$ and $u_L = u_R = 0$. The discontinuity is applied at $x = 0$ in the domain $[-1, 1]$. The mesh contains 2000 cells and the time of computation is $t = 0.2$. Note that $\rho_{2,L}$ belongs to the metastable liquid region and $\rho_{2,R}$ belongs to the pure liquid region such that $p(\rho_{2,R}) = p(\rho^-) = p(\rho_{1,L})$ and $\rho_{1,R}$ belongs to the pure gaseous region. Figure 2 presents the results for $\varepsilon = 10^{-3}$ and $\varepsilon = 10^{-4}$. The main feature to notice here is that the relaxation approximation introduces a mixture zone on both sides of the interface, which eventually fills the whole domain as time evolves. The velocity of the waves delimiting this zone is faster when $\varepsilon$ goes to 0. Within this zone, there are variations of the velocity, which remains compressive ($u > 0$ for $x < 0$, $u < 0$ for $x > 0$).

## 5 Conclusion and Prospects

The first tests with this model show that it is able to cope with phase transitions with metastable states using a van der Waals EoS, even though the behaviour of the mixture zone around the interface has to be investigated in more details. The explicit treatment of the stiff relaxation term enforces tough constraints on the time step, we plan to implement a semi-implicit scheme in the spirit of [3]. Finally, we attend to include temperature dependance to obtain a fully heat, mass and mechanical transfer model in order to compare our results to those of [7] and [8].

## References

1. Ambroso, A., Chalons, C., Coquel, F., Galié, T.: Relaxation and numerical approximation of a two-fluid two-pressure diphasic model. M2AN. Math. Model. Numer. Anal. **43**(6), 1063–1097 (2009)
2. Drew, D.: Mathematical modeling of two-phase flow. Ann. Rev. Fluid Mech. **15**, 261–291 (1983)
3. Gallouët, T., Hérard, J.M., Seguin, N.: Numerical modeling of two-phase flows using the two-fluid two-pressure approach. Math. Models Methods Appl. Sci. **14**(5), 663–700 (2004)
4. Gibbs, J.W.: The Collected Works of J. Willard Gibbs, Vol. I: Thermodynamics. Yale University Press, New Haven (1948)
5. Landau, L., Lifschitz, E.: A Course of Theoretical Physics, Vol. 5, Statistical Physics, Chap. 8. Pergamon Press, New York (1969)
6. Murrone, A., Guillard, H.: A five equation reduced model for compressible two phase flow problems. J. Comput. Phys. **202**(2), 664–698 (2005)
7. Saurel, R., Petitpas, F., Abgrall, R.: Modelling phase transition in metastable liquids: application to cavitating and flashing flows. J. Fluid Mech. **607**, 313–350 (2008)
8. Zein, A., Hantke, M., Warnecke, G.: Modeling phase transition for compressible two-phase flows applied to metastable liquids. J. Comp. Phys. **229**, 1964–2998 (2010)

# Almost Parallel Flows in Porous Media

**Alaa Armiti-Juber and Christian Rohde**

**Abstract**   This paper considers a reduced two-phase model for mostly unidirectional porous media flows. It is a nonlinear conservation law, in which velocity depends nonlocally on the unknown saturation, see [6]. We aim to construct and analyze a finite-volume scheme for the model. For the analysis, the main difficulty is the reduced regularity in the transverse velocity component. The upwind finite-volume scheme is used to prove the existence of weak solutions of a regularized Cauchy problem in the framework of functions of bounded variations. Then, we consider the limit of vanishing regularization parameter. Numerical examples that analyze the efficiency of the approach are also presented.

## 1 Introduction

The process of fluid displacement in a heterogeneous porous medium by another immiscible fluid belongs to the general field of multiphase flow. Assuming incompressible fluids, constant medium's porosity, and negligible gravity and capillary pressure forces, the standard two-phase flow model, in fractional flow formulation, see [2], is given by:

$$\partial_t S + \mathrm{div}(\mathbf{v} f(S)) = 0, \qquad \text{in } D \times (0, T), \tag{1}$$
$$\mathbf{v} = \lambda(S)\mathbf{K}\nabla p, \quad \mathrm{div}(\mathbf{v}) = 0$$

A. Armiti-Juber (✉) · C. Rohde
Institute für Angewandte Analysis und Numerische Simulation, Universität Stuttgart, Stuttgart, Germany
e-mail: a.armiti@mathematik.uni-stuttgart.de

C. Rohde
e-mail: crohde@mathematik.uni-stuttgart.de

where the spatial domain $D$ is defined as $D := (0, H) \times (0, L)$, with $H, L > 0$ being the domain's width and length, respectively. The unknowns in this model are saturation $S = S(x, y, t) \in [0, 1]$, global pressure $p = p(x, y, t) \in \mathbb{R}$ and the total velocity $\mathbf{v} = \mathbf{v}(x, y, t) \in \mathbb{R}^2$, for $(x, y, t) \in D \times [0, T]$, $T > 0$. The intrinsic permeability tensor $\mathbf{K} = \mathbf{K}(x, y) := \begin{pmatrix} K_1(x, y) & 0 \\ 0 & K_2(x, y) \end{pmatrix}$ with $K_1, K_2$ being the permeabilities in $x$-, $y$-direction, respectively, is given. The total mobility $\lambda : [0, 1] \rightarrow (0, \infty)$ is given by $\lambda(S) := \dfrac{k_{rw}(S)}{\mu_w} + \dfrac{k_{rn}(S)}{\mu_n}$, where $k_{rw}, k_{rn}$ are the wetting, nonwetting relative permeabilities and the constants $\mu_w, \mu_n > 0$ are the wetting-, nonwetting-phase viscosities, respectively. The flux $f : [0, 1] \rightarrow \mathbb{R}$ is also a given function of the unknown $S$.

   Many interesting porous media formations satisfy the geometrical property $H \ll L$, where most of the fluid flows in horizontal direction. In other words, a vertical (transverse) equilibrium assumption can be applied. Yortsos in [6] uses this assumption to derive a reduced model. By setting $\gamma := H/L$ and rescaling the variables $x, y, t, \mathbf{v}, K_1, K_2, p$ in (1) into other corresponding dimensionless ones denoted in the same way, model (1) is transformed into the dimensionless model,

$$\partial_t S^\gamma + \partial_x (u^\gamma f(S^\gamma)) + \frac{1}{\gamma} \partial_y (w^\gamma f(S^\gamma)) = 0,$$
$$\partial_x u^\gamma + \frac{1}{\gamma} \partial_y w^\gamma = 0, \qquad \text{in } D \times (0, T), \quad (2)$$
$$u^\gamma = -\lambda(S^\gamma) \, K_1 \, \partial_x p^\gamma, \quad \gamma w^\gamma = -\lambda(S^\gamma) \, K_2 \, \partial_y p^\gamma$$

where $D = (0, 1) \times (0, 1) \subset \mathbb{R}^2$ and $u^\gamma, w^\gamma$ are the flux velocity components in the main-, transverse-direction (or equivalently, x-, y-direction), respectively. This model still has saturation $S^\gamma$ and pressure $p^\gamma$ as unknowns, but depends on the scaling parameter $\gamma > 0$. So, Yortsos in [6] applies formal asymptotic analysis for $\gamma \rightarrow 0$ and eliminates pressure $p^\gamma$ from the flux velocity $(u^\gamma, w^\gamma)^T$. He derives a nonlinear nonlocal transport equation for the limit saturation $S := \lim_{\gamma \rightarrow 0} S^\gamma$ alone:

$$\partial_t S + \partial_x (u[x, y; S] \, f(S)) + \partial_y (w[x, y; S] \, f(S)) = 0 \quad \text{in } D \times (0, T), \quad (3)$$

$$u[x, y; S] = \frac{\lambda(S) K_1(x, y)}{\int_0^1 \lambda(S) K_1(x, y) dy}, \quad w[x, y; S] = -\partial_x \frac{\int_0^y \lambda(S) K_1(x, q) dq}{\int_0^1 \lambda(S) K_1(x, y) dy}. \quad (4)$$

One can easily check that the velocity vector $(u, w)^T$ defined in (4) is divergence free. We call (3), (4) vertical equilibrium (VE)-model.

   The goal of this paper is to design and analyze a finite volume scheme for the initial value problem (3), (4). For the sake of simplicity in this contribution, we choose the domain $D = \mathbb{R}^2$. The main analytical difficulty is the reduced regularity of the velocity component $w$ due to the existence of the derivative $\partial_x$. Therefore, in Sect. 2, we convolute the velocity components (4) in the x-direction. Then, in Sect. 3, we

suggest a finite-volume scheme for the regularized Cauchy problem (7), (8) with the property of being locally mass conservative. The existence of weak solutions for the regularized Cauchy problem is established in Sect. 4 in the framework of functions of bounded variation. After that, we consider the limit of vanishing regularization parameter. Finally, in Sect. 5 numerical examples, that illustrate the efficiency of the VE-model as a reduced model are presented.

## 2 Regularized Model and Weak Solutions

The conditions (A1), (A2) are supposed to hold throughout the paper:

(A1) The initial saturation $S_0$ satisfies $S_0 \in BV(\mathbb{R}^2)$ and $S_0(x, y) \in [0, 1]$ for almost all $(x, y) \in \mathbb{R}^2$.
(A2) The flux function satisfies $f \in C^1([0, 1])$. The functions $\lambda \in C^1([0, 1])$, $K_1 \in L^\infty(\mathbb{R}^2)$ are strictly positive with $K_1(x, .) \in L^1(\mathbb{R})$ for almost all $x \in \mathbb{R}$.

The transverse velocity component $w$ has reduced regularity because of the $\partial_x$ derivative in its definition together with the high heterogeneity of the medium and the expected low regularity of solutions of the VE-model. Therefore, we suggest to convolute the velocity components in the x-direction. Consider a smooth kernel function $\psi \in C_0^\infty(\mathbb{R})$ such that supp$(\psi) \subset (-1, 1)$. By (A2), $u[., y; Z] \in L^1_{loc}(\mathbb{R})$ holds for almost all $y \in \mathbb{R}$, $Z \in [0, 1]$; therefore,

$$u^\varepsilon[x, y; Z] := (\psi^\varepsilon * u)[x, y; Z] = \int_\mathbb{R} \psi^\varepsilon(x - \eta)u[\eta, y; Z]d\eta, \qquad (5)$$

$$w^\varepsilon[x, y; Z] := -\frac{1}{\varepsilon} \int_{-\infty}^\infty \int_{-\infty}^y (\psi')^\varepsilon(x - \eta)u[\eta, q; Z]dq d\eta \qquad (6)$$

are well-defined for all $y \in \mathbb{R}$, $Z \in [0, 1]$, where $\psi^\varepsilon(x) := \frac{1}{\varepsilon}\psi(\frac{x}{\varepsilon})$. For $\varepsilon > 0$, the regularized Cauchy problem for the unknown $S^\varepsilon$ is given by:

$$\partial_t S^\varepsilon + \partial_x \left(u^\varepsilon[x, y; S^\varepsilon] f(S^\varepsilon)\right) + \partial_y \left(w^\varepsilon[x, y; S^\varepsilon] f(S^\varepsilon)\right) = 0 \quad \text{in } \mathbb{R}^2 \times (0, T), \quad (7)$$

$$S^\varepsilon(x, y, 0) = S_0(x, y) \quad \text{in } \mathbb{R}^2. \qquad (8)$$

The regularized velocity satisfies also div$(u^\varepsilon[x, y; S^\varepsilon], w^\varepsilon[x, y; S^\varepsilon])^T = 0$.

**Definition 1** For $\varepsilon > 0$, a function $S^\varepsilon = S^\varepsilon(x, y, t)$, with $S^\varepsilon(., ., t) \in BV(\mathbb{R}^2)$ $\forall t \in (0, T)$ and $S^\varepsilon(x, y, t) \in [0, 1]$ for almost all $(x, y, t) \in \mathbb{R}^2 \times (0, T)$ is called a weak solution of the regularized Cauchy problem (7), (8) if

$$\int_0^T \int_{\mathbb{R}^2} S^\varepsilon(x, y, t) \partial_t \phi(x, y, t) dx dy dt + \int_{\mathbb{R}^2} S_0(x, y) \phi(x, y, 0) dx dy$$

$$= - \int_0^T \int_{\mathbb{R}^2} \left( u^\varepsilon[x, y; S^\varepsilon] f(S^\varepsilon) \partial_x \phi(x, y, t) + w^\varepsilon[x, y; S^\varepsilon] f(S^\varepsilon) \partial_y \phi(x, y, t) \right) dx dy dt$$

holds for every function $\phi \in C_0^\infty(\mathbb{R}^2 \times [0, T))$.

## 3 The Finite-Volume Scheme

For $\varepsilon > 0$, we construct a finite-volume scheme for the regularized Cauchy problem (7), (8). For $h > 0$, consider the uniform Cartesian grid

$$\mathscr{T} = \{T_{i,j} = [(i - 1/2)h, (i + 1/2)h) \times [(j - 1/2)h, (j + 1/2)h) \mid (i, j) \in \mathbb{Z}^2)\},$$

with $|T_{i,j}| := h^2$. The set of edges of the cell $T_{i,j}$ is denoted by $\{E_l \mid l \in \theta_{i,j}\}$ for $\theta_{i,j} := \{(i - \frac{1}{2}, j), (i + \frac{1}{2}, j), (i, j - \frac{1}{2}), (i, j + \frac{1}{2})\}$. We also define the set of neighbor cells of $T_{i,j}$ as $\{T_{(i,j)_l} \mid l \in \theta_{i,j}\}$, where $T_{(i,j)_l}$ is the neighbor cell to $T_{i,j}$ with the common edge $E_l$. Then, for $\triangle t > 0$, the *Finite-Volume Scheme* in $T_{i,j}$ is given by:

$$\begin{aligned} S_{i,j}^{\varepsilon,n+1} &= S_{i,j}^{\varepsilon,n} - \frac{\triangle t}{h} \sum_{l \in \theta_{i,j}} \mathscr{F}_l^\varepsilon(S_{i,j}^{\varepsilon,n}, S_{(i,j)_l}^{\varepsilon,n}), \\ S_{i,j}^0 &= \frac{1}{h^2} \int_{T_{i,j}} S_0(x, y) dx dy, \end{aligned} \tag{9}$$

where $S_{(i,j)_l}^{\varepsilon,n}$ is the saturation in the cell $T_{(i,j)_l}$ and $\mathscr{F}_l^\varepsilon$ is an upwind numerical flux function. For any $P, Q \in [0, 1]$, it is defined by

$$\mathscr{F}_l^\varepsilon(P, Q) := \max\{\mathbf{n}_l \cdot \mathbf{v}_l^{\varepsilon,n}, 0\} f(P) + \min\{\mathbf{n}_l \cdot \mathbf{v}_l^{\varepsilon,n}, 0\} f(Q), \tag{10}$$

where, $\mathbf{n}_l$ is the outer normal to the edge $E_l$ of $T_{i,j}$ and $\mathbf{v}_l^{\varepsilon,n} = (u_l^{\varepsilon,n}, w_l^{\varepsilon,n})^T$ is the discrete velocity vector. At the edge $E_{i+\frac{1}{2},j}$, we define:

$$u_{i+\frac{1}{2},j}^{\varepsilon,n} := h \sum_{r=-\infty}^{\infty} \psi_{i+\frac{1}{2}-r}^\varepsilon u_{r,j}^n, \qquad w_{i+\frac{1}{2},j}^{\varepsilon,n} := 0, \tag{11}$$

for $\psi_{i+\frac{1}{2}-r}^\varepsilon = \frac{1}{2}(\psi_{i-r}^\varepsilon + \psi_{i+1-r}^\varepsilon)$, and $\psi_{i-r}^\varepsilon$ is the averaged integral of the Kernel $\psi^\varepsilon$ in the sub-interval $[(i - \frac{1}{2} - r)h, (i + \frac{1}{2} - r)h]$. We choose the mid-point quadrature formula to discretize the $y$-integral and the centered difference quotient to discretize $(\psi')^\varepsilon$. Then, at the edge $E_{i,j+\frac{1}{2}}$, we define:

$$w_{i,j+\frac{1}{2}}^{\varepsilon,n} := h \sum_{r=-\infty}^{\infty} \sum_{k=-\infty}^{j} \left( \psi_{i+\frac{1}{2}-r}^\varepsilon - \psi_{i-\frac{1}{2}-r}^\varepsilon \right) u_{r,k}^n, \qquad u_{i,j+\frac{1}{2}}^{\varepsilon,n} := 0. \tag{12}$$

Using the finite-volume scheme (9), we introduce the discrete solution

$$S_h^\varepsilon(x, y, t) := S_{i,j}^{\varepsilon,n} \quad \forall (x, y) \in T_{i,j}, \ t \in [t^n, t^{n+1}).$$

It is straightforward to prove the following properties:

**Lemma 1** (Mass-Conservation) *If the numerical flux function is defined as in* (10)*, then the finite-volume scheme is mass conservative. i.e.,*

$$\sum_{i,j} S_{i,j}^{\varepsilon,n+1} = \sum_{i,j} S_{i,j}^{\varepsilon,n}, \quad n = 0, 1, 2, ... N_T. \tag{13}$$

**Lemma 2** (Lipschitz-Continuity) *For all* $(Z_1, Z_2)$*,* $(Q_1, Q_2) \in (0, 1)^2$*, there exists a constant* $C = C \left( \sup_{S \in [0,1]} \|u[., .; S]\|_{L^\infty(\mathbb{R}^2)}, \|f'\|_{L^\infty([0,1])}, \|(\psi')^\varepsilon\|_{L^\infty(\mathbb{R})} \right) > 0$*, such that the following properties are satisfied*

$$|\mathscr{F}_{i+\frac{1}{2},j}^\varepsilon(Q_1, Q_2) - \mathscr{F}_{i-\frac{1}{2},j}^\varepsilon(Z_1, Z_2)| \le C \left( |Q_1 - Z_1| + |Q_2 - Z_2| + \frac{h}{\varepsilon} \right),$$

$$|\mathscr{F}_{i,j+\frac{1}{2}}^\varepsilon(Q_1, Q_2) - \mathscr{F}_{i,j+\frac{1}{2}}^\varepsilon(Z_1, Z_2)| \le \frac{C}{\varepsilon} \left( |Q_1 - Z_1| + |Q_2 - Z_2| + h \right).$$

**Lemma 3** (Incompressibility) *If the discrete modified velocity* $v_l^{\varepsilon,n} = (u_l^{\varepsilon,n}, w_l^{\varepsilon,n})^T$*,* $l \in \theta_{i,j}$*, is defined as in* (11)*,* (12)*, then* $\sum_{l \in \theta_{i,j}} \boldsymbol{n}_l \cdot \boldsymbol{v}_l^{\varepsilon,n} = 0$*.*

We define the map $G : [0, 1]^5 \to \mathbb{R}$ such that,

$$G_{i,j}^n = G(S_{i,j}^{\varepsilon,n}, S_{i+1,j}^{\varepsilon,n}, S_{i-1,j}^{\varepsilon,n}, S_{i,j+1}^{\varepsilon,n}, S_{i,j-1}^{\varepsilon,n}) := S_{i,j}^{\varepsilon,n} - \frac{\Delta t}{h} \sum_{l \in \theta_{i,j}} \mathscr{F}_l^\varepsilon(S_{i,j}^{\varepsilon,n}, S_{(i,j)_l}^{\varepsilon,n}). \tag{14}$$

Then, Eq. (9) can be rewritten as follows:

$$S_{i,j}^{\varepsilon,n+1} = G(S_{i,j}^{\varepsilon,n}, S_{i+1,j}^{\varepsilon,n}, S_{i-1,j}^{\varepsilon,n}, S_{i,j+1}^{\varepsilon,n}, S_{i,j-1}^{\varepsilon,n}). \tag{15}$$

By Lemma 3 and for any $\mathbf{S} := (S, S, S, S, S) \in [0, 1]^5$ the map $G$ satisfies the consistency property:

$$G(\mathbf{S}) = S. \tag{16}$$

**Lemma 4** (Monotonicity) *If the CFL-condition*

$$\Delta t \le \frac{h}{2 \max_{i,j} |u_{i+\frac{1}{2},j}^{\varepsilon,n}| \max_{S \in [0,1]} |f'(S)|} + \frac{h}{2 \max_{i,j} |w_{i,j+\frac{1}{2}}^{\varepsilon,n}| \max_{S \in [0,1]} |f'(S)|} \tag{17}$$

*is satisfied, then the map G in* (14) *is increasing with respect to its arguments.*

# 4 Convergence Analysis for $\varepsilon > 0$

In this section we study the convergence of the finite-volume scheme to the regularized Cauchy problem (7), (8) as $h \to 0$. The main theorem is stated as follows:

**Theorem 1** *Assume that the CFL condition* (17) *is satisfied for each $\varepsilon > 0$. Then, there exists a function $S^\varepsilon$ satisfying $S^\varepsilon(.,.,t) \in BV(\mathbb{R}^2)$ for all $t \in (0, T)$ and $S^\varepsilon(x, y, t) \in [0, 1]$ for almost all $(x, y, t) \in \mathbb{R}^2 \times (0, T)$ such that, up to a subsequence, $\{S_h^\varepsilon\}_{h>0}$ converges to $S^\varepsilon$ in $L_{loc}^1(\mathbb{R}^2 \times (0, T))$. Moreover, $S^\varepsilon$ is a weak solution of the Cauchy problem* (7), (8).

A sketch of the proof follows, which is based on a classical BV analysis. Therefore, we prove the following a priori-estimates on the discrete solution $S_h^\varepsilon$.

**Lemma 5** ($L^\infty$-Estimate) *If the conditions in Theorem 1 are satisfied, then the discrete solution satisfies $S_h^\varepsilon(x, y, t) \in [0, 1] \ \forall (x, y, t) \in \mathbb{R}^2 \times (0, T)$.*

*Proof* We prove the upper bound, the lower follows similarly. Define the vectors $\mathbf{S} = (S_1, S_2, S_3, S_4, S_5) := (S_{i,j}^{\varepsilon,n}, S_{i+1,j}^{\varepsilon,n}, S_{i-1,j}^{\varepsilon,n}, S_{i,j+1}^{\varepsilon,n}, S_{i,j-1}^{\varepsilon,n})$, $\mathbf{S}_{max} := (S_k, S_k, S_k, S_k, S_k)$ such that $S_k := \max\{S_i, \ i = 1, 2, ..., 5\}$. Then, Eq. (15), Lemma 4 and Eq. (16), yield:

$$S_{i,j}^{\varepsilon,n+1} = G(\mathbf{S}) \leq G(\mathbf{S}_{max}) = S_{r,k}^{\varepsilon,n} \quad \forall (i, j) \in \mathbb{Z}^2.$$

By induction, we get:

$$\sup\nolimits_{i,j} S_{i,j}^{\varepsilon,n+1} \leq \sup\nolimits_{i,j} S_{i,j}^{\varepsilon,n} \leq \cdots \leq \sup\nolimits_{i,j} S_{i,j}^0 \leq 1 \quad \forall (i, j) \in \mathbb{Z}^2, \ n = 0, 1, ..., N_T.$$
□

**Lemma 6** (BV-Estimate) *If the conditions in Theorem 1 are satisfied, then the discrete solution satisfies $|S_h^\varepsilon(.,.,t)|_{BV(\mathbb{R}^2)} \leq |S_0|_{BV(\mathbb{R}^2)}$ for all $t \in (0, T)$ and $h > 0$.*

*Proof* We prove the statement for variation in x-direction, the variation in y-direction is similar. Define the vectors:

$$\mathbf{S} = (S_1, S_2, S_3, S_4, S_5) := (S_{i+1,j}^{\varepsilon,n}, S_{i+2,j}^{\varepsilon,n}, S_{i,j}^{\varepsilon,n}, S_{i+1,j+1}^{\varepsilon,n}, S_{i+1,j-1}^{\varepsilon,n}),$$
$$\hat{\mathbf{S}} = (\hat{S}_1, \hat{S}_2, \hat{S}_3, \hat{S}_4, \hat{S}_5) := (S_{i,j}^{\varepsilon,n}, S_{i+1,j}^{\varepsilon,n}, S_{i-1,j}^{\varepsilon,n}, S_{i,j+1}^{\varepsilon,n}, S_{i,j-1}^{\varepsilon,n}),$$
$$\bar{\mathbf{S}} = (\bar{S}_1, \bar{S}_2, \bar{S}_3, \bar{S}_4, \bar{S}_5), \quad \text{such that } \bar{S}_i := \max\{S_i, \hat{S}_i\}, \ i = 1, 2, 3, 4, 5.$$

Then using Eq. (15), we write

$$|S_{i+1,j}^{\varepsilon,n+1} - S_{i,j}^{\varepsilon,n+1}| = |G(\mathbf{S}) - G(\hat{\mathbf{S}})|.$$

The definition of $\bar{\mathbf{S}}$ and Lemma 4 give:

$$
\begin{aligned}
|G(\mathbf{S}) - G(\hat{\mathbf{S}})| &\leq |G(\mathbf{S}) - G(\bar{\mathbf{S}})| + |G(\bar{\mathbf{S}}) - G(\hat{\mathbf{S}})| = (G(\bar{\mathbf{S}}) - G(\mathbf{S})) + (G(\bar{\mathbf{S}}) - G(\hat{\mathbf{S}})) \\
&= (\bar{S}_1 - S_{i+1,j}^{\varepsilon,n}) + (\bar{S}_1 - \hat{S}_1) + [(G(\bar{\mathbf{S}}) - \bar{S}_1) - (G(\mathbf{S}) - S_{i+1,j}^{\varepsilon,n})] \\
&\quad + [(G(\bar{\mathbf{S}}) - \bar{S}_1) - (G(\hat{S}) - \hat{S}_1)] \\
&= |S_{i+1,j}^{\varepsilon,n} - S_{i,j}^{\varepsilon,n}| + [(G(\bar{\mathbf{S}}) - \bar{S}_1) - (G(\mathbf{S}) - S_{i+1,j}^{\varepsilon,n})] \\
&\quad + [(G(\bar{\mathbf{S}}) - \bar{S}_1) - (G(\hat{S}) - \hat{S}_1)].
\end{aligned}
$$

Then, Lemma 1 and the induction assumption yield:

$$
\sum_{i,j} |S_{i+1,j}^{\varepsilon,n+1} - S_{i,j}^{\varepsilon,n+1}| \leq \sum_{i,j} |S_{i+1,j}^{\varepsilon,n} - S_{i,j}^{\varepsilon,n}| \leq \cdots \leq \sum_{i,j} |S_{i+1,j}^0 - S_{i,j}^0|.
$$

The estimate now follows using

$$
|S_h^\varepsilon(.,.,t)|_{BV(\mathbb{R}^2)} = h \sum_{i,j} \left( |S_{i+1,j}^{\varepsilon,n} - S_{i,j}^{\varepsilon,n}| + |S_{i,j+1}^{\varepsilon,n} - S_{i,j}^{\varepsilon,n}| \right)
$$

$$
h \sum_{i,j} \left( |S_{i+1,j}^0 - S_{i,j}^0| + |S_{i,j+1}^0 - S_{i,j}^0| \right) \leq |S_0|_{BV(\mathbb{R}^2)}. \tag{18}
$$

$\square$

**Lemma 7** (Time-Lipschitz Estimate) *If the conditions in Theorem 1 are satisfied, then, for all $t_1, t_2 \in (0, T)$ there exists a constant $C > 0$ such that,*

$$
\|S_h^\varepsilon(.,.,t_1) - S_h^\varepsilon(.,.,t_2)\|_{L^1(\mathbb{R}^2)} \leq \frac{C}{\varepsilon} \left( |t_1 - t_2| + h^3 \right).
$$

*Proof* Assume without loss of generality that $t_1 > t_2$ with $t_1 = m\triangle t$, $t_2 = n\triangle t$ for $m, n \in \mathbb{N}^+$. Then, (9), Lemma 2, Lemma 6 and the second part of (18) yield:

$$
\sum_{(i,j)\in\mathbb{Z}^2} |S_{i,j}^{\varepsilon,m} - S_{i,j}^{\varepsilon,n}| \leq \sum_{k=n+1}^{m} \sum_{(i,j)\in\mathbb{Z}^2} |S_{i,j}^{\varepsilon,k} - S_{i,j}^{\varepsilon,k-1}| \leq \frac{C}{\varepsilon} \left( \frac{\triangle t}{h^2}(m - n)|S_0|_{BV(\mathbb{R}^2)} + h \right).
$$

$\square$

*Proof of Theorem 1* The uniform estimates in Lemmas 5, 6 together with Kolmogorov Compactness theorem imply that for each $t \in [0, T]$, there exists a function $S^\varepsilon(.,.,t)$ such that, up to a subsequence,

$$
\|S_h^\varepsilon(.,.,t) - S^\varepsilon(.,.,t)\|_{L_{loc}^1(\mathbb{R}^2)} \to 0. \tag{19}
$$

as $h \to 0$. The set $\mathscr{P}_m := \{t \in [0, T]| \ t = t^n = n\triangle t_m, \text{ for } n \in \{0, 1, ..., T/\triangle t_m\}\}$ satisfies $\cup_{m\in\mathbb{N}}\mathscr{P}_m$ is dense in $[0, T]$. Hence, using a standard diagonalization

process, there exists a diagonal subsequence of $\{S_h^\varepsilon(., ., t)\}_{h>0}$ such that the convergence (19) is valid for all $t \in \mathscr{P}_m$, $m \in \mathbb{N}$. Now, using Lemma 7, the convergence follows for all $t \in [0, T]$. This convergence and assumption (A2) imply a pointwise almost everywhere convergence of the nonlinear operators $u^\varepsilon[S_h^\varepsilon] f(S_h^\varepsilon)$, $w^\varepsilon[S_h^\varepsilon]$ $f(S_h^\varepsilon)$. Finally, using a Lax-Wendroff type theorem, we show that $S^\varepsilon$ is a weak solution of the regularized Cauchy problem (7), (8).                                          $\square$

It is also possible by the uniform $L^\infty$ and BV estimates to deduce a compactness result on the sequence of weak solutions $\{S^\varepsilon\}_{\varepsilon>0}$ similar to (19). The key point is to prove an $\varepsilon$-independent Lipschitz bound on the numerical flux function, see [1]. The limit $S \in L^\infty \cap BV(\mathbb{R}^2 \times (0, T))$ satisfies $u^\varepsilon[S^\varepsilon] \to u[S]$ pointwise a.e., but up to now it is not proven to have sufficient regularity (in x-direction). Consequently, we have only $w^\varepsilon[S^\varepsilon] \to \bar{w}$ pointwise a.e. for some function $\bar{w}$. As a result, $S$ is not a standard weak solution, but can be interpreted as a measure-valued solution, see e.g., [3].
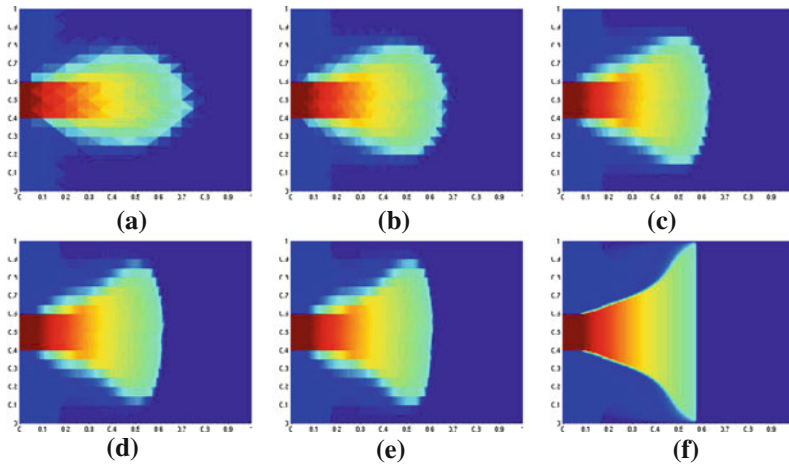
## 5 Numerical Results

In this section, we present numerical results that display the practical efficiency of the VE-model, see [5]. First, the standard two-phase model (1) is considered in five domains $D = (0, 1) \times (0, L)$, $L = 1, 2, 4, 6, 8$. The domains are rescaled into $(0, 1) \times (0, 1)$ and discretized into triangular-grids of $(800 \times L)$ elements, respectively. Using the IMPES-method, see [4], the numerical solutions of (1) are presented in Fig. 1a, e. Then, the VE-model is considered in the scaled domain $(0, 1) \times (0, 1)$ which is discretized into a uniform Cartesian grid of 400 elements. Using the upwind finite-volume scheme, the numerical solution of (3) is presented in Fig. 1f.

A Dirichlet-boundary condition on the inflow boundary $\{(0, y)|y \in (0, 1)\}$ equals to the initial data

$$S_0(0, y) = \begin{cases} 0.1 & : & y \leq \frac{2}{5} \text{ and } y > 1 - \frac{2}{5}, \\ 0.9 & : & \frac{2}{5} < y \leq 1 - \frac{2}{5}. \end{cases}$$

Zero-Neumann conditions on the edges $\{(x, 0) \text{ and } (x, 1)|x \in (0, 1)\}$ are also chosen. The end time $T = 0.3$ is chosen such that a zero-Dirichlet condition on $\{(1, y)|y \in (0, 1)\}$ can be satisfied. Moreover, solutions in Fig. 1 correspond to viscosity ratio $\mu_w/\mu_n = 1/5$.

The numerical solutions of the standard model (1), in Fig. 1a, e, converge to the numerical solution of the VE-model, in Fig. 1f, as the domain parameter $\gamma = (1/L)$ tends to 0. Also, the spreading speed of the wetting front is captured very good by the

**Fig. 1** Solutions of the standard two-phase model versus solution of the VE-model. **a** Sol. of (1) for $\gamma = 1$. **b** Sol. of (1) for $\gamma = 1/2$. **c** Sol. of (1) $\gamma = 1/4$. **d** Sol. of (1) for $\gamma = 1/6$. **e** Sol. of (1) for $\gamma = 1/8$. **f** Sol. of VE-model

reduced VE-model. Moreover, with the used codes and for grids with equal number of elements as in Fig. 1f, the computational time of the solution of the VE-model would be $10^5$ faster than that of the standard model (1).

# References

1. Evje, S., Karlsen, K.H.: Monotone difference approximations of BV solutions to degenerate convection-diffusion equations. SIAM J. Numer. Anal. **37**(6), 1838–1860 (2000)
2. Helmig, R.: Multiphase Flow and Transport Processes in the Subsurface. Springer, New york (1997)
3. Holden, H., Karlsen, K.H., Mitrovic, D.: Zero diffusion-dispersion-smoothing limits for a scalar conservation law with discontinuous flux function. Int. J. Differ. Equ. (2009)
4. Huber, R., Helmig, R.: Multiphase flow in heterogeneous porous media: a classical finite element method versus an implicit pressure-explicit saturation-based mixed finite element-finite volume approach. Int. J. Numer. Meth. Fl. **29**(8), 899–920 (1999)
5. Strohmer, V.: Numerische Analysis von nahezu parallelen Strömungen in porösen Medien (Diploma Thesis). University of Stuttgart (2012)
6. Yortsos, Y.: A theoretical analysis of vertical flow equilibrium. Transp. Porous Media **18**, 107–129 (1995)

# Towards a Stochastic Closure Approach for Large Eddy Simulation

**Th. von Larcher, R. Klein, I. Horenko, P. Metzner, M. Waidmann, D. Igdalov, A. D. Beck, G. J. Gassner and C. -D. Munz**

**Abstract**  We present a  stochastic sub grid scale modeling strategy currently under development for application in Finite Volume Large Eddy Simulation (LES) codes. Our concept is based on the integral conservation laws for mass, momentum and energy of a flow field that are universally valid for arbitrary control volumes. We model the space-time structure of the fluxes to create a discrete formulation. Advanced methods of time series analysis for the data-based construction of stochastic models with inherently non-stationary statistical properties and concepts of information theory for the model discrimination are used to construct stochastic

Dr. P. Metzner was formerly associated with the University of Lugano, Switzerland.

Th. von Larcher (✉) · R. Klein · M. Waidmann
Institute of Mathematics, Freie Universität Berlin, Berlin, Germany
e-mail: larcher@math.fu-berlin.de

R. Klein
e-mail: rupert.klein@math.fu-berlin.de

M. Waidmann
e-mail: waidmann@math.fu-berlin.de

I. Horenko · D. Igdalov · P. Metzner
Institute of Computational Science, Universita della Swizzerà italiana, Lugano, Switzerland
e-mail: illia.horenko@usi.ch

D. Igdalov
e-mail: dimitri.igdalov@usi.ch

A. D. Beck · C.-D. Munz
Institute of Aerodynamics and Gas Dynamics, University of Stuttgart, Stuttgart, Germany
e-mail: beck@iag.uni-stuttgart.de

C.-D. Munz
e-mail: munz@iag.uni-stuttgart.de

G. J. Gassner
Mathematical Institute, University of Cologne, Cologne, Germany
e-mail: ggassner@math.uni-koeln.de

surrogate models for the non-resolved fluctuations. Vector-valued auto-regressive models with external influences (VARX-models) form the basis for the modeling approach. The reconstruction capabilities of the modeling ansatz are tested against fully three dimensional turbulent channel flow data computed by direct numerical simulation (DNS). We present here the outcome of our reconstruction tests.

# 1 Introduction

The LES Navier-Stokes equations in their mathematical formulation incorporate the so-called sub grid scale stress tensor, $(\tau_{ij})$, which links the resolved eddies on the large scales (larger than a specific filter width) and the unresolved eddies on the small scales (smaller than that filter width), see, e.g. [9]. The sub grid scale stress tensor, written as

$$\tau_{ij} = \widetilde{u_i u_j} - \tilde{u}_i \tilde{u}_j, \tag{1}$$

with ~ as the filtered quantities, implements the unfiltered velocity field, $u$, which is not known a priori. It, therefore, has to be prescribed by an appropriate model function. Despite the progress that has been made, determining a suitable sub grid scale model remains a challenging task.

In this paper, we present our stochastic modelling ansatz based on the integral conservation laws developed in preparation of a novel LES closure approach. The reconstruction capabilities of the data-based modeling approach are tested against three dimensional direct numerical simulation (DNS) turbulent channel flow data computed for an incompressible, isothermal fluid at Reynolds number $Re_\tau = 590$. Our approach is similar in spirit to earlier propositions, e.g., [10], but differs in terms of both the stochastic modelling ansatz, and in terms of the underlying combined Finite Volume—Discontinuous Galerkin approximation framework, e.g., [1]. We, here, mention that our approach particularly allows for the analysis of non-stationary and non-homogeneous data, resp. That is in contrast to stationary (homogeneous) patterns, e.g. first order and second order statistics, that could lead to biased results due to their inability to characterize inhomogeneous (instationary) data sufficiently.

The integral conservation laws for mass, momentum and energy of a flow field are universally valid for arbitrary control volumes. These laws describe the time evolution of the integral values of the conserved quantities per control volume as a function of the associated fluxes across its bounding surfaces written as

$$\int_\Omega u_t \, dx + \oint_{\partial\Omega} f_n(u) \, ds = 0, \tag{2}$$

whith $\Omega$ and $\partial\Omega$ as the control volume and its surface, respectively, $u$ as the conserved quantity and $f_n$ as its normal flux across $\partial\Omega$. The *exact* evolution for the mean value of $\bar{u}$ is given by

$$\bar{u}_t = -\frac{1}{|\Omega|} \oint_{\partial \Omega} f_n(u)\, ds. \tag{3}$$

Thus, if the associated fluxes $f_n(u)$ across its bounding surfaces are determined exactly, the equations capture the underlying physics of conservation correctly and guarantee an accurate prediction of the temporal evolution of the integral mean values.

In the discrete view, the discretization basis for the Finite Volume approach are generally written as

$$\bar{u}_t + \frac{1}{|\Omega|} \oint_{\partial \Omega} H(\bar{u}^+, \bar{u}^-) ds = 0, \tag{4}$$

with a numerical flux function $H(\bar{u}^+, \bar{u}^-)$. Note that the arguments $\bar{u}^+$ and $\bar{u}^-$ result from a suitable spatial reconstruction of the mean values. Thus, a discretization error is introduced into the evolution of $\bar{u}$ in case $H(\bar{u}^+, \bar{u}^-) \neq f_n(u)$.

With respect to the LES method, an induced flux separation can be written as

$$f_n(u) := H(\bar{u}^+, \bar{u}^-) + \Delta f, \tag{5}$$

with $H(\bar{u}^+, \bar{u}^-)$ as the coarse flux, $^+$ and $^-$ as the specific cell face side, and $\Delta f$ as the flux correction. Thus, the reconstruction of the flux correction terms, i.e. of the sub grid fluxes, are necessary to yield the exact evolution of the mean values. Starting from this concept, we model the temporal structure of the fluxes to create a discrete formulation.

## 2 Stochastic Modeling Approach

We use advanced methods of time series analysis for the data-based construction of stochastic models with inherently non-stationary statistical properties to construct stochastic surrogate models for the non-resolved fluxes and flux corrections from specific time series (cf. [7]). Vector-valued auto-regressive models with external influences (VARX-models) form the basis for the modeling approach. We realize non-stationary statistical properties of these models by allowing for time dependent switches between different fluctuation regimes which are represented by different, but fixed, sets of the stochastic model parameters. The LES-grid-averaged conserved quantities on the coarse grid cells in the immediate vicinity of a given LES grid cell interface are interpreted as external influences in constructing the VARX surrogate model. In this fashion the stochastic models incorporate the information available from a typical numerical discretization stencil as would be used, e.g., in formulating a classical Smagorinsky closure.

The ansatz of the VARX model reads

$$\Delta f_{t,\mathbf{x}} = \mu(t, \mathbf{x}) + A(t, \mathbf{x})\phi_1(\Delta f_{t-\tau}, \ldots, \Delta f_{t-m\tau})_{\mathbf{x}} + B(t, \mathbf{x})\phi_2(u_{t,\mathbf{x}}) + \varepsilon_{t,\mathbf{x}}, \tag{6}$$

where $\Delta f_{t,\mathbf{x}}$ is the flux correction term ($t$ as time, $\mathbf{x}$ as 3D space vector), $(\mu, A, B)$ $(t, \mathbf{x})$ are time-dependent model parameters ($m$ as the memory depth), $\phi_1, \phi_2$ are model ansatz functions which are generally non-linear, $u_{t,\mathbf{x}}$ denotes the external influences (here the coarse-grid stencil data), and $\varepsilon_{t,\mathbf{x}}$ is the model-data discrepancy.

The basic idea of the approach is to detect the switching processes between the fluctuation regimes and their parameters, here named $(\gamma_j(t, \mathbf{x}), \Theta_j)$ ($j$ as the cell index), which characterize the local models. $\Theta_j$ denotes $K$ sets of $k$ parameters $\{\Theta_j \equiv (\Theta_1, \ldots, \Theta_k)_j\}_{j=1}^K$ representing the model parameters $(\mu, A, B)$, and $\gamma_j$ are model affiliation functions with $\gamma_j(t, \mathbf{x}) \in [0, 1]$ and $\sum_{j=1}^K \gamma_j(t, \mathbf{x}) = 1$. The total variation $TV$ of $\gamma_j$ is bounded

$$TV_t(\gamma_j(t, \mathbf{x})) \leq C.$$

With $(k, K, C)$ given, the best-fit and therefore the optimal parameters are found with minimization of the model-data distance, i.e.

$$\int_t \int_{\mathbf{x}} \delta_{t,\mathbf{x}} \, d\mathbf{x} \, dt \to \min_{\gamma, \Theta} \qquad \text{where } \delta_{t,\mathbf{x}} = \sum_{j=1}^K \gamma_j(t, \mathbf{x}) \left\| \varepsilon_{t,\mathbf{x}}^{\Theta_j} \right\|.$$

A balance between the requirements of high representation quality and low number of free parameters (Occam's razor) is achieved by involving criteria from information theory.

An extensive study based on the information criteria on a high performance computer cluster using about 280,000 CPU-hours at CSCS, Switzerland, shows that the general model ansatz, (6), simplifies in the context of our approach and that the optimal model is the VX approach instead of the VARX-ansatz, i.e. our modelling approach should not include auto-regressive terms. It follows, also, that it is sufficient to fit the cells' time series of exact LES corrections $\Delta F_{exact}^j(t)$ by means of an affine linear function that depends only of the cells' available LES observables $u^j(t)$. With that, our model approach now reads

$$\Delta f_{t,\mathbf{x}}^{\gamma,\Theta} = \sum_{j=1}^K \gamma_j(\mathbf{x}) \left[ \mu_j + \mathbf{B}_j \phi_2(\mathbf{u}_{t,\mathbf{x}}) \right]. \tag{7}$$

Thus, the resulting best-fit model takes the form of a stencil-based LES-closure that determines the turbulent flux corrections just from the cell-averages in a finite number of grid cells surrounding the considered grid cell interface.

The term $\mathbf{B}\phi_2(\mathbf{u})$ comprises the already mentioned flux correction terms derived from classical finite volume numerical flux functions of different order

$$\mathbf{B}\phi_2(\mathbf{u}) = \left( b_1 \Delta F^{1\text{st}} + b_2 \Delta F^{2\text{nd}} + b_3 \Delta F^{\text{WENO}} + b_{\text{lin}} \Delta F^{\text{lin}} \right)(\mathbf{u}), \tag{8}$$

with a next-neighbour stencil for $\Delta F^{\text{lin}}(\mathbf{u})$, and the term $(\mathbf{u}_{t,\mathbf{x}})$ in (7) incorporates coarse-grid stencil data, cf. Sect. 3.

## 3 Data Preprocessing

The computations of the turbulent channel flow data, described in detail in [12], made use of the instationary Navier-Stokes equations which were solved for the wall normal vorticity $\eta = \partial_z u - \partial_x w$, with $(u, v, w)$ as the velocity components in $(x, y, z)$, and the Laplacian of the wall-normal velocity $\varphi = \nabla^2(v)$. The boundary condition at the side walls in $y$-direction, normal to the main flow in $x$-direction, were rigid wall no-slip conditions, and along the other $(x\text{-}, z)$-axes periodic boundaries were defined. A pseudo-spectral Fourier-Chebyshev method similar to [3, 4] and [8] with Chebyshev-tau formulation in wall-normal $y$-direction and Fourier representation in the other directions was used. A third order Runge-Kutta based approach was used for time discretization of the non-linear convective terms and an implicit Euler approach was used for the viscous terms, cf. [13].

The data set that we use here consists of snapshots at 240 particular time steps in terms of the wall-unit time $\Delta t^+ = 1$, where $t^+ = \frac{t \cdot u_\tau^2}{\nu}$, with $u_\tau$ as the shear velocity and $\nu$ as the kinematic viscosity. The spatial resolution is $600 \times 385 \times 600$ in $(x, y, z)$.

To focus on the aim of our study, the 3D velocity field has to be recomputed and resampled resulting in particular spatial and temporal resolutions corresponding to typical LES simulations. The DNS data are averaged on a coarse LES grid which is a cartesian finite volume grid with equidistant spacing in all coordinates. Resolved LES-grid fluxes are determined from these averages using a straight-forward finite volume approximation for the Euler equations.

By subtracting these resolved fluxes from the DNS fluxes averaged over the cell faces of the LES-grid, we obtain one time series of non-resolved fluxes for each cell interface of the LES grid. In particular, we compute a so-called exact flux $(F_{ex})$ based on the preprocessed DNS data, and, furthermore, a reference flux $(F_{ref})$ and numerical fluxes of particular order from the average velocity data on the coarse grid. Once those flux data have been generated, flux correction terms are calculated as follows.

For each cell $C^j$ and each face $a = 1, \ldots, 6$ on the coarse-grid, a time series of the following LES-observables is calculated: (a) the exact flux corrections $\Delta F_{ex}^{j,a}(t) \in R^3$, (b) the 1st order flux corrections $\Delta F_1^{j,a}(t) \in R^3$, (c) the 2nd order flux corrections $\Delta F_2^{j,a}(t) \in R^3$, and (d) the 3rd order flux corrections $\Delta F_3^{j,a}(t) \in R^3$, where

$$\Delta F_i = F_i - F_{ref}, \ i = 1, 2, 3, \qquad \Delta F_{ex} = F_{ex} - F_{ref},$$

**Fig. 1** Exp A: Snapshot of the DNS velocity data and the resampled LES velocity data of the turbulent channel flow. Two *left panels*: main flow component (x-direction), two *right panels*: secondary flow component (z-direction)
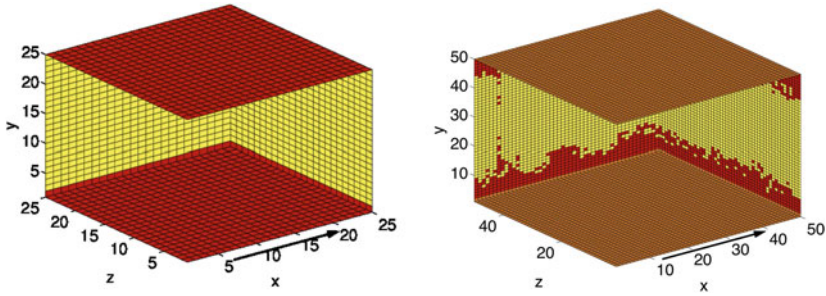
and, finally, the velocity field $V^j(t) \in R^{21}$ consisting of the average velocity field of cell $C^j$ and of the average velocity fields of all cells sharing a common face with $C^j$ (neighbored cell).

Generally, the numerical flux function proposed by [2] is used to compute fluxes of the particular order. For $F_1$, the cell average state value is assumed to cover the whole grid cell and, thus, the values at the cell faces are assumed to be equal to the cell center value. Consequently, no state reconstruction is needed. However, for $F_2$, piecewise linear state reconstruction within the grid cells is performed direction by direction based on the cell center values as in standard second order $FV$ methods using a monotonized central limiter, [6], for slope limiting during the reconstruction. This yields higher-order accurate cell interface data. Finally, for $F_3$, state recovery at the cell faces is obtained via a third order WENO scheme proposed in [11]. For $F_{ref}$, no state reconstruction or specific numerical flux function is used but the simple flux average is calculated.
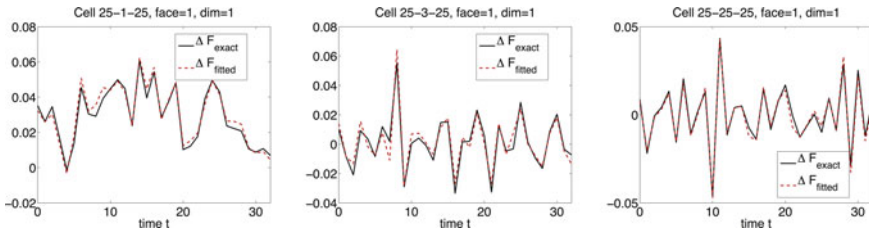
The flux correction data sets are then applied to the above mentioned stochastic model framework. In our approach, the pre-processing is a crucial procedure which here can be described in compressed form only, but is described more extensively in a forthcoming paper, cf. [5].

## 4 Results

In our study, we focus on two particular LES grid resolutions. Therefore, the DNS data are coarsened to a grid resolution of $25 \times 25 \times 25$ cells and to $50 \times 50 \times 50$ cells, here after referred to as *Exp A* and *Exp B*. For *Exp A*, Fig. 1 shows a snapshot of the recomputed (original) DNS velocity data and the corresponding LES velocity data for the principal flow direction $(x-)$ and for the secondary $(z-)$ flow component. Obviously, the main flow structures are captured in the LES data.

**Fig. 2** Fluctuation regimes for *Exp* A (*left*) and *Exp* B (*right*). The arrow indicates the main flow direction



**Fig. 3** Exp B: reconstruction of the exact flux correction terms, $\Delta F_{exact}$, determined with the optimal model parameters. From *left* to *right*: a LES cell located in the boundary layer, a near-boundary layer cell, a cell located in the flow interior. Panels show the $x$-component of the flux corrections for cell face 1 that is normal to the main flow ($x-$) component

For both LES grid resolutions, the analysis of the data with the stochastic modeling approach results in stationary, i.e. time-independent, models ($C = 0$), but the number of the fluctuation regimes is different as two and three clusters are determined for *Exp A* and *Exp B*, resp. For *Exp A*, the two clusters represent the boundary layer and the flow interior (Fig. 2, left). For *Exp B*, the boundary layer and the flow interior are also represented by two clusters, similar to *Exp A*, and the third regime is associated to cells which are located close to the rigid boundary wall (Fig. 2, right). We, therefore, call the third regime *transition model*.

The optimal model parameters are used to reconstruct the flux correction terms, cf. (7), shown exemplarily for *Exp B* in Fig. 3. The fitted flux correction terms show good agreement for cells located in the flow interior as well as for those cells located in the boundary layer.

# 5 Concluding Remarks

In this paper, we have presented the outcome of our reconstruction test, and we show specifically results of the non-trivial time series data analysis. We found resolution-dependent closure regimes as a third fluctuation regime has been detected when the

data are coarsened to $50 \times 50 \times 50$ grid cells but only two fluctuation regimes have been detected in *Exp A*.

Due to lack of space we, here, were able to present a compressed summary of our results. We, also, have tested our approach against data of a Taylor-Green vortex flow showing a transition from laminar to fully turbulent flow. In that test case, our approach also captures non-stationary regimes. We refer the reader to a more extensive description of our work, given in a forthcoming paper which has been currently submitted to *Meteorologische Zeitschrift*, [5].

These results encourages us for the ambitious attempt at dynamic LES closure along these lines.

# References

1. Gassner, G.J., Beck, A.D.: On the accuracy of high-order discretizations for underresolved turbulence simulations. Theoret. Comput. Fluid Dyn. **27**(3–4), 221–237 (2013). doi:10.1007/s00162-011-0253-7, http://dx.doi.org/10.1007/s00162-011-0253-7
2. Hickel, S., Adams, N., Domaradzki, J.: And adaptive local deconvolution method for implicit les. J. Comput. Phys. **213**, 413–436 (2006)
3. Jimenez, J., Moin, P.: The minimal flow unit in near-wall turbulence. J. Fluid Mech. **225**, 213–240 (1991)
4. Kim, J., Moin, P., Moser, R.: Turbulence statistics in fully developed channel flow at low reynolds number. J. Fluid Mech. **177**, 133–166 (1987)
5. von Larcher, T., Beck, A., Klein, R., Horenko, I., Metzner, P., Waidmann, M., Igdalov, D., Gassner, G., Munz, C.D.: A framework for the stochastic modelling of subgrid scale fluxes for large eddy simulation. Meteorol. Z. (Submitted)
6. van Leer, B.: Towards the ultimate conservative difference scheme. ii. monotonicity and conservation combined in a second-order scheme. J. Comput. Phys. **14**, 361–370 (1974)
7. Metzner, P., Putzig, L., Horenko, I.: Analysis of persistent non-stationary time series and applications. CAMCoS **7**, 175–229 (2012)
8. Moser, R., Kim, J., Mansour, N.: Direct numerical simulation of turbulent channel flow up to $Re_\tau = 590$. Phys. Fluids **11**(4), 943–945 (1999)
9. Sagaut, P.: Large Eddy Simulation for Incompressible Flows. Springer, Berlin (2006)
10. Scotti, A., Meneveau, C.: A fractal model for large eddy simulation of turbulent flows. Physica D **127**, 198–232 (1999)
11. Shu, C.W.: Essentially non-oscillatory and weighted essentially non-oscillatory schemes for hyperbolic conservation laws. Technical Report 97–65, ICASE, NASA Langley Research Center, Hampton, VA (1997)
12. Uhlmann, M.: Generation of initial fields for channel flow investigation. Intermediate Report (2000). http://www-turbul.ifh.uni-karlsruhe.de/uhlmann/home/report.html. Accessed 15 Jan 2014
13. Uhlmann, M.: The need for de-aliasing in a Chebyshev pseudo-spectral method. Technical Note No. 60 (2000). http://www-turbul.ifh.uni-karlsruhe.de/uhlmann/reports/dealias.pdf. Accessed 15th Jan 2014

# A Well Balanced Scheme for a Transport Equation with Varying Velocity Arising in Relativistic Transfer Equation

**T. Leroy, C. Buet and B. Després**

**Abstract** We are interested in the study of numerical schemes for the homogeneous in space asymptotic limit in the non equilibrium regime of the relativistic transfer equation. This limit leads to a frequency drift term modeling the Doppler effects for photons, and our aim is to design costless well-balanced schemes. One difficulty is that wave speed may vanish, which implies that standard well-balanced schemes constructed by discretizing the source term at the interfaces and by using a Godunov scheme may become inconsistent in this limit. This is indeed observed numerically.

## 1 Introduction

Our model equation comes from photons transport models with Doppler effects [1, 2, 6]. These Doppler effects are modeled by a frequency drift term:

$$\begin{cases} \partial_t \rho = \frac{\kappa}{3} \nu \partial_\nu \rho + \sigma(\nu)(B(\nu) - \rho) & \text{in } \mathbb{R}_t^+ \times \mathbb{R}_\nu^+, \\ \rho(0, \nu) = \rho^{in}(\nu), \end{cases} \tag{1}$$

with no need of boundary condition at $\nu = 0$ since the equation is degenerate. Here $\rho = \rho(t, \nu)$ represents the density of photons and $\kappa$ is the divergence of the velocity

---

T. Leroy (✉) · B. Després
LJLL/UPMC, 75252 Paris Cedex 05, France
e-mail: thomas.leroy@upmc.fr

B. Després
e-mail: despres@ann.jussieu.fr

T. Leroy · C. Buet
CEA, DAM, DIF, 91297 Arpajon Cedex, France

C. Buet
e-mail: christophe.buet@cea.fr

of the matter. For the simplicity of the presentation, we consider in this work that $0 < \kappa \le \kappa^*$. The coefficient $\sigma$ is the emission absorption coefficient. It is known to be very irregular with respect to the frequency. The function B is the Planck's function: $B(\nu) = \nu^3(e^\nu - 1)^{-1}$. Numerical methods for the coupling of hydrodynamics and transfer equations have been extensively studied (see for example [5, 7]) and are still an active field of research. The model problem (1) is also representative of asymptotic preserving issues, due to the parameter $\kappa$, the limit system being

$$\begin{cases} \partial_t \rho = \sigma(\nu)(B(\nu) - \rho) & \text{in } \mathbb{R}_t^+ \times \mathbb{R}_\nu^+, \\ \rho(0, \nu) = \rho^{in}(\nu). \end{cases} \tag{2}$$

As we explain in the next section, system (1) has stationary solutions, and our aim is to design numerical schemes which preserve these solutions. This could be interesting for kinetic equations for which the frequency discretization is known to be very costly. A first approach consists to use Greenberg-Leroux [4] type schemes. These schemes are well-balanced, but analytical and numerical studies show that they are not consistent in the limit regime $\kappa \to 0$. We propose a new scheme, for which we prove a uniform (according to the parameter $\kappa$) convergence result. We present numerical results which confirm this study.

## 2 Well-Balanced Schemes

For technical reasons, we restrict the frequency domain to $\mathbb{D} = [0, \nu^*]$, for a given $0 < \nu^* < +\infty$. Since we want to design well-balanced schemes, we are interested in the stationary solutions of (1). We thus solve the following Cauchy problem

$$\begin{cases} \frac{\kappa}{3}\nu\partial_\nu\rho + \sigma(\nu)(B(\nu) - \rho) = 0, \\ \rho(\nu^*) = \rho^*. \end{cases} \tag{3}$$

The positivity of the parameter $\kappa$ yields a transport of the photons toward the frequency $\nu = 0$, and is the reason of the boundary data at $\nu^*$. This is a simple O.D.E., and one can find the analytical solution, given by

$$\rho(\nu; \rho^*, \nu^*) = \rho^* e^{-\frac{3}{\kappa}\int_\nu^{\nu^*} \frac{\sigma(s)}{s}ds} + \frac{3}{\kappa}\int_\nu^{\nu^*} \frac{\sigma(s)B(s)}{s} e^{-\frac{3}{\kappa}\int_\nu^s \frac{\sigma(\tau)}{\tau}d\tau}ds. \tag{4}$$

Noting that $3\sigma(s)(s\kappa)^{-1} e^{-\frac{3}{\kappa}\int_\nu^s \frac{\sigma(\tau)}{\tau}d\tau} = -\frac{d}{ds}(e^{-\frac{3}{\kappa}\int_\nu^s \frac{\sigma(\tau)}{\tau}d\tau})$, one finds

$$\lim_{\kappa \to 0} \rho(\nu; \rho^*, \nu^*) = B(\nu). \tag{5}$$

For $1 \leq j \leq N$, we consider an irregular mesh defined by $(N + 1)$ points $0 = v_{\frac{1}{2}} < ... < v_{N+\frac{1}{2}} = v^*$. We define $v_j$ as the middle of the j-th frequency band, i.e. $v_j = (v_{j-\frac{1}{2}} + v_{j+\frac{1}{2}})/2$ and we denote $\Delta v_j$ its length. We also define the dual $(j + \frac{1}{2})$-th frequency band as the cell $[v_j, v_{j+1}]$, which length is denoted $\Delta v_{j+\frac{1}{2}}$. We denote $h = \max_j \Delta v_j$. We assume that there exists a constant C such that $\forall j \in \{1, ..N\}, \ 0 < Ch \leq \Delta v_j$.

## 2.1 A First Class of Well-Balanced Schemes

As presented in the introduction, we study a class of well-balanced schemes in the spirit of what was introduced by Greenberg-Leroux [4] (see [3] for a recent state of the art on the topic). It consists to localized the source term at the interfaces and to use a Godunov method to construct a scheme for the resulting equation. For equation (1), it yields the following scheme, denoted as WB1 in the numerical results:

$$\frac{d}{dt}\rho_j = \frac{\kappa}{3}v_j \frac{\rho(v_j; \rho_{j+1}; v_{j+1}) - \rho_j}{\Delta v_j}, \quad 1 \leq j \leq N. \tag{6}$$

This scheme is well-balanced by construction. Actually, for a stationary solution $\rho_j = \rho(v_j; \rho^*, v^*)$, the semigroup property yields

$$\rho(v_j; \rho_{j+1}; v_{j+1}) = \rho(v_j; \rho(v_{j+1}; \rho^*; v^*); v_{j+1}) = \rho(v_j; \rho^*; v^*),$$

and thus $\frac{d}{dt}\rho_j = 0$. On the other hand, this scheme is not consistent in the regime $\kappa \to 0$. Actually taking into account the limit as $\kappa$ tends to 0 of the analytical stationary solution (5), one finds for this scheme $\lim_{\kappa \to 0} \frac{d}{dt}\rho_j = 0$, which obviously is not a consistent discretization of the limit equation (2). We propose a new construction strategy to avoid this consistency problem.

## 2.2 Spectrally Well-Balanced Scheme

We study and prove several properties for the following scheme, denoted as the Spectrally Well-Balanced (SWB) scheme:

$$\frac{d}{dt}\rho_j = \frac{\sigma(v_j)}{1 - M(v_{j+1}; v_j)}\left(\rho(v_j; \rho_{j+1}, v_{j+1}) - \rho_j\right), \quad 1 \leq j \leq N, \tag{7}$$

with the natural boundary condition $\rho_{N+1} = \rho(t, v_{N+1})$, where $\rho(t, v_{N+1})$ is defined by the formula (15). The scheme is built using the integrating factor. It is defined,

for an arbitrary $v_0 \in \mathbb{D}$, by $M(v; v_0) = e^{-\frac{3}{\kappa} \int_{v_0}^{v} \frac{\sigma(s)}{s} ds}$. Multiplying equation (1) by $3\sigma(v)M(v; v_0)/\kappa v = -M'(v; v_0)$ yields

$$-M'(v; v_0)\partial_t \rho = \sigma(v)\left(\partial_v \big(M(v; v_0)\rho\big) + \frac{3\sigma(v)}{\kappa v} M(v; v_0)B(v)\right).$$

Integrating this equation between $v_j$ and $v_{j+1}$ and discretizing each term conveniently, one finds

$$-\Big[M(v; v_0)\Big]_{v_j}^{v_{j+1}} \frac{d}{dt}\rho_j = \sigma(v_j)\left(\Big[M(v; v_0)\rho\Big]_{v_j}^{v_{j+1}} + \int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_0)B(s)ds\right).$$

Taking $\rho(v_{j+1}) = \rho_{j+1}$, $\rho(v_j) = \rho_j$ and dividing this equation by $M(v_j; v_0)$ and $1 - M(v_{j+1}; v_j)$, one finds, using the relation $M(v; v_0)/M(s; v_0) = M(v; s)$,

$$\frac{d}{dt}\rho_j = \frac{\sigma(v_j)}{1 - M(v_{j+1}; v_j)}\left(M(v_{j+1}; v_j)\rho_{j+1} - \rho_j + \int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j)B(s)ds\right).$$

The definition of the stationary solution $\rho(v_j; \rho_{j+1}, v_{j+1}) = M(v_{j+1}; v_j)\rho_{j+1} + \int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j)B(s)ds$ yields the SWB scheme (7). The same argument than for the scheme (6) shows that this scheme is well-balanced. We prove a uniform (in $\kappa$) convergence result for this scheme. We define $\rho_h = (\rho_j)_{1 \leq j \leq N}$ and the discrete norm $\| \cdot \|_{L_d^2}$ such that $\|\rho_h\|_{L_d^2}^2 = \sum_j \Delta v_{j+\frac{1}{2}} \rho_j^2$. We make some assumptions:

- The initial data satisfies $\rho^{in} \in H^2(\mathbb{D})$.
- The emission absorption coefficient satisfies $\sigma_a \in W^{2,\infty}(\mathbb{D})$. Moreover, there exists a constant $\sigma_* > 0$ such that $\forall v \in \mathbb{D}, \sigma(v) \geq \sigma_*$.

We need the following stability result

**Lemma 1** (*$L^2$ Stability*) *Under these assumptions, the following estimate holds, where the constants $C(T)$ depends on all the parameters and the boundary condition but is uniform in $\kappa \in (0, \kappa^*]$:* $\|\rho_h(t)\|_{L_d^2} \leq C(T)\sqrt{1 + \|\rho_h(0)\|_{L_d^2}^2}$, $0 < t < T$.

*Proof* Since the proof is rather classical, we only develop the main ideas. We want to reveal in the SWB scheme (7) a consistent discretization of Eq. (1). Injecting the expression of the stationary solution (4), one can write it as

$$\frac{d}{dt}\rho_j = \frac{\kappa}{3}v_j \frac{\rho_{j+1} - \rho_j}{\Delta v_{j+\frac{1}{2}}} + \sigma(v_j)\big(B(v_j) - \rho_j\big) + \sigma(v_j)\big(\rho_{j+1} - \rho_j\big)R_{j,1} + R_{j,2}, \quad (8)$$

with

$$\begin{cases} R_{j,1} = \frac{1}{1 - M(v_{j+1}; v_j)} - \frac{\kappa v_j}{3\sigma(v_j)\Delta v_{j+\frac{1}{2}}} - 1, \\ R_{j,2} = \frac{\sigma(v_j)}{1 - M(v_{j+1}; v_j)} \int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j)\big(B(s) - B(v_j)\big)ds, \end{cases} \quad (9)$$

where we used, by definition of $M(s; v_j)$, $\int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j) ds = 1 - M(v_{j+1}; v_j)$. We introduce $B_h = (B(v_j))_{1 \le j \le N}$. Using the positivity of the coefficients $\kappa$ and $\sigma$ and the Cauchy-Schwarz inequality, one finds by multiplying Eq. (8) by $\Delta v_{j+\frac{1}{2}} \rho_j$ and adding on all the cells a positive constant C such that,

$$\frac{d}{dt} \|\rho_h\|_{L_d^2}^2 \le C\left(1 + \max_j |R_{j,1}|\right) \|\rho_h\|_{L_d^2}^2 + \|\sigma\|_{L^\infty} \|B\|_{L_d^2}^2 + \|R_{j,2}\|_{L_d^2}^2 + \frac{\kappa}{6} \rho_{N+1}^2. \quad (10)$$

We thus need to control $R_{j,1}$ and the $L_d^2$ norm of $R_{j,2}$. Denoting $z_j = \frac{3}{\kappa} \int_{v_j}^{v_{j+1}} \frac{\sigma(s)}{s} ds$ and using the definition of $M(v_{j+1}; v_j)$, one can write $R_{j,1}$ as

$$R_{j,1} = \left(\frac{1}{1 - e^{-z_j}} - \frac{1}{z_j}\right) - 1 + \left(\frac{1}{z_j} - \frac{\kappa v_j}{3\sigma(v_j)\Delta v_{j+\frac{1}{2}}}\right).$$

For the first term one has $\frac{1}{1-e^{-z_j}} - \frac{1}{z_j} \le 1$. Using a Taylor expansion of the function $v \mapsto \sigma(v)v^{-1}$ at the frequency $v_j$, one finds

$$\max_j |R_{j,1}| \le C \frac{\kappa}{3} \frac{\|\sigma\|_{W^{1,\infty}}}{\sigma_*^2}, \quad (11)$$

where the constant C depends on the mesh but is independent of $\kappa$. We now turn to the term $R_{j,2}$. A Taylor expansion of the function $v \mapsto B(v)$ shows

$$|R_{j,2}| \le \Delta v_{j+\frac{1}{2}} \|B'\|_{L^\infty} \frac{\sigma(v_j)}{1 - M(v_{j+1}; v_j)} \int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j) ds.$$

Using the relation $\int_{v_j}^{v_{j+1}} \frac{3\sigma(s)}{\kappa s} M(s; v_j) ds = 1 - M(v_{j+1}; v_j)$, one finds

$$\|R_{j,2}\|_{L_d^2} \le h\sqrt{v^*} \|B'\|_{L^\infty} \|\sigma\|_{L^\infty}. \quad (12)$$

Using all these results in (10), one finds a constant C such that

$$\frac{d}{dt} \|\rho_h(t)\|_{L_d^2}^2 \le C\left(1 + \frac{\kappa}{3} \frac{\|\sigma\|_{W^{1,\infty}}}{\sigma_*^2}\right) \|\rho_h(t)\|_{L_d^2}^2 + \|\sigma\|_{L^\infty} \|B\|_{L_d^2}^2 + h^2 v^* \|B'\|_{L^\infty}^2 \|\sigma\|_{L^\infty}^2$$
$$+ \frac{\kappa}{6} \rho_{N+1}^2.$$

The Gronwall lemma finally gives the result.

The key point was to prove a uniform estimate for the consistency errors $R_{j,1}$ and $R_{j,2}$. Actually, estimates (11) and (12) are no longer true for the WB1 scheme (6). We now turn to the uniform (in $\kappa$) convergence result of the scheme (7). Without loss of generality, we assume that $\forall j, \rho_j(0) = \rho(v_j, 0)$. We define $\rho_{ex}(t) = (\rho(t, v_j))_{1 \le j \le N}$.

**Lemma 2** (Uniform convergence) *Under the same assumptions, the numerical solution of the scheme* (7) *satisfies the following estimate, where the constant $C(T)$ is uniform in $\kappa \in (0, \kappa^*]$: $\|\rho_{ex}(t) - \rho_h(t)\|_{L_d^2} \le C(T)h$, $0 < t < T$.*

*Proof* Evaluating the solution of the P.D.E. (1) at the frequency $v_j$ and using a Taylor expansion of the function $v \mapsto \rho(t, v)$ with integral remainder, one has

$$
\begin{aligned}
\frac{d}{dt}\rho(t, v_j) = \frac{\kappa}{3}v_j\left(\frac{\rho(t, v_{j+1}) - \rho(t, v_j)}{\Delta v_{j+\frac{1}{2}}} - \int_{v_j}^{v_{j+1}}\frac{v_{j+1} - s}{\Delta v_{j+\frac{1}{2}}}\partial_v^2\rho(t, s)ds\right) \\
+ \sigma(v_j)\left(B(v_j) - \rho(t, v_j)\right).
\end{aligned}
\tag{13}
$$

We obtain an equation on the unknown $e_j(t) := \rho_j(t) - \rho(t, v_j)$ by deducting to the expression (8) of the SWB scheme this equation. Multiplying the obtained equation by $\Delta v_{j+\frac{1}{2}}e_j(t)$ and adding on all the cells, one gets

$$
\begin{aligned}
\frac{1}{2}\frac{d}{dt}\|e_h(t)\|_{L_d^2}^2 = \sum_j \Delta v_{j+\frac{1}{2}}e_j(t)S_j(t) + \sum_j \Delta v_{j+\frac{1}{2}}e_j(t)\sigma(v_j)\left(\rho(v_{j+1}) - \rho(v_j)\right)R_{j,1} \\
+ \sum_j \Delta v_{j+\frac{1}{2}}e_j(t)R_{j,2} + \sum_j \Delta v_{j+\frac{1}{2}}e_j(t)\int_{v_j}^{v_{j+1}}\frac{v_{j+1} - s}{\Delta v_{j+\frac{1}{2}}}\partial_v^2\rho(t, s)ds,
\end{aligned}
\tag{14}
$$

where $S_j(t) = \frac{\kappa}{3}v_j\frac{e_{j+1}(t) - e_j(t)}{\Delta v_{j+\frac{1}{2}}} + \sigma(v_j)\left(e_{j+1}(t) - e_j(t)\right)R_{j,1} - \sigma(v_j)e_j(t)$ and $R_{j,1}$ and $R_{j,2}$ are defined in (9). We control successively each of these terms. First, the term $S_j(t)$ has already been studied. Using the estimate (11), one has

$$
\sum_j \Delta v_{j+\frac{1}{2}}e_j(t)S_j(t) \le C\frac{\kappa}{3}\frac{\|\sigma\|_{W^{1,\infty}}^2}{\sigma_*^2}\|e_h(t)\|_{L_d^2}^2,
$$

where the constant C depends on the mesh but is independent of $\kappa$. The term $R_{j,2}$ have also been controlled in the previous part. The estimate (12) and the inequality $ab \le (a^2 + b^2)/2$ gives

$$
\sum_j \Delta v_{j+\frac{1}{2}}e_j(t)R_{j,2} \le \frac{1}{2}\|e_h(t)\|_{L_d^2}^2 + \frac{1}{2}h^2v^*\|B'\|_{L^\infty}^2\|\sigma\|_{L^\infty}^2.
$$

Similar arguments associated to a Taylor expansion of the function $v \mapsto \rho(t, v)$ leads to

$$
\sum_j \Delta v_{j+\frac{1}{2}}e_j(t)\sigma(v_j)\left(\rho(v_{j+1}) - \rho(v_j)\right)R_{j,1} \le C\left(\|e_h(t)\|_{L_d^2}^2 + h^2\|\partial_v\rho(t)\|_{L^2(\mathbb{D})}^2\right),
$$

where, once again, the constant C is independent of $\kappa$. In the same way, one has for the last term

$$\sum_j \Delta v_{j+\frac{1}{2}} e_j(t) \int_{v_j}^{v_{j+1}} \frac{v_{j+1} - s}{\Delta v_{j+\frac{1}{2}}} \partial_v^2 \rho(t, s) ds \leq \frac{1}{2} \left( \|e_h(t)\|_{L_d^2}^2 + h^2 \|\partial_{vv} \rho(t)\|_{L^2(\mathbb{D})}^2 \right).$$

As $\rho$ is solution of a simple linear P.D.E., one easily controls its $H^2$ norm. Actually, using the regularity of $\sigma$, one finds a constant $C$ such as $\|\partial_v \rho(t)\|_{L^2}^2 \leq C(1 + \|\partial_v \rho(0)\|_{L^2}^2)$ and $\|\partial_{vv} \rho(t)\|_{L^2}^2 \leq C(1 + \|\partial_{vv} \rho(0)\|_{L^2}^2)$. Using all these results in (14), one finds another constant C, once again uniform in $\kappa$, such that

$$\frac{d}{dt} \|e_h(t)\|_{L_d^2}^2 \leq C \left( \|e_h(t)\|_{L_d^2}^2 + h^2 \right).$$

As before, the Gronwall lemma and the assumption on the initial data gives the announced result.

## 3 Numerical Results

In this part we present some numerical results, computed in the $L^1$ norm. As the model problem is a linear transport equation with damping and a source term, one finds the analytical solution:

$$
\begin{aligned}
\rho(t, v) = {} & \rho^{in}(v e^{\frac{\kappa}{3}t}) e^{-\frac{3}{\kappa} \int_v^{v e^{\frac{\kappa}{3}t}} \sigma(\tau)\tau^{-1}d\tau} \\
& + \int_0^t \sigma\left(v e^{\frac{\kappa}{3}s}\right) B\left(v e^{\frac{\kappa}{3}s}\right) e^{-\frac{3}{\kappa} \int_v^{v e^{\frac{\kappa}{3}s}} \sigma(\tau)\tau^{-1}d\tau} ds,
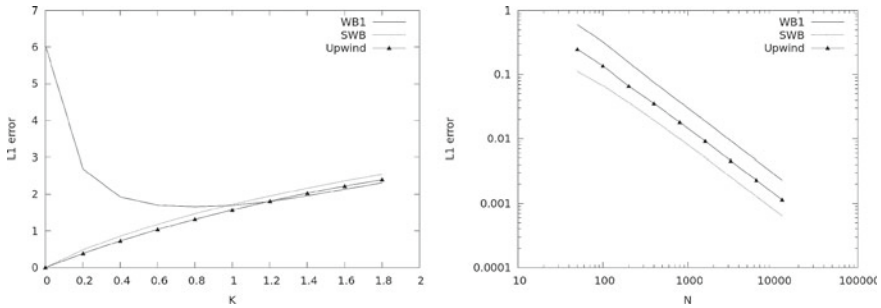\end{aligned}
\tag{15}
$$

which is used to compute error estimates. Numerically, the integrals are approximated by classical five points Gauss Legendre formulae. All the numerical tests are performed on a random mesh on a frequency domain $\mathbb{D} = [0, 30]$ with the following data: $\rho^{in}(v) = 0$ and $\sigma(v) = 1$. We use an explicit Euler discretization for the time derivatives. In all the following graphics we also plotted the upwind scheme, with the source term discretized in the middle of the cell :

$$\frac{d}{dt}\rho_j = \frac{\kappa}{3} \frac{\rho_{j+1} - \rho_j}{\Delta v_j} v_{j+\frac{1}{2}} + \sigma(v_j)\left(B(v_j) - \rho_j\right). \tag{16}$$

In Fig. 1 we displayed the relative $L^1$ error between the numerical solutions of the schemes and the analytical solution, with $N = 50$ cells and $\kappa = 1$. As expected, the SWB scheme and the WB1 scheme converge toward the analytical solution as time goes on.

**Fig. 1** Evolution of the $L^1$ relative error between the numerical and the analytical solutions, N = 50, $\kappa = 1$



**Fig. 2** *Left*: $L^1$ error versus K, N = 50, t = 2. *Right*: $L^1$ error versus N in a Log-Log scale plan, t = 2, $\kappa = 1$

As we have seen previously, the WB1 scheme is not consistent in the regime $\kappa \to 0$. Figure 2 plots on the left side the evolution, as $\kappa$ tends to 0 and at time $t = 2$, of the $L^1$ error between the solutions of the WB1, SWB and upwind schemes and the numerical solution of the following scheme, consistent with Eq. (2):

$$\frac{d}{dt}\rho_j = \sigma(v_j)\big(B(v_j) - \rho_j\big), \tag{17}$$

and confirms the theoretical study. On the right side we plotted the $L^1$ norm in a Log-Log scale between the analytical solution and the numerical solution of the schemes at time $t = 2$ and with $\kappa = 1$.

# References

1. Buet, C., Despres, B.: Asymptotic analysis of fluid models for the coupling of radiation and hydrodynamics. J. Quant. Spectrosc. Radiat. Transfer **85**, 385–418 (2004)
2. Godillon-Lafitte, P., Goudon, T.: A coupled model for radiative transfer: doppler effects, equilibrium, and nonequlibrium diffusion asymptotics. Multiscale Model. Simul. **4**(**4**), 1245–1279 (2005)
3. Gosse, L.: Computing qualitatively correct approximations of balance laws. In: Springer Series (2013)
4. Greenberg, J.M., Leroux, A.Y.: A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. J. Numer. Anal. **33**(1), 1–16 (1996)
5. Lowrie, R.B., Morel, J.E., Hittinger, J.A.: The coupling of radiation and hydrodynamics. Astrophys. J. **521**, 423–450 (1999)
6. Mihalas, D., Weibel-Mihalas, B.: Foundations of Radiation Hydrodynamics. ISBN 0-486-40925-2 (1999)
7. Turpault, R.: Properties and frequential hybridisation of the multigroup m1 model for radiative transfer. Non Linear Anal. Real World Appl. **11** (2012)

# An Arbitrary Space-Time High-Order Finite Volume Scheme for Gas Dynamics Equations in Curvilinear Coordinates on Polar Meshes

**Bertrand Meltz, Stéphane Jaouen and Frédéric Lagoutière**

**Abstract** We are interested in the study of numerical schemes for the resolution of gas dynamics equations which preserve symmetric (or axisymmetric) flows. A simple way to achieve this is to derive a numerical scheme whose mesh and coordinates system are aligned with the flow. Typically, for the simulation of cylindrical implosions of gas, the cylindrical coordinate system and a polar mesh are well-suited. But such coordinates systems introduce geometrical singularities as well as geometrical source terms. In this paper, we investigate an arbitrary high-order space-time Finite Volume (FV) scheme in cylindrical coordinates. Test-cases with and without polar symmetries are studied in order to confirm the order of the scheme as well as its robustness.

## 1 Introduction and Governing Equations

Symmetric (or axisymmetric) flows arise in many applications such as Inertial Confinement Fusion (ICF) experiments, or astrophysics. Usual FV methods built in Cartesian (or axisymmetric) coordinates and operating on slab meshes can capture these symmetries thanks to artificial viscosity models as in [6]. But such models are quite costly. We propose to derive a numerical solver in cylindrical coordinates using a FV method on a polar mesh. The solved equations are the Euler's equations which in cylindrical coordinates $(r, \varphi)$ write:

B. Meltz (✉) · S. Jaouen
CEA, DAM, DIF, 91297 Arpajon Cedex, France
e-mail: meltz.bertrand@gmail.com; bertrand-julien.meltz@cea.fr

S. Jaouen
e-mail: stephane.jaouen@cea.fr

F. Lagoutière
Laboratoire de Mathématiques d'Orsay, Université Paris-Sud, 91405 Orsay Cedex, France
e-mail: frederic.lagoutiere@math.u-psud.fr

$$\partial_t \begin{pmatrix} r\rho \\ r\rho u_r \\ r\rho u_\varphi \\ r\rho e \end{pmatrix} + \partial_r \begin{pmatrix} r\rho u_r \\ r\rho u_r^2 + rp \\ r\rho u_r u_\varphi \\ r(\rho e + p)u_r \end{pmatrix} + \partial_\varphi \begin{pmatrix} \rho u_\varphi \\ \rho u_r u_\varphi \\ \rho u_\varphi^2 + p \\ (\rho e + p)u_\varphi \end{pmatrix} = \begin{pmatrix} 0 \\ p + \rho u_\varphi^2 \\ -\rho u_r u_\varphi \\ 0 \end{pmatrix}. \quad (1)$$

The system is closed with an Equation of State (EOS) $p = p(\tau = \frac{1}{\rho}, \varepsilon)$, with $\varepsilon$ denoting the internal energy. The geometric source terms come from the divergence operator "$\nabla\cdot$" in polar coordinates $(r, \varphi)$ under conservative form. The terms $\rho u_\varphi^2$ and $\rho u_r u_\varphi$ are related to the centrifugal and Coriolis forces respectively. The system of balance laws (1) is solved using a FV method, based on a Lagrange-remap formalism together with a Directional Splitting Method (DSM) as in [2]. The DSM allows us to build efficient 1D schemes using centered discretizations on regular structured grid.

## 2 Numerical Scheme

Let $U$ be the vector of unknowns: $U = (\rho, \rho u_r, \rho u_\varphi, \rho e)$. Let $F^r(U)$ be the vector of fluxes along the radial direction, $F^\varphi(U)$ the vector of fluxes along the azimuthal direction, and $G(U)$ the vector of source terms. Using a DSM, the two systems to be alternatively solved are:

$$\partial_t(rU) + \partial_r(rF^r(U)) = G(U), \quad (2a)$$

$$\partial_t(rU) + \partial_\varphi(F^\varphi(U)) = 0. \quad (2b)$$

Each 1D scheme, called a sweep, is based on a Lagrange-remap solver. In the sequel, we will focus on the radial direction.

**Lagrangian step:** Let us introduce the Euler-Lagrange change of variable $(r, t) \rightarrow (R, t)$ defined by: $dr = J dR + u_r dt$ and $dt = dt$. We can rewrite the system (2a) with Lagrangian coordinates:

$$\partial_t(RU^0) + \partial_R(rF^{r,0}(U^0)) = G^0(U^0), \quad (3)$$

$$U^0 = \begin{pmatrix} \rho_0 \tau \\ \rho_0 u_r \\ \rho_0 u_\varphi \\ \rho_0 e \end{pmatrix}, \quad F^{r,0}(U^0) = \begin{pmatrix} -u_r \\ p \\ 0 \\ pu_r \end{pmatrix}, \quad G^0(U^0) = \begin{pmatrix} 0 \\ Jp + J\rho u_\varphi^2 \\ -J\rho u_r u_\varphi \\ 0 \end{pmatrix}, \quad J = \frac{R\rho_0}{r\rho}.$$

$\rho_0$ denotes the density at the beginning of each time step on the regular grid: $\rho_0(R, \varphi) = \rho(R, \varphi, t^n)$.

From now on, we drop the 0 superscript indexing the Lagrangian quantities. The system (3) is integrated over a space-time cell $\Omega_{i,j} \times [t^n; t^{n+1}]$ and divided by the control volume $|\Omega_{i,j}| \times \Delta t = R_j \Delta R \Delta \varphi \Delta t$:

$$\frac{\overline{U}^{n+1}_{i,j} - \overline{U}^n_{i,j}}{\Delta t} + \frac{F^r_{i,j+\frac{1}{2}} - F^r_{i,j-\frac{1}{2}}}{R_j \Delta R} = \frac{G_{i,j}}{R_j}, \tag{4}$$

with $\overline{U}^n_{i,j}$ denoting the 2D cell-average of $U$ over $\Omega_{i,j}$ at time $t^n$, $F^r_{i,j+\frac{1}{2}}$ is the 2D flux which can be defined from the 1D flux $F^r_{j+\frac{1}{2}}(\varphi)$ with a transverse integration:

$$F^r_{i,j+\frac{1}{2}} = \frac{1}{\Delta \varphi} \int_{\varphi_{i-\frac{1}{2}}}^{\varphi_{i+\frac{1}{2}}} F^r_{j+\frac{1}{2}}(\varphi) \mathrm{d}\varphi, \quad F^r_{j+\frac{1}{2}}(\varphi) = \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} (r F^r(U))(R_{j+\frac{1}{2}}, \varphi, t) \, \mathrm{d}t.$$

The same applies to $G_{i,j}$ which is the space-time average of $G(U)$ over $\Omega_{i,j} \times [t^n; t^{n+1}]$. The 1D fluxes are approximated using a local Taylor expansion of the solution around $(t^n)$. For each line $\varphi_i$, we have:

$$F^r_{j+\frac{1}{2}}(\varphi_i) = \sum_{k=0}^{N-1} \frac{\Delta t^k}{(k+1)!} \frac{\partial^k (r F(U))}{\partial t^k}(R_{j+\frac{1}{2}}, \varphi_i, t^n) + O(\Delta t^N), \tag{5}$$

Note that the Eulerian coordinate $r$ depends on time in the Lagrangian description and must also be expanded around $(t^n)$. The numerical source terms are also computed using a local expansion of the solution around $(R_j, t^n)$. The expansions introduce time-derivatives of thermodynamic quantities. These are replaced by spatial derivatives using the so-called Cauchy-Kovaleskaya procedure as in [5]. The 1D numerical fluxes operate on point-wise values of thermodynamic quantities at interfaces. We first compute an $N$-th order approximation of point-wise values of $U$ at cell-centers using a conservative polynomial interpolation on a centered stencil:

$$U^n_{ij} = U(R_j, \phi_i, t^n) + O(\Delta R^N) = p^l_0 \overline{U}^n_{ij} + \sum_{k=1}^{l} p^l_k \left( \frac{R_{j-k}}{R_j} \overline{U}^n_{i,j-k} + \frac{R_{j+k}}{R_j} \overline{U}^n_{i,j+k} \right), \tag{6}$$

with $l = \lfloor \frac{N}{2} \rfloor$ and the $p^l_k$ coefficients can be found in [2]. Any thermodynamic quantity $\psi$ and its space derivatives are then evaluated at the appropriate order of accuracy using centered finite differences operators with $s = \lceil \frac{N}{2} \rceil$:

$$\left( \partial^m_R \psi \right)^n_{i,j+\frac{1}{2}} = \left( \partial^m_R \psi \right)(R_{j+\frac{1}{2}}, \phi_i, t^n) + O(\Delta R^{N-m}),$$

$$= \begin{cases} \frac{1}{\Delta R^m} \sum_{k=1}^{s} d^s_{m,k} (\psi^n_{i,j+k} + \psi^n_{i,j-k+1}) & \text{if } m \text{ is even}, \\ \frac{1}{\Delta R^m} \sum_{k=1}^{s} d^s_{m,k} (\psi^n_{i,j+k} - \psi^n_{i,j-k+1}) & \text{if } m \text{ is odd}. \end{cases} \tag{7}$$

The $d^s_{m,k}$ coefficients can be found in [2]. Eulerian coordinates evolve according to $\partial_t(r)(R, \varphi, t) = u_r$. A time-integration over $[t^n; t^{n+1}]$ gives:

$$r^{n+1}_{i,j+\frac{1}{2}} - r_{j+\frac{1}{2}} = \int_{t^n}^{t^{n+1}} u_r(R_{j+\frac{1}{2}}, \varphi_i, t) \mathrm{d}t = (u_r)^*_{i,j+\frac{1}{2}} \Delta t. \tag{8}$$

Since $(u_r)^*$ is computed at desired order of accuracy, either are the positions of the interfaces $r^{n+1}_{i,j+\frac{1}{2}}$ at time $t^{n+1}$.

**Remap step:** Once the Lagrangian step is applied, conservative variables $(\rho\psi)$ for $\psi \in \{1, \boldsymbol{u}, e\}$ are remapped on the initial grid according to the following integral splitting:

$$
\begin{aligned}
|\Omega_{i,j}| \overline{(\rho\psi)}^{n+1}_{i,j} &= \int_{\varphi_{i-\frac{1}{2}}}^{\varphi_{i+\frac{1}{2}}} \int_{r_{j-\frac{1}{2}}}^{r_{j+\frac{1}{2}}} (r\rho\psi)(r, \varphi, t^{n+1}) \mathrm{d}r \, \mathrm{d}\varphi, \\
&= \int_{\varphi_{i-\frac{1}{2}}}^{\varphi_{i+\frac{1}{2}}} \left[ \int_{r_{j-\frac{1}{2}}}^{r^{n+1}_{i,j-\frac{1}{2}}} (r\rho\psi) + \int_{r^{n+1}_{i,j-\frac{1}{2}}}^{r^{n+1}_{i,j+\frac{1}{2}}} (r\rho\psi) + \int_{r^{n+1}_{i,j+\frac{1}{2}}}^{r_{j+\frac{1}{2}}} (r\rho\psi) \right] \mathrm{d}\varphi.
\end{aligned}
\tag{9}
$$

By definition of the Euler-Lagrange change of variable and introducing the exact remap fluxes, we have:

$$|\Omega_{i,j}| \overline{(\rho_0\psi)}^{n+1}_{i,j} = \int_{\varphi_{i-\frac{1}{2}}}^{\varphi_{i+\frac{1}{2}}} \int_{R_{j-\frac{1}{2}}}^{R_{j+\frac{1}{2}}} (R\rho_0\psi), \quad \text{and} \quad (\rho\psi)^*_{i,j+\frac{1}{2}} = \int_{\varphi_{i-\frac{1}{2}}}^{\varphi_{i+\frac{1}{2}}} \int_{r_{j+\frac{1}{2}}}^{r^{n+1}_{i,j+\frac{1}{2}}} (r\rho\psi).$$

The numerical remap fluxes are evaluated using a conservative polynomial interpolation together with an up-winding given by the sign of $(u)^*_{i,j+\frac{1}{2}}$ as in [6]. The remap step then writes:

$$|\Omega_{i,j}| \overline{(\rho\psi)}^{n+1}_{i,j} = |\Omega_{i,j}| \overline{(\rho_0\psi)}^{n+1}_{i,j} - \left( (\rho\psi)^*_{i,j+\frac{1}{2}} - (\rho\psi)^*_{i,j-\frac{1}{2}} \right). \tag{10}$$

**Multidimensional extension:** In order to preserve accuracy in the multidimensional case, one has to compute an high-order approximation of the transverse integrations. We use a conservative polynomial interpolation. For instance, computing an $N$-th order approximation of the 2D numerical Lagrangian fluxes writes:

$$\boldsymbol{F}^r_{i,j+\frac{1}{2}} = q^l_0 \boldsymbol{F}^r_{j+\frac{1}{2}}(\phi_i) + \sum_{k=1}^{l} q^l_k \left( \boldsymbol{F}^r_{j+\frac{1}{2}}(\phi_{i-k}) + \boldsymbol{F}^r_{j+\frac{1}{2}}(\phi_{i+k}) \right), \tag{11}$$

with $l = \lfloor \frac{N}{2} \rfloor$, and the $q^l_k$ coefficients can be found in [2].

One must also use high-order splitting sequences in order to preserve accuracy. Beyond the well-known Godunov (1st order) and Strang (2nd order) splitting sequences, high-order sequences contain necessarily negative time steps. Since

all computations are centered, no particular treatment is needed to handle such an unusual case. Such sequences can be found in [6] up to the 6th order.

The stability constraint is a CFL constraint which take the following form:

$$\max_{i,j} \left( \frac{\max(\Delta R, R_j \Delta \phi)(\|\boldsymbol{u}_{i,j}\| + c_{i,j})}{|\Omega_{i,j}|} \right) \Delta t \leq 1.$$

Most of the time, the velocity field is regular and the physical domain includes the pole axis. In this case, the stability constraint is computed in the vicinity of $r = 0$ and becomes parabolic:

$$\max_{i,j} \left( \frac{(\|\boldsymbol{u}_{i,j}\| + c_{i,j})}{R_j \Delta \varphi} \right) \Delta t \leq 1.$$

## 3 Hyperviscosity Model

It is well-known that high-order schemes are subject to Gibbs phenomenon: oscillations appear in the vicinity of discontinuities and can corrupt physical values. Slope limiters are often used in order to reduce the oscillations, but these techniques require many conditional tests and usually break the Experimental Order of Convergence (EOC). We rather propose to add artificial viscosity by mean of an hyperviscosity model as in [7]. One most important feature of this technique is that EOC is preserved. Practically, we take into account the viscous strain tensor $\underline{\tau}$ in the momentum equation as well as its work in the energy equation. If we consider a Newtonian gas (Stokes assumption), then the viscous strain tensor writes:

$$\underline{\tau} = \nu \left[ (\nabla \boldsymbol{u}) + (\nabla \boldsymbol{u})^T \right] + \left( \beta - \frac{2}{3} \nu \right) (\nabla \cdot \boldsymbol{u}) \underline{I},$$

$\nu$ and $\beta$ denote the dynamic viscosity and the volumic viscosity respectively. Hyperviscosity techniques consist in replacing $\nu$, $\beta$ by artificial constants $\nu^*$, $\beta^*$. These are computed in the same way as [1]. In order to remain consistent with the Euler's equations, these constants are designed to tend to 0 as the mesh is refined:

$$\beta^* = C_\beta \langle \rho | \nabla^k S | \rangle \Delta^{k+2}, \quad \nu^* = C_\nu \langle \rho | \nabla^k S | \rangle \Delta^{k+2},$$

with $S$ denoting $S = \sqrt{\underline{S} : \underline{S}}$, $\underline{S} = [(\nabla \boldsymbol{u}) + (\nabla \boldsymbol{u})^T]/2$. $C_\beta$ and $C_\nu$ are user-defined constants, and $\Delta$ is the characteristic size of a cell. In order to preserve accuracy, we set $k$ to 2 for the second-order scheme, and 4 for third and fourth order schemes, and so on. The operator $\langle \cdot \rangle$ denotes a truncated Gaussian filter discretized at cell-centers:

$$\langle \psi \rangle_{i,j} = \sum_{|k|,|l| \leq 4} f_{|k|} f_{|l|} \psi_{i+k,j+l},$$

with $f_0 = \frac{3565}{10368}$, $f_1 = \frac{3091}{12960}$, $f_2 = \frac{1997}{25920}$, $f_3 = \frac{149}{12960}$, $f_4 = \frac{107}{103680}$. Once the artificial constants have been computed at cell-centers, one can evaluate the components of the viscous strain tensor $\underline{\tau}$. Artificial viscosity is then added in the Lagrangian numerical scheme. Note that computing $\nabla \cdot \underline{\tau}$ produces viscous source terms. We choose to follow the same splitting strategy as the convective part.

## 4 Numerical Results

**Vortex** The first numerical study deals with a stationary isentropic vortex. Since an analytical solution is available in [8], we can assess the EOC of the scheme. Let $(x_0, y_0)$ be the coordinates of the eye of the vortex in the x-y plane. Three configurations are investigated:

1. centered configuration: the eye of the vortex is aligned with the pole of the mesh, physical space: $(r, \varphi) \in [0; 8] \times [0; 2\pi]$, $(x_0, y_0) = (0, 0)$,
2. shifted configuration 1: the eye of the vortex is located far away from the pole, such that we have an hydrostatic state in the vicinity of $r = 0$, physical space: $(r, \varphi) \in [0; 8(1 + \sqrt{2})] \times [0; \frac{\pi}{2}]$, $(x_0, y_0) = (8, 8)$,
3. shifted configuration 2: the eye of the vortex is close from the pole leading to a non null but singular velocity field at the pole, physical space: $(r, \varphi) \in [0; 11] \times [0; 2\pi]$, $(x_0, y_0) = (1, 1)$.

Unless otherwise stated, the Courant number is set to 0.9 and all simulations are carried out with the 3rd order scheme with hyperviscosity ($C_\beta = C_\nu = 1$) till a time $t = 1$. The error is measured by a space-time $L^1$ norm.

Table 1 reports the EOCs for the three configurations. For the centered configuration, we get an EOC of 4. Indeed, in this configuration, the remap step has no influence. Moreover, since the CFL constraint is parabolic, and all interpolations are done on centered stencils, the overall order of the Lagrangian step is always even. In this case, we get better results than the same simulation done with Cartesian coordinates on a slab mesh, but the restitution time is much longer. Regarding the first shifted configuration, we see that the theoretical order is reached. We can take advantage of the parabolic type CFL by increasing the order of the projection step. We ran the same simulation with a 5-th order projection step (only 2 more cells are added to the stencil) and we get an EOC of 4. Finally, for the second shifted configuration, the second order is achieved (results are not presented in this paper). But for the 3-rd order scheme (with a 5-th order projection step), using classical Courant number doesn't lead to a satisfying high-order scheme. Indeed, the geometric singularity together with the singular velocity field pollute measures of the $L^1$ norm in the vicinity of $r = 0$. Reducing the Courant number lead to an effectively high-order scheme. Another satisfying point is that hyperviscosity doesn't affect the EOCs.

**Sod shock tube** We consider the well-known Sod test-case in a two-dimensional cylindrically symmetric geometry. The physical space is: $(r, \varphi) \in [0; 1] \times [0; \frac{\pi}{2}]$, and the initial data are configured to initially let the waves propagate toward $r = 0$.

**Table 1** EOCs for the vortex test-case

| Mesh (N × N) | Config. 1 | Config. 2 Proj. order 3 | Config. 2 Proj. order 5 | Config. 3 (CFL=0.9) | Config. 3 (CFL=0.009) |
|---|---|---|---|---|---|
| 16 | 2.5e-1 | 3.0e-0 | 2.8e-0 | 7.9e-1 | 7.7e-1 |
| 32 | 2.0e-2 3.63 | 5.3e-1 2.51 | 4.1e-1 2.79 | 9.7e-2 3.02 | 7.1e-2 3.44 |
| 64 | 1.4e-3 3.89 | 7.8e-2 2.76 | 3.3e-2 3.61 | 1.7e-2 2.53 | 4.5e-3 3.97 |
| 128 | 8.9e-5 3.96 | 1.1e-2 2.86 | 2.2e-3 3.90 | 3.8e-3 2.17 | 2.9e-4 3.98 |
| 256 | 5.6e-6 3.99 | 1.4e-3 2.95 | 1.4e-4 3.94 | 9.1e-4 2.05 | 2.3e-5 3.61 |
| 512 | 3.5e-7 3.99 | 1.7e-4 2.97 | 9.2e-6 3.97 | 2.3e-4 1.99 | |



**Fig. 1** Density at time $t = 0.5$ for the Sod test-case at third order without (*left*) and with (*right*) hyperviscosity, zoom on the $[0; 0.4]^2$ domain

We run computations with the 3rd-order scheme until $T_f = 0.5$ in order to let the shock be reflected and cross the contact discontinuity. The Courant number is set to 0.7. When hyperviscosity is enabled, we choose the following constants: $C_\beta = 2$, $C_\nu = 1$. Figure 1 plots the density at final time. We see that cylindrical symmetry is preserved. Moreover, when hyperviscosity is enabled, the oscillations in the vicinity of discontinuities are noticeably reduced.

**Two-dimensional Riemann problem** This test-case has been proposed in [4] and studied in [3] in the case of cylindrical coordinates with a polar mesh. The simulation is initialized using piece-wise constant data in each of the four quadrants defined by the x- and y-axis. Let the north-eastern quadrant having the index 1, the others are labeled counter-clockwise with ascending index. The initial data are:

$$
\begin{bmatrix} \rho = 1 \\ u_x = 0 \\ u_y = 1 \\ p = 1 \end{bmatrix}_1, \quad
\begin{bmatrix} \rho = 2 \\ u_x = 0 \\ u_y = -0.3 \\ p = 1 \end{bmatrix}_2, \quad
\begin{bmatrix} \rho = 1.0625 \\ u_x = 0 \\ u_y = 0.2145 \\ p = 0.4 \end{bmatrix}_3, \quad
\begin{bmatrix} \rho = 0.5197 \\ u_x = 0 \\ u_y = 0.2741 \\ p = 0.4 \end{bmatrix}_4.
$$

**Fig. 2** Density at final time for the two-dimensional Riemann problem at first (*left*) and third (*right*) order

With this initial data, a shock-wave propagates between quadrants 2 and 3, a rarefaction between 1 and 4, and two contact discontinuities between quadrants 3 and 4, and 1 and 2. In the vicinity of $r = 0$, the four solutions join each other and a vortex-like flow appears. The test-case is set on the disk of $\frac{\sqrt{2}}{2}$ radius. The Courant number is set to 0.4 and the final time to 0.2. The mesh has 282 cells in the radial direction and 360 cells in the azimuthal direction. Figure 2 plots 30 isolines of the density between 0.525 and 2.025 in cylindrical geometry with the 1-st order scheme (Acoustic Riemann Solver and Godunov splitting) and the 3-rd order scheme with hyperviscosity ($C_\beta = 2, C_v = 1$). We see that the shock-wave is sharper for the 3-rd order scheme than for the 1-st order scheme. Moreover, the complex flow is correctly resolved with the 3-rd order scheme.

## 5 Conclusions

An arbitrary space-time high-order scheme for the resolution of Euler's equations in cylindrical coordinates has been proposed. Various test-cases assess the EOCs for various order of the scheme as well as robustness. To our knowledge, it is the first attempt to derive such a Lagrange-remap scheme for cylindrical coordinates on a polar mesh without any polar symmetry assumptions. In following works, spherical coordinates will be studied.

# References

1. Cook, A., Cabot, W.: Hyperviscosity for shock-turbulent interactions. J. Comput. Phys. **203**, 379–385 (2005)
2. Duboc, F., Enaux, C., Jaouen, S., Jourdren, H., Wolff, M.: High-order dimensionally split Lagrange-remap schemes for compressible hydrodynamics. C.R. Acad. Sci. Paris 348, 105–110 (2010)
3. Illenseer, T.F., Duschl, W.J.: Two-dimensional central-upwind schemes for curvilinear grids and application to gas dynamics with angular momentum. Comp. Phys. Comm. **180**, 2283 (2009)
4. Kurganov, A., Tadmor, E.: Solution of two-dimensional riemann problems for gas dynamics without riemann problem solvers. Num. Meth. Partial Differ. Equ. **18**, 584–608 (2002)
5. Titarev, V., Toro, E.: Ader: arbitrary high-order Godunov approach. J. Sci. Comput. **17**, 609–618 (2002)
6. Wolff, M.: Mathematical and numerical analysis of the resistive magnetohydrodynamics system with self-generated magnetic field terms. Ph.D. thesis, Université de Strasbourg (2011)
7. Wolff, M., Jaouen, S., Jourdren, H., Sonnendrücker, E.: High-order dimensionally split Lagrange-remap schemes for ideal magnetohydrodynamics. In: Proceedings of Numerical Models for Controlled Fusion (NMCF'09), Porquerolles, France (2009)
8. Yee, H., Sandham, N., Djomehri, M.: Low-dissipative high-order schock-capturing methods using characteristic-based filters. J. Comput. Phys. **150**, 199–238 (1999)

# A Combined Finite Volume Discontinuous Galerkin Approach for the Sharp-Interface Tracking in Multi-Phase Flow

**Stefan Fechter and Claus-Dieter Munz**

**Abstract** In this paper, a numerical method for the simulation of compressible two-phase flows is presented. The multi-scale approach consists of several components that allow to sharply resolve the discontinuous nature of multi-phase flow: A discontinuous Galerkin solver for the macroscopic scales of the flow, a micro-scale Riemann solver at the interface that supplies the necessary interfacial jump conditions, a ghost-fluid based coupling of the interfacial conditions to the flow, and a level-set interface tracking formalism. To be able to locally guarantee a sharp and stable resolution at the interface, a finite volume technique on an adaptive sub-cell refinement is applied. The capabilities of the method are demonstrated for a three-dimensional shock-droplet interaction problem.

## 1 Introduction

In many technical applications, multi-phase flows meet conditions such as high pressure environments and/or high velocities that prohibit the popular assumption of incompressibility. Important examples for such extreme ambient conditions include fuel injection systems of aeronautical, automotive and rocket engines that are used at high-pressure operating conditions. The numerical simulation of compressible multi-phase flow is much more difficult than the incompressible treatment, because all conservation equations are coupled via the equation of state (EOS) and have to be solved simultaneously in a consistent way. In the commonly used incompressible treatment hydrodynamics and thermodynamics are decoupled.

S. Fechter · C.-D. Munz (✉)
Institute of Aerodynamics and Gas Dynamics, University of Stuttgart,
Pfaffenwaldring 21, 70569 Stuttgart, Germany
e-mail: munz@iag.uni-stuttgart.de

S. Fechter
e-mail: fechter@iag.uni-stuttgart.de

Three elements are crucial for compressible multi-phase solvers: The first is a numerical method to cope with the large discontinuities in the flow field. The second is the sharp resolution of the interface and the determination of the proper interfacial conditions. Here, we apply a ghost-fluid method as in [2] and supply the interface jump conditions by local Riemann solvers [1]. The third includes the accurate description and tracking of the phase interface that allows for the localisation as well as an estimation of the interface curvature. Here, a level-set method [6] is chosen, as it is easily applicable in the context of high-order methods.

In Sect. 2 these building blocks of the numerical method are described. In Sect. 3 the shock droplet interaction problem is shown, followed by a short conclusion.

## 2 Building Blocks for Sharp Interface Tracking

To be able to include local interfacial phenomena such as surface tension or phase change into the flow simulation, a heterogeneous multi-scale approach is considered, which is based on the solution of the Riemann problem. In the following we neglect viscosity and consider, for simplicity, the Euler equations as macro scale model together with suitable equation of states for the accurate description of multi-phase flows.

The numerical solution of the macro-scale model is provided by a discontinuous Galerkin spectral element method (DGSEM) with an explicit time approximation as described in [3]. At the interface position a Riemann solution is calculated. In case of phase transition the usual solution of the Riemann problem, consisting of four constant states separated by simple waves, brakes down. Information from the micro scale has to be used to get a thermodynamical consistent approximation. With micro scale we denote information from smaller scales that are not resolved by the numerical scheme, e. g. from molecular theory at the interface. Hence, the coupling between the micro and macro scale model is done via such a solution of the Riemann problem. The sharp approximation of the interface is accomplished in the flow solver by shifting the interface always away from the grid cell boundary and calculating the numerical flux within the flow solver by the classical Riemann problem only. Note that this ghost-cell approach does not preserve the conservativity locally.

In the following we describe the basic steps of the sharp interface treatment:

Step 1:    Computation of the interfacial curvature based on the level-set solution.
Step 2:    Solution of the multi-phase Riemann problem at the interface, identified by the level-set function. The data of the Riemann solver is given by interpolated values before and behind the interface. The Riemann solver takes the interface curvature into account, allows a general equation of state, and in the case of phase transition the local solution is established by additional local information from the micro-scala, see [1]. We call this the micro Riemann solver. The local interface velocity is an additional output parameter of the micro Riemann solver, which is then used to describe the interface motion in the level-set equation.

Step 3: The explicit DGSEM is used to advance the flow field to the next time level $t^{n+1}$. Via the ghost-cell approach the flow solver has only grid cells within the bulk phases and rely on standard Riemann solvers for general equation of state. The flux at the phase interface is given by the micro Riemann solver.

Step 4: The new position of level-set zero is used to determine the new position of the interface. In case the interface has moved across a grid cell, the new state is extrapolated using the adjacent grid cells with the same fluid.

Step 5: Based on the refinement indicator that takes the local level-set value as well as a Persson oscillation indicator into account, the refinement is updated.

These are the basic steps of the sharp interface treatment, which are explained in the following in more details.

## *2.1 The Discontinuous Galerkin Spectral Element Method*

In this section we describe the discontinuous Galerkin spectral element method for the flow equations. The description of the method is kept short, for more details we refer to Hindenlang et al. [3].

The key properties of the method are the following. The three-dimensional domain is divided into non-overlapping hexahedral elements, each mapped onto a reference cube element $E := (-1, 1)^3$ by a mapping $\zeta(x)$. The conservation equations on this reference element read as

$$JU_t + \operatorname{div} F(U) = 0 \tag{1}$$

with a flux $F(U) = \big(F^1(U), F^2(U), F^3(U)\big)$ and with Jacobian $J$ of the transformation onto the reference cube. The approximate solution has the form

$$U_h(t, \zeta) = \sum_{i,j,k=0}^{N} \hat{U}(t)_{ijk} \psi_{ijk}(\zeta), \qquad \psi_{ijk}(\zeta) = l_i(\zeta^1) l_j(\zeta^2) l_k(\zeta^3), \tag{2}$$

where $l_j(\zeta)$ are 1D Lagrange polynomials of degree $N$ defined as:

$$l_j(\zeta) = \prod_{\substack{i=0 \\ i \neq j}}^{N} \frac{\zeta - \zeta_i}{\zeta_j - \zeta_i}, \qquad j = 0, \dots, N. \tag{3}$$

Here the points $\zeta_j, j = 0, \dots, N$ are the Gauss-Legendre or Gauss-Legendre-Lobatto points in the reference cube $E$. These points are named interpolation points, the basis is the corresponding tensor product basis, and $\hat{U}(t)_{ijk}$ are the time-dependent degrees of freedom. Multiplying (1) with a test function $\phi$ and integration by parts of the second term leads to three contributions: A volume integral of the

time derivative term (*a*), a surface integral term (*b*) and a volume integral term (*c*), which now contains the gradient of the test function $\phi$:

$$\underbrace{\frac{\partial}{\partial t} \int_E JU_h \phi \, d\zeta}_{a} + \underbrace{\int_{\partial E} \left( F^* \cdot N \right) \phi \, dS}_{b} - \underbrace{\int_E F(U_h) \cdot \mathrm{grad}(\phi) \, d\zeta}_{c} = 0 \, . \qquad (4)$$

Volume as well as surface integrals are approximated by Gauss-Legendre or Gauss-Legendre-Lobatto quadrature. Hence, the quadrature points coincide with the interpolation points. As no continuity constraint is enforced between the elements, the flux function $F(U)$ at the cell boundaries is replaced by a numerical flux function $F^*(U^-, U^+)$ depending on the left and right adjacent states $U^-$ and $U^+$.

In the Galerkin approach the test functions are identical to the basis functions $\phi = \psi_{ijk}$. All the integrals are split in the different coordinate directions and approximated by Gauss-Legendre (-Lobatto) quadrature, which introduces the integration weights $\omega_i$, $i = 0, \ldots, N$. As the quadrature points are chosen to be the same as the interpolation points, the Lagrange property $l_j(\zeta_i) = \delta_{ij}$; $i, j = 0, \ldots, N$ can be exploited. The final semi-discrete form of the DGSEM scheme reads as

$$\left( \hat{U}_{ijk} \right)_t = - (J_{ijk})^{-1} \left[ - \sum_{\lambda=0}^{N} \frac{\omega_\lambda}{\omega_i} D_{i\lambda} F^1_{\lambda jk} - \sum_{\mu=0}^{N} \frac{\omega_\mu}{\omega_j} D_{j\mu} F^2_{i\mu k} - \sum_{\nu=0}^{N} \frac{\omega_\nu}{\omega_k} D_{k\nu} F^3_{ij\nu} \right.$$

$$+ \left( [F^* \hat{s}]_{jk}^{+\zeta^1} \frac{l_i(1)}{\omega_i} - [F^* \hat{s}]_{jk}^{-\zeta^1} \frac{l_i(-1)}{\omega_i} \right) + \left( [F^* \hat{s}]_{ik}^{+\zeta^2} \frac{l_j(1)}{\omega_j} - [F^* \hat{s}]_{ik}^{-\zeta^2} \frac{l_j(-1)}{\omega_j} \right)$$

$$+ \left. \left( [F^* \hat{s}]_{ij}^{+\zeta^3} \frac{l_k(1)}{\omega_k} - [F^* \hat{s}]_{ij}^{-\zeta^3} \frac{l_k(-1)}{\omega_k} \right) \right] .$$

The numerical fluxes $F^*$ are evaluated at the faces of the reference element in each coordinate direction. These terms are denoted by $[]^{-\zeta^1}$ and $[]^{+\zeta^1}$ for the left and right face in $\zeta^1$-direction and analogously for $\zeta^2$ and $\zeta^3$. With $D_{ij} = dl_j(\zeta)/d\zeta|_{\zeta=\zeta_i}$ a differentiation matrix is denoted, which is needed for the integrand of the volume integral. This semi-discrete formulation is then approximated in time by an explicit fourth-order Runge-Kutta scheme.

## 2.2 Adaptive Mesh Refinement (AMR) and Finite Volume Subcells

The advantage of the DG scheme is that high-order approximations can be applied on coarse grids, which is very efficient with respect to the computational effort. With this approach difficulties occur at any discontinuity as e.g. at a phase interface. The continuous in-cell resolution of the DG scheme is not favorable to resolve jump at the phase interface. Our approach to overcome this problem is to replace the coarse DG grid cells by multiple subcells on which a second-order finite volume scheme is applied in the vicinity of the phase interface. The subcell refinement is done such that the number of the degrees of freedom remain the same to avoid a negative impact

**Fig. 1** Schematic representation of a typical setting in our sharp interface approach, involving the liquid-vapor interface, approximated by the zero level-set $\Phi = 0$, the computationally approximated computational interface aligned with the element boundaries and the different ways to apply the numerical fluxes provided by either a standard Riemann solver (bulk phase) or the micro Riemann solver at the computational interface. The white dots visualize the surface integration points

to the global time step restriction. The refinement and the flux calculation at the interface are visualized in Fig. 1 by showing the difference in the interface resolution with and without use of finite volume subcells.

This approach can be efficiently included into the DGSEM description. The coarse grid cell is now considered as a subdomain, in which a finite volume scheme is applied. The coupling to the neighbors is simply the weak coupling of the DG approach. Inside the grid cell the spectral scheme is replaced by a finite volume scheme on the sub-grid. The subcell ansatz enters the DGSEM description (4) in terms of a modified volume integral only. Instead of the continuous DG volume integral we calculate the sum of surface contributions for the equidistant subcell FV cells. This can be written in the following way

$$\int_E F(U) \cdot \text{grad}(\phi) \, d\zeta = \sum_{i,j,k=1}^N \int_{\partial e} \tilde{F} \cdot \tilde{n}_\xi \, dS, \tag{5}$$

where $\partial e$ is the surface area of the subcell FV-cell $e$. This approach allows discontinuities between each subcell as no continuity constraint is enforced. At the interface between DG and FV cells, a conservative flux projection and interpolation method is chosen.

## 2.3 The Level-Set Interface Tracking Method

For interface tracking an additional conservation equation is solved. The level-set advection equation as introduced by Sussman [6] is recast to a conservation equation. This is done to be able to solve this equation with the DGSEM allowing for a high

order approximation. As compressible flow is considered here, an additional right-hand side term has to be solved that can be estimated using the high-order ansatz polynomials

$$\frac{\partial \Phi}{\partial t} + \text{div} \left( s_{\text{PB}} \Phi \right) = \Phi \, \text{div}(s_{\text{PB}}) . \tag{6}$$

The level-set distribution is solely important within a small region around the interface, where geometry and the secondary interface quantities are needed. Outside this region only the sign of the level-set function is important but not the magnitude. A narrow-band approach is used for the advection of the level-set function $\Phi$.

Every 50–100 iterations (depending on the problem), the level-set is redistanced to a signed distance function to be able to accurately estimate the curvature $\kappa$ (second derivative of the level-set function). This procedure resets the level-set function to a numerically preferable shape. The used algorithm is based on the constrained level-set reinitialization equation [6] that is discretized with a 5th order WENO scheme as described by Jiang and Peng [5].

For an accurate estimation of the curvature, the discontinuous Galerkin level-set solution is reconstructed using a $P_n P_m$ method to a polynomial of $M = 3N$. This is done to enhance the accuracy of the curvature calculation. The reconstruction reduces the negative impact of the discontinuous states at the element boundaries and allows for a element-local gradient estimation.

## 2.4 Coupling at the Phase Interface

The consistent numerical and thermodynamic approximation of the phase interface is provided by the solution of an approximate Riemann problem as described in [1]. The used linearized Lax curve Riemann solver has comparatively low computational costs and solely needs an estimation of the sound speed in both bulk phases at the interface. The effects of phase transfer can be included into the Riemann solution as described by Zeiler and Rohde in [7] for the isothermal case.

Surface tension effects are taken into account by a pressure jump according to the Young-Laplace law, for which the mean curvature at the interface is needed as input parameter. The user interface approximation is based on the use of non-conservative fluxes to ensure a sharp interface at all times. The interface propagation velocity $s_{\text{PB}}$ is an additional output parameter and this approach allows for a more general treatment of the interface.

## 3 Computational Results

We show the capabilities of the numerical method by a shock-droplet interaction problem. The equation of state of a perfect gas is applied in the gaseous phase, while the Tait equation is used in the liquid phase. The initial conditions of the

**Fig. 2** Result of a 3D water droplet interacting with a planar shock at various time instances. *Left* pressure contours in the range of $-20$ to $40$ atm. The *solid white line* indicates the interface position. *Right* Schlieren type image of the logarithmic density gradient $\log(\nabla\rho + 1)$

pre- and post-shock states are chosen according to Hu [4] featuring a $Ma = 3$ shock wave impacting on a initially spherical droplet. We consider here a three-dimensional shock-droplet interaction problem. At the domain boundaries in $y$ and $z$-direction a wall boundary condition is assumed. The droplet's initial position is $(0.55, 0, 0)$, the initial position of the planar $M = 3$ shock is set to $x = 0.35$ inside a computational domain that extends $(0, -0.7, -0.7) \times (1.2, 0.7, 0.7)$. The non-dimensional parameters as described by Hu [4] are used in this test case. The chosen

numerical resolution was $80 \times 90 \times 90$ grid cells in the respective axis directions with a DG approximation of order four.

A Persson indicator based on the density is used for shock capturing purposes, which reliably detects the shock position. Depending on the indicator value, the time update is calculated using the fourth order accurate DG scheme in smooth regions or otherwise the TVD stable second order finite-volume scheme, which copes with strong discontinuities and shocks.

The results for simulation times of 2, 4 and $8\,\mu s$ are shown in Fig. 2 in terms of a pressure plot and a Schlieren-type density gradient visualization. The introduced deformations of the droplet as well as the pressure and density-gradient visualization are in agreement to the reference simulations of Hu [4]. They conducted a higher resolved 2D simulation of the problem whereas here a slightly lower resolved 3D problem is considered.

## 4 Conclusion

In this paper, we introduced a numerical method for compressible two-phase flows using a sharp interface method. We applied the method to the simulation of a three-dimensional shock-droplet interaction. The present numerical approach allows for high order of accuracy as well as efficient calculations. The high order is especially advantageous in smooth parts of the flow and for the resolution of the interface as well as its curvature within the level-set approach. The sharp resolution of the interface is established by ideas from the ghost-fluid approach, adapted to the discontinuous Galerkin framework. At the interface the solution of a two-phase Riemann problem is used to get information about the interface states and propagation velocity.

## References

1. Fechter, S., Jaegle, F., Schleper, V.: Exact and approximate Riemann solvers at phase boundaries. Comput. Fluids **75**, 112–126 (2013)
2. Fedkiw, R.P., Aslam, T., Merriman, B., Osher, S.: A non-oscillatory Eulerian approach to interfaces in multimaterial flows (the ghost fluid method). J. Comput. Phys. **152**(2), 457–492 (1999)
3. Hindenlang, F., Gassner, G., Altmann, C., Beck, A., Staudenmaier, M., Munz, C.D.: Explicit discontinuous Galerkin methods for unsteady problems. Comput. Fluids **61**, 86–93 (2012)
4. Hu, X.Y., Adams, N.A., Iaccarino, G.: On the HLLC Riemann solver for interface interaction in compressible multi-fluid flow. J. Comput. Phys. **228**(17), 6572–6589 (2009)
5. Jiang, G.S., Peng, D.: Weighted ENO schemes for Hamilton-Jacobi equations. SIAM J. Sci. Comput. **21**(6), 2126–2143 (2000)
6. Sussman, M., Smereka, P., Osher, S.: A level set approach for computing solutions to incompressible two-phase flow. J. Comput. Phys. **114**(1), 146–159 (1994)
7. Zeiler, C., Rohde, C.: A relaxation Riemann solver for compressible two-phase flow with phase transition and surface tension (2013)

# Numerical Simulation of an Incompressible Two-Fluid Model

**Michael Ndjinga, Thi-Phuong-Kieu Nguyen and Christophe Chalons**

**Abstract**  We investigate some finite volume methods for the numerical simulation of a flow involving two incompressible phases in mechanical disequilibrium. The model consists of two hyperbolic equations with characteristic fields that are neither linearly degenerate nor genuinely nonlinear. We show that the system may involve sonic points, hence the importance of using entropic schemes to accurately capture the volume fraction waves. We propose a Godunov scheme and a Roe scheme with a Harten type correction and compare them on test cases involving the transition between two phase and single phase flows.

## 1 The Model

The flow regime involved in nuclear reactor thermalhydraulics may be single or two phase. More precisely the flow in the reactor core is purely liquid in normal operating condition, a liquid-gas mixture in incidental conditions or purely gaseous in the case of a severe accident involving a total core dewatering. The simulation of the single phase/two phase transition is numerically challenging and has been a major difficulty in the design of new simulation platforms based on advanced two-fluid models (see [3, 5]). Roe type schemes give unphysical solutions, sometimes with negative volume fraction.

M. Ndjinga (✉) · T.-P.-K. Nguyen
CEA-Saclay DEN/DM2S/STMF/LMEC, 91191 Gif-sur-Yvette, France
e-mail: michael.ndjinga@cea.fr

T.-P.-K. Nguyen
e-mail: thi.nguyen@cea.fr

C. Chalons
LMV, UMR 8100, Université Versailles St-Quentin-en-Yvelines, 78035 Versailles, France
e-mail: christophe.chalons@uvsq.fr

We consider a one dimensional isentropic two phase flow involving two fluids 1 and 2 with pressures $P_1$ and $P_2$, densities $\rho_1(P_1)$ and $\rho_2(P_2)$, sound speeds $c_1(P_1)$ and $c_2(P_2)$, volume fractions $\alpha_1$ and $\alpha_2$ (with $\alpha_1 + \alpha_2 = 1$), and velocities $u_1$ and $u_2$. The phasic mass and momentum balance equations yield the following four equation model (see [3, 5, 6, 9])

$$\begin{cases} \partial_t \alpha_1 \rho_1 + \partial_x (\alpha_1 \rho_1 u_1) = \Gamma_1, \\ \partial_t (\alpha_1 \rho_1 u_1) + \partial_x (\alpha_1 \rho_1 u_1^2) + \alpha_1 \partial_x P_1 = \alpha_1 \rho_1 g + \Gamma_1 u_{int}, \\ \partial_t \alpha_2 \rho_2 + \partial_x (\alpha_2 \rho_2 u_2) = \Gamma_2 = -\Gamma_1, \\ \partial_t (\alpha_2 \rho_2 u_2) + \partial_x (\alpha_2 \rho_2 u_2^2) + \alpha_2 \partial_x P_2 = \alpha_2 \rho_2 g + \Gamma_2 u_{int}, \end{cases} \quad (1)$$

where $g$ is the gravitational acceleration. The phase change is considered through the function $\Gamma_1(x) = -\Gamma_2(x)$, and the interfacial velocity is chosen to be $u_{int} = \alpha_1 u_2 + \alpha_2 u_1$. Unlike [3, 5, 9] we do not introduce an interfacial pressure default $\triangle p \partial_x \alpha_k$, but instead a non zero pressure difference of the form $P_1 - P_2 = \frac{\rho_1 \rho_2}{2(\rho_1 - \rho_2)}(u_1 - u_2)^2$ which yields a hyperbolic system. This pressure gap corresponds to a dynamic surface tension model accounting for the fact that velocity shear yields an increase of the microscale interfacial curvature via the well-known Kelvin-Helmholtz instability (see [1]). Taking into account surface tension, the increase of local curvature results in a pressure difference via the Laplace law $P_1 - P_2 = \gamma \sigma$ which should vanish only when $u_1 = u_2$. The kinetic energy gap $\frac{1}{2} \rho_1 u_1^2 - \frac{1}{2} \rho_2 u_2^2 = \frac{1}{2} \frac{(\rho_1 u_1 - \rho_2 u_2)^2}{\rho_1 - \rho_2} - \frac{1}{2} \frac{\rho_1 \rho_2 (u_1 - u_2)^2}{\rho_1 - \rho_2}$, is related to the momentum gap $\rho_1 u_1 - \rho_2 u_2$ and to the velocity gap $u_1 - u_2$. In this first study, we make the simple assumption that the pressure gap exactly compensates the contribution of the velocity gap to the kinetic energy gap.

The system (1) has four main unknowns: $\alpha_1$, $P_1$, $u_1$, $u_2$. The other unknowns can be obtained using the equations of state $\rho_k(P_k)$, $c_k(P_k)$ and the pressure gap law $P_1 - P_2 = \frac{\rho_1 \rho_2}{2(\rho_1 - \rho_2)}(u_1 - u_2)^2$.

Defining the mixture sound wave $c_m = \sqrt{\frac{(\alpha_1 \rho_2 + \alpha_2 \rho_1) c_2^2 c_1^2}{\alpha_1 \rho_2 c_2^2 + \alpha_2 \rho_1 c_1^2}}$, we can compute the Taylor expansion of the system eigenvalues when $u_1 - u_2 \ll c_m$ and the system has four real eigenvalues: two acoustic waves $\frac{\alpha_1 \rho_2 u_1 + \alpha_2 \rho_1 u_2}{\alpha_1 \rho_2 + \alpha_2 \rho_1} \pm c_m + O\left(\frac{u_1 - u_2}{c_m}\right)$ and two volume fraction waves $\frac{\rho_1 u_1 - \rho_2 u_2}{\rho_1 - \rho_2}\left(1 - \frac{\rho_1 \rho_2}{(\alpha_1 \rho_2 + \alpha_2 \rho_1)^2}\right) + O\left(\frac{u_1 - u_2}{c_m}\right)$ and $\frac{\rho_1 u_1 - \rho_2 u_2}{\rho_1 - \rho_2} + O\left(\frac{u_1 - u_2}{c_m}\right)$.

In order to study more precisely the volume fraction waves involved in our applications, we follow [6] and assume that both phases are incompressible with constant densities $\rho_1$ and $\rho_2$. It is then possible to reduce the number of equations to two by setting

$$K = \alpha_1 u_1 + \alpha_2 u_2, \quad \beta = \alpha_1 \rho_2 + \alpha_2 \rho_1, \quad \omega = \rho_1 u_1 - \rho_2 u_2.$$

Combining the two mass balance equations in (1) and using $\alpha_1 + \alpha_2 = 1$ yields $\partial_x K = \frac{\Gamma_1}{\rho_1} + \frac{\Gamma_2}{\rho_2}$. We thus obtain $K(x, t) = \int_0^x (\frac{1}{\rho_1} - \frac{1}{\rho_2})\Gamma_1(x, t) + \alpha_1(0, t)u_1(0, t) + \alpha_2(0, t)u_2(0, t)$ which is entirely determined by the boundary conditions. Using the new unknowns $\beta$ and $\omega$, we can rewrite the system (1) as:

$$\partial_t U + \partial_x F(U) = G(U), \quad \text{with} \tag{2}$$

$$U = \begin{pmatrix} \beta \\ \omega \end{pmatrix}, \, F(U) = \begin{pmatrix} \frac{-K\rho_1\rho_2}{\beta} + \frac{(\beta-\rho_1)(\beta-\rho_2)\omega}{\beta(\rho_1-\rho_2)} \\ \frac{\omega^2}{2(\rho_1-\rho_2)} \end{pmatrix}, \quad G(U) = \begin{pmatrix} 0 \\ g(\rho_1 - \rho_2). \end{pmatrix}.$$

In the theoretical analysis [8], the existence of a positive solution ($\alpha_1, \alpha_2 \in [0, 1]$) to the Riemann problem was proven for the model (2) in the case where $\Gamma_1 = \Gamma_2 = 0$. In this case, $\partial_x K = 0$ and assuming that the boundary conditions are constant we obtain that $K$ is a constant function of time and space. Using a Galilean change of coordinate $u_1 \to u_1 - K, u_2 \to u_2 - K$ with the constant velocity $K$ we can assume that $K = 0$. In this case the Jacobian matrix $\nabla F$ has two real eigenvalues

$$\lambda_1 = \frac{\omega}{\rho_1 - \rho_2} \left( 1 - \frac{\rho_1\rho_2}{\beta^2} \right), \quad \lambda_2 = \frac{\omega}{\rho_1 - \rho_2},$$

and is diagonalisable provided $(\beta, \omega)$ belongs to the state space $(]\rho_1, \rho_2[ \times \mathbb{R}^*) \cup (\{\rho_1\} \times \mathbb{R}) \cup (\{\rho_2\} \times \mathbb{R})$. $\lambda_1, \lambda_2$ correspond to the volume fraction waves of (1) when $c_1, c_2 \to \infty$ and the corresponding eigenvectors are

$$\mathbf{r}_1 = {}^t(1, 0), \quad \mathbf{r}_2 = {}^t(\beta(\beta - \rho_1)(\beta - \rho_2), \rho_1\rho_2\omega).$$

Since the sign of $\lambda_1$ and $\lambda_2$ is not clear we may expect numerous sonic points ($\lambda_1 = 0$ or $\lambda_2 = 0$) during the numerical simulation of the system. Moreover the signs of $\nabla\lambda_1 \cdot \mathbf{r}_1 = \frac{2\rho_1\rho_2\omega}{(\rho_1-\rho_2)\beta^3}$ and $\nabla\lambda_2 \cdot \mathbf{r}_2 = \frac{\rho_1\rho_2\omega}{\rho_1-\rho_2}$ are not clear, so the characteristic fields associated to $\lambda_1$ and $\lambda_2$ are neither genuinely non linear nor linearly degenerate in general. We may therefore expect a non classical wave structure in the solutions for the Riemann problem. This is for instance the case when a pure phase appears in the solution of the Riemann problem (see Fig. 1 in the next section).

## 2 Numerical Schemes

We now investigate the numerical simulation of the system (2) and show that the basic Roe scheme fails to capture the expected dynamics whereas the Godunov scheme and the Roe scheme with a Harten type correction capture the analytical solution.

We consider a uniform mesh of the computational domain $[0, 1]$ whose $N$ cells are centered at $x_i, i = 1, \ldots, N$. The space step $\Delta x = x_i - x_{i-1}$ is constant whereas the time step $\Delta t(U^n) > 0$ depends on the discrete field $U^n = (U^n_i)_{i=1,\ldots,N}$ which approximates the exact solution $U(x, t)$ at cells $i$ and time $t^n = \sum_{k=0}^{n-1} \Delta t(U^k)$. The time step should satisfy the following CFL condition in order to ensure the stability of the explicit schemes: $\Delta t \leq \frac{\Delta x}{\max_i\{\lambda_1(U_i, U_{i+1}), \lambda_2(U_i, U_{i+1})\}}$, where $\lambda_k(U_i, U_{i+1})$ is the largest value of $|\lambda_k|$ on the path connecting $U_i$ to $U_{i+1}$ using the rarefactions and admissible shock waves defined in [8]. We point out that $\lambda_k(U_i, U_{i+1})$ may

be different from $|\lambda_k(U_i)|$ and $|\lambda_k(U_{i+1})|$ because the characteristic fields are non genuinely nonlinear.

We consider conservative finite volume schemes in the following explicit form:

$$U_i^{n+1} = U_i^n - \frac{\Delta t}{\Delta x}\left(\Phi_{i+1/2}^n - \Phi_{i-1/2}^n\right) + \Delta t\, G(U_i^n), \qquad (3)$$

where $\Phi_{i+1/2}^n$ is the numerical flux function at the interface between cells $i$ and $i+1$, and at time $t^n$. We compute the numerical flux $\Phi_{i+1/2}^n$ using one of the following Riemann solvers.

**Godunov scheme**

$$\Phi_{i+1/2}^n = F(U^*(U_i^n, U_{i+1}^n)),$$

where $U^*(U_i^n, U_{i+1}^n))$ is the value taken by the solution of the Riemann problem between the left state $U_i^n$ and the right state $U_{i+1}^n$ at $x = 0$.

**Roe scheme with a Harten type correction**

$$\Phi_{i+1/2}^n = \frac{F(U_i^n) + F(U_{i+1}^n)}{2} - \left(|A^{Roe}(U_i^n, U_{i+1}^n)| + \mathrm{har}_{i,i+1}^n \mathrm{Id}\right) \cdot \left(\frac{U_{i+1}^n - U_i^n}{2}\right),$$

where $A^{Roe}(U_i^n, U_{i+1}^n)$ is the Roe matrix, (see the Appendix for its expression), and $\mathrm{har}_{i,i+1}^n = C \max\left(|\lambda_1(U_i^n) - \lambda_1(U_{i+1}^n)|, |\lambda_2(U_i^n) - \lambda_2(U_{i+1}^n)|\right)$. If $C = 0$ we recover the standard Roe scheme. However it is well-known that the Roe scheme may capture non admissible solutions (see [4]). Hence we used a constant value $C = \frac{1}{5}$ to include a Harten type entropic correction in the Roe scheme.

## 3 Numerical Results

We present some numerical results obtained with the constant densities $\rho_1 = 1$, $\rho_2 = 3$, which give a good overview of the wave structure. We first show that the Godunov scheme and the Roe scheme with Harten type correction are able to capture the non classical wave structure arising in the Riemann problem involving a pure phase intermediate state. Then we take into account a momentum source term through the classical and challenging case of two phase sedimenting under gravity. In the last test we consider a mass source term modeling the drying out of a liquid occurring in the central part of a nuclear reactor vessel. In all these cases the numerical values of $\alpha_1$ and $\alpha_2$ remain between 0 and 1.

### 3.1 The Riemann Problem

The Riemann problem consists in solving the system (2) with $K = g = 0$ and the initial data

**Fig. 1** The solution of the Riemann problem at time $t = 0.1$ for the initial data $\alpha_1 = \alpha_2 = 0.5$ and $\omega_L = -\omega_R = -5$ (*left*), $\omega_L = -\omega_R = 5$ (*right*)

$$U(x, 0) = \begin{cases} (\beta_L, \omega_L) & \text{if } x \leq 0, \\ (\beta_R, \omega_R) & \text{if } x > 0. \end{cases} \qquad (4)$$

In [8] we proved that this problem admits a unique admissible solution satisfying Liu's criterion (see [7]) with $\alpha_1, \alpha_2 \in [0, 1]$. In the special case where $\omega_L = -\omega_R$, the solution involves a pure phase: the heavier if $\omega_L < 0$, and the lighter if $\omega_L > 0$. It consists of three shocks waves in the former case and two transonic rarefactions in the latter (see Appendix). We present in Fig. 1 the numerical results obtained using the Godunov scheme, the Roe scheme, and the Roe scheme with the Harten type entropy correction presented at Sect. 2. In the second case ($\omega_L > 0$, the original Roe scheme is unable to capture the admissible solution and captures instead an undercompressive shock.

## 3.2 The Sedimentation Problem

This is a classical test case in the assessment of numerical methods in the modelling of counter-current two phase flows with steep transition (see [5]). We consider the model (1) with $K = 0$, $\Gamma_1 = -\Gamma_2 = 0$, $g = -10 \, \text{m/s}^2$ and the following initial and boundary data for $x \in [0, 1]$

**Fig. 2** Volume fraction $\alpha_1$ for the sedimentation problem, transient (*left*) and stationary (*right*) solutions

**Initial data:** $u_1(x, 0) = 0$, $u_2(x, 0) = 0.1$, $\alpha_1(x, 0) = 0.5$, $\alpha_2(x, 0) = 0.5$.
**Wall boundary conditions:** $u_1(0, t) = u_2(0, t) = u_1(1, t) = u_2(1, t) = 0$.

The stationary state expected is $\alpha_1 = 0$ on $[0, 0.5]$ and $\alpha_1 = 1$ on $[0.5, 1]$. The transient result in Fig. 2 (left) shows that the Roe scheme captures an undercompressive shock departing from $x = 1$. This is consistent with the results shown in the previous section since the Riemann problems at the walls yield pure phases intermediate states and a transonic rarefaction fan for the lighter phase. However, the Roe scheme with Harten entropic correction gives a similar result to the Godunov scheme, both of them being consistent with the analysis of the Riemann problem.

We remark that the pure liquid wave and the pure gas wave have different structures, the former being a shock wave and the latter a rarefaction wave. However in some publications [2, 5] the "analytical" solution for this problem is claimed to be composed of two shock waves and used to study the convergence of the numerical methods. We do not believe this statement is true and this is confirmed by the theoretical results in [8] and by other numerical results obtained with the compressible model (1) where we used an interfacial pressure term $\triangle p$ similar to [2, 3, 5, 9].

## 3.3 The Boiling Channel Problem

The boiling channel test case is a simplified description of a nuclear vessel thermalhydraulics in incidental conditions. The inlet water is assumed at saturation and remains liquid in the lower part of the vessel. Due to the heating source term in the core the liquid undergoes phase change and may be purely gaseous in the upper part of the vessel (see for example [3]).

We consider the model (1) with $g = 0$ and the piecewise constant phase change function $\Gamma_1(x) = -\Gamma_2(x) = \Gamma_0 1_{[\frac{1}{3}, \frac{2}{3}]}(x)$, for $x \in [0, 1]$. This is a simple 1D description of a nuclear core dewatering, where we do not detail the energy transfers involved in the phase change but only consider a non zero mass source term $\Gamma_k \neq 0$.

**Fig. 3** The stationary state for the boiling channel problem

In the following numerical test, we choose $\Gamma_0 = 3\rho_2$ with the following initial and boundary conditions

**Initial data:** $\alpha_1(x, 0) = 0$, $u_1(x, 0) = 1$, $u_2(x, 0) = 1$, $\forall x \in [0, 1]$.

**Boundary conditions:** inlet at $x = 0$ with $u_1(0, t) = u_2(0, t) = 1$, $\alpha_1(0, t) = 0$ and outlet at $x = 1$ with Neumann condition.

Figure 3 shows that the Roe scheme and the Roe scheme with a Harten type correction give similar results very close to the analytic solution. This could be expected since this problem involves no transonic rarefaction.

## 4 Conclusion

We have shown in this paper that finite volume Riemann solvers are able to solve systems of balance laws in complex configurations. Our system has non genuinely nonlinear characteristic fields and many sonic points but the Godunov scheme as well as the Roe scheme with Harten type correction give satisfactory results with positive volume fractions, provided the time step is carefully chosen. The ability of Riemann solvers to accurately propagate waves in the computational domain is an important advantage when it comes to simulating boiling or condensation fronts in the nuclear energy thermalhydraulics.

## Appendix: The Roe matrix

We used the Roe matrix $A^{Roe}(U_R, U_L) = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, where

$$\begin{cases} a = \dfrac{w_L + w_R}{2(\rho_1 - \rho_2)}\left(1 - \dfrac{\rho_1 \rho_2}{\beta_L \beta_R}\right), \\ b = \dfrac{1}{2(\rho_1 - \rho_2)}\left[\dfrac{(\beta_L - \rho_1)(\beta_L - \rho_2)}{\beta_L} + \dfrac{(\beta_R - \rho_1)(\beta_R - \rho_2)}{\beta_R}\right], \\ c = 0, \\ d = \dfrac{\omega_L + \omega_R}{2(\rho_1 - \rho_2)}. \end{cases}$$

If $\omega_L + \omega_R \neq 0$, the Roe matrix is diagonalisable. However when $\omega_L + \omega_R = 0$, the eigenvalues are real $(0, 0)$ but the matrix is not diagonalisable because of a Jordan block. This is consistent with the continuous model (2) being not hyperbolic for $\omega = 0$. Whenever $\omega_L + \omega_R = 0$ we take a Roe matrix $A^{Roe}(\omega = 0) = 0$.

## References

1. Chandrasekhar, S.: Hydrodynamic and Hydromagnetic Stability. Dover Publication, New York (1981)
2. Cho, H.K., Lee, H.D., Park, I.K. Jeong, J.J.: Implementation of a second-order upwind method in a semi-implicit two-phase flow code on unstructured meshes. Ann. Nucl. Energy **37**(4), 3403–3412 (2010)
3. Cordier, F., Degond, P., Kumbaro, A.: Phase appearance or disappearance in two-phase flows. J. Sci. Comput. **58** (2014)
4. Harten, A., Hyman, J.M.: Self adjusting grid methods for one-dimensional hyperbolic conservation laws. J. Comput. Phys. **50**, 235–269 (1983)
5. Jeong, J.J., Yoon, H.Y., Cho, H.K., Jim, J.: A semi-implicit numerical scheme for transient two-phase flows on unstructured grids. Nucl. Eng. Des. **238**(12), 3403–3412 (2008)
6. Keyfitz, B.L., Sanders, R., Sever, M.: Lack of hyperbolicity in the two-fluid model for two-phase incompressible flow. Discrete Cont. Dyn. Syst. B **3**, 541–563 (2003)
7. Leveque, R.J.: Numerical methods for conservation laws. ETH Zürich, Birkhäuser, Basel (1990)
8. Ndjinga, M., Chalons, C., Nguyen, T.P.K.: On the Riemann problem for an incompressible two-fluide model. To be submitted.
9. Ndjinga, M., Kumbaro, A., De Vuyst, F., Laurent-Gengoux, P.: Numerical simulation of hyperbolic two-phase flow models using a Roe-type solver. Nucl. Eng. Des. **238**(8), 2075–2083 (2008)

# On Boundary Approximation for Simulation of Granular Flow

**David Neusius, Sebastian Schmidt and Axel Klar**

**Abstract** We introduce a Cartesian cut-cell method to numerically solve a system of granular equations in complicated domains. A non-Newtonian Navier-Stokes model is used, which covers both the dense and dilute regime of granular flow. In a Cartesian cut-cell method, one starts from a Cartesian grid and modifies cells that intersect the boundary. In contrast to adaptive or boundary fitting grids, the cutting process yields only local modifications. Thus, the simple Cartesian finite volume structure can be sustained on the interior. To ensure stability in the presence of arbitrarily small cut cells, a merging process will be used, which will result in a combination of the discretization equations on the algebraic level. An interpolation is used to ensure first order convergence near the boundary. We restrict the presentation of numerical examples to two dimensions, while the method derivation includes the three dimensional case.

## 1 Introduction

Designing robust methods for simulations of complex non-Newtonian fluids on complicated geometries is not trivial. It is a common approach to choose a simple Cartesian or rectilinear grid, possibly with local refinement. When this is used on complicated domains, a smooth boundary is discretized as a "stair"-like structure, if the boundary is not parallel to the Cartesian grid. For a compressible fluid this results

D. Neusius (✉) · S. Schmidt
Fraunhofer ITWM, Fraunhoferplatz 1, 67663 Kaiserslautern, Germany
e-mail: neusius@itwm.fhg.de; david.neusius@itwm.fraunhofer.de

S. Schmidt
e-mail: schmidts@itwm.fhg.de

A. Klar
Technische Universität Kaiserslautern, Fachbereich Mathematik, Postfach 3049,
67653 Kaiserslautern, Germany

in an error in density. Since the fluid properties of the macroscopic granular flow ([6] and Sect. 2) are dependent on the density, this density error easily results in a global error. Thus, in order to obtain a consistent Cartesian grid method on a complicated domain it is necessary to use some modification at the boundary. These approaches are summarized under the expression "Immersed Boundary Methods". Our aim is to develop an immersed boundary method that is applicable to the three dimensional macroscopic granular equations.

Mittal [8] has categorized these methods into continuous and discrete forcing approaches. Methods of the former type have either issues with stiffness or become increasingly complicated with the complexity of the model. Since we apply our work to complex fluids we want to avoid this and have chosen direct imposition. The finite volume version of this, which is called the cut-cell method, is, furthermore, the only of the named approaches that retains strict conservation of all state variables [8].

There have been many applications of the Cartesian cut-cell method in the last decade. It has mostly been used for compressible non-viscous flows, e.g. [5]. Many papers have also shown its applicability to incompressible viscous flows, see [1, 3, 11]. Only very recently people have started using it on the compressible viscous Navier Stokes Equations, see [4]. As far as we know, it has not been applied to non-Newtonian fluids, yet. For a more detailed technical report of our cut-cell method see [9].

## 2 Simplified Hydrodynamic Granular Model

In general, simulation of granular materials is interesting due to its widespread use in industrial processes. When an inside view into a production process is not possible, a simulation is important for better understanding and optimization.

Particle-based simulations of granular material are limited in particle numbers by computation time and storage. This limit may be too small to simulate a complete production process. Furthermore, complex non-spherical and non-uniform particles pose difficult modeling challenges.

The continuum model [6], which is shortly presented here, does not scale in runtime with the number of particles. Moreover, the granular properties required for the simulation, e.g. shear stresses, can be obtained from macroscopic lab experiments. Multiphase flows involving fluids are also possible. For ease of presentation, the method is described via a simplified model.

### 2.1 Continuous Equations

The general framework of the model is the isothermal compressible viscous Navier-Stokes equations, having as unknowns the density $\rho$ and the momentum $\rho u$. The density is scaled to a dimensionless volume fraction, such that $\rho \in [0, \rho_C)$, with

$\rho_C < 1$, see Sect. 2.1. Including some volume force $f$, e.g. the gravity, the general continuous equations are

$$\partial_t \rho + \nabla \cdot (\rho u) = 0$$
$$\partial_t (\rho u) + \nabla \cdot (\rho u u) - \nabla \cdot \sigma + \nabla p - f = 0, \tag{1}$$

with the asymmetric (2) stress strain relation

$$\sigma = \eta \kappa \quad \kappa_{ij} = \frac{\partial u_i}{\partial x_j}. \tag{2}$$

**Closure**

We first require the concept of a granular temperature $T$, as introduced in [2, 6]. $T$ resembles the energy of "random movement" of particles, similar to the temperature being a measure of the random movement of molecules. This similarity indicates that the temperature dominates the dilute regimes, where granulate behaves in many respects like a gas. Thus, it will be a major contribution to all kinetic terms that are introduces later. The higher the granular temperature the more often particles, just as molecules would, will interact.

Interactions between granular particles are non-elastic collisions. This leads to a constant loss of energy. Without outer sources the granular temperature will always converge to zero.

The equation for $T$ is omitted in the simplified model and is in the numerical results replaced by an application dependent constant. One could also use an asymptotic temperature formula in a dense slow regime as given in [12].

A further magnitude required is that of a maximum density $\rho_C$. As the material of which the particles are composed is assumed to be incompressible, there is a maximum packing one can reach without destruction of particles. This value is mainly used within the radial distribution function $g(\rho) = \left(1 - \frac{\rho}{\rho_C}\right)^{-1}$. Having only a continuously resolved particle distribution we need this function to measure the probability of having particles in collision range. An infinite $g(\rho)$ implies a probability of one. Using these, we can define the first part of the pressure $p_k = T\rho g(\rho)$, the kinetic pressure. For low density and constant temperature this resembles the ideal gas law.

With increasing density, the finite radius of our particles requires additional forces that are not present in standard fluid equations. Any material whose temperature approaches zero Kelvin will contract strongly. This does not apply, if a granulate comes to a rest, i.e. the granular temperature converges to zero. Instead, if the gravity is the only external force, it will have a density of no more than half the possible maximum density $\rho_C$. This equilibrium is numerically not reproducible as long as the pressure is always proportional to $T$. Thus, a second part of the pressure, $p_y$ or the yield pressure is introduced:

$$p = p_y + p_k \quad \text{where} \quad \begin{aligned} p_k &= T\rho g(\rho) \\ p_y &= \Theta(\rho - \rho_{C_0})T_0(\rho - \rho_{C_0})g(\rho) \end{aligned} \tag{3}$$

The parameters $T_0$ and $C_0$ are material dependent and may be functional. Among other criteria, they have to be determined by the existence of the previously described equilibrium at the correct density.

Similarly, one can derive a kinetic [2] and yield [10] viscosity.

$$\eta = \eta_k + \eta_y = \eta_k \left( 1 + \frac{p_y}{p_k} \right) \quad \text{where} \quad \begin{aligned} \eta_k &= \eta_0 \sqrt{T} \rho g(\rho) \\ \eta_y &= \eta_0 \Theta(\rho - \rho_{C_0}) \frac{T_0}{\sqrt{T}} (\rho - \rho_{C_0}) g(\rho). \end{aligned} \tag{4}$$

**Stress Strain Relation**

The advantage and main reasoning for using a symmetric stress strain relation is that it ensures conservation of angular momentum. Physically, this conservation cannot be observed in granular flow. Angular momentum can be converted into rotation of a single particle since rotation of single particles in our model will be a temperature rather than a velocity. Thus, as both, the symmetric and the asymetric stress strain relation are not completely correct, we choose the easier asymmetric one which leads to a decoupling of the velocities in the implicit part of the numerical method.

## 2.2 Discrete Equations

Using the asymmetric stress-strain relation, we derive the discrete equations

$$V_C \frac{\partial \rho_C}{\partial t} + \sum_{f \in \text{faces}} \left( A_f \rho (n_f \cdot u_f) \right) = 0 \tag{5}$$

$$V_C \frac{\partial (\rho u)_C}{\partial t} + \sum_{f \in \text{faces}} A_f \left( \eta_f(\rho, T) \frac{\partial u}{\partial n} - p(\rho, T) n_f \cdot I - \rho u (n_f \cdot u_f) \right) = 0, \tag{6}$$

where $V_C$ is the volume of the cell, $A_f$ the area of a face and $n_f$ and $u_f$ are the normal and velocity on a face. Regarding the temporal derivative we use a partly implicit scheme. The advection and pressure are using pure explicit Euler, while the velocity in the diffusion is split as $u = (\rho u)/\rho$. The former part is linear and we can thus easily apply the implicit Euler method, while the latter part is again treated explicitly.

As long as only a Cartesian grid without cut-cell is used, the face velocity and the arguments of the pressure function will be the average of the two adjacent cell-center

values. The state variables in the advection term are chosen according to the upwind-scheme. The temporal and the normal derivatives require the cell-center values for the current and the next time step. Simple cell-center values do not suffice to achieve a higher convergence using the cut-cell method, see Sect. 3.3.

# 3 Cut Cell Method

To apply the cut-cell method we have to consider three main tasks. The first is the creation of a boundary representing mesh, this is called the cutting procedure. The second task is to ensure stability properties similar to the standard Finite Volume methods. As small cut cells would prevent this, they are merged with larger cells. Since the changes done in the first two tasks lead to shifted center points of cells and faces and even the creation of new boundary faces, it is furthermore necessary to apply an interpolation to ensure first order convergence near the boundary.

## 3.1 Mesh Creation

The target of the cutting process is a simple automatically constructed grid, which represent the actual boundary up to a continuous and piecewise linear accuracy, see Fig. 1. The simplest way would be to restrict ourselves to one additional face per cell. The most complicated on the other hand would allow faces to have corners in the interior of the Cartesian cell, as e.g. done by Ahmadi [1]. As an intermediate approach, we allow the triangulation possibilities of the Marching Cubes/Squares algorithm [7]. Thus, a slightly modified version of this algorithm can be applied to compute the cut-faces. In other words our cutting should have the following properties:

We assume that each cut cell is a subset of the underlying Cartesian cell. In 2D, we further assume that each cell is a polygon whose corners are located on the edges of the Cartesian cell and no more than one corner plus the two endpoints are allowed per edge. In 3D, the cell is assumed to be simply connected and its boundary must be a union of polygons, with the same properties as the polygon in the 2D case. The limit on the number of *distinct* corners per edge is applied on the set of corners of *all* these polygons.

## 3.2 Merging of Small Cells

At first sight it might be useful to keep each cut-cell as a separate cell in order to achieve a good approximation at a complicated boundary. Yet, a few very small cells are created by the cutting and would globally require a very low time-step, as e.g. seen in the CFL-condition. There are many ways to counter this: Klein et al. [5] use

**Fig. 1** Illustration of a cut-cell grid around one cell

a flux balancer to stabilize these areas. Tucker et al. [11] sacrifice exact conservation by simply interpolating the values on small cut-cells instead of using the governing equations. As most people, we will use a cell merging instead, where it is necessary to state a lower bound on the volume on every non-merged cell.

This lower bound is a design parameter: One can see that by decreasing the smallest allowed cell size the $L_\infty$-error decreases. This improvement varies for different examples, depending on how many cells are affected. On the other hand, the larger differences in cell sizes will increase the $L_2$ and $L_1$ errors. A short convergence study for a Newtonian example is done in [9].

The merging of two neighboring cells can be seen as the unification of their volumes. The new center point, where the state variables are saved, will be in the centroid of the unified cell. The fluxes on the faces of both cells will now contribute to the same discrete equations. Fluxes on the face between the two merged cells cancel out.

### 3.3 Interpolation

Let us motivate the necessity of an interpolation of quantities by illustrating a very noticeable error that can be observed if we just use cell wise constant quantities.

This can for instance be seen in an example, where we have a straight channel with no-slip boundary condition. In the analytical solution, the velocity would increase with the distance to the channel boundary, but would remain constant if we move parallel along the boundary. Thus, velocity gradients that are almost parallel to the boundary should be very small. This should hold for the normal gradient of the edges $w$ and $e$ in Fig. 2.

Using no interpolation, the discretization of the gradient $\frac{\partial u_w}{\partial n_w} \approx \frac{u_{P_2} - u_{P_1}}{h}$ would indeed be very small. Yet, the discrete normal gradient at $e$, $\frac{\partial u_e}{\partial n_e} \approx \frac{u_M - u_{P_2}}{h}$ would

**Fig. 2** Example, where interpolating the velocity on a cut-face by two center points fails at a non-slip boundary



| $h$ | Cartesian | | Cut cell | |
|---|---|---|---|---|
| $h$ | $\frac{1}{10\pi}L_1$ | $\frac{1}{10}L_\infty$ | $\frac{10^{-5}}{\pi}L_1$ | $\frac{1}{100}L_\infty$ |
| Density Errors | | | | |
| 1/50 | 2.09 | 2.25 | 9.03 | 1.00 |
| 1/100 | 1.54 | 3.19 | 5.82 | 0.852 |
| 1/200 | 1.06 | 3.54 | 2.81 | 0.513 |
| 1/400 | 0.760 | 3.76 | 1.60 | 0.249 |
| Density Convergence Rate | | | | |
| 1/100 | 0.445 | -0.500 | 0.634 | 0.235 |
| 1/200 | 0.539 | -0.150 | 1.05 | 0.731 |
| 1/400 | 0.476 | -0.0873 | 0.810 | 1.04 |

**Fig. 3** Rotation of granular material induced by a uniform tangential velocity on a circular no-slip boundary. Uniform initial condition $\rho = 0.4, u = 0$. Granular parameters $T = 10^{-5}, \rho_C = 0.8, \rho_{C_0} = 0.3, T_0 = 1$. *Left* Cartesian grid. *Right* Presented cut-cell method with lower volume bound (see Sect. 3.2) 1/4, yielding the expected result. *Table* Errors and convergence rates for this example using different uniform grid sizes $h$ in both coordinate directions. A Cut-cell solution with $h = 10^{-3}$ has been used as reference. A Cartesian reference solution cannot be used since there is convergence in $L_\infty$ norm. We have used $h = 0.01$ for the plots

be very large: The centroid $M$ has a much higher velocity than $P_2$, since it is further away from the boundary.

Since these discrete normal gradients have to be used in the diffusion term, this discrepancy would lead to an overly strong diffusion flux on edge $e$. It would cause an artificially high velocity in cell $P_2$ and in the long term an increased density in cell $M$. We have also proofed, see [9], that this error does not decrease with smaller spatial step sizes. Similar errors can be found in other terms of the equation.

An interpolation will be necessary in all face flux terms, if the according face is adjacent to any non-Cartesian cell. It does not disturb conservation, as it is only used for these fluxes that are used in both adjacent cells with a different sign.

The necessary interpolation can be done by a polynomial fit, see e.g. [3], or by a nearest neighbor interpolation, which utilizes a Voronoi partition and is described here shortly: We want to interpolate a quantity $x$ at some arbitrary point $P$ while we know its value on a number of data points $D_0, \ldots, D_n$, which are all center points and possibly boundary face values. We first construct a Voronoi partition $V_0$ from these data points, which can be used multiple times for any $P$. Adding $P$ to $V_0$ we receive a new Voronoi partition $V_P$ and the interpolation formula:

$$x(P) = \sum_{i=0}^{n}(A_{D_i,V_0} - A_{D_i,V_P})x(D_i), \tag{7}$$

where $A_{D,V}$ is the Voronoi cell area/volume (3D) of data point $D$ in Voronoi partition $V$. Since adding $P$ only changes the Voronoi partition locally, it can be inserted with effort $\mathcal{O}(1)$ and only very few summands are non-zero.

## 4 Numerical Example: Rotating Cylinder

In many applications if suffices to add some normal information of the curved domain to a standard Cartesian Finite Volume solver in order to achieve satisfying result. Yet, there are some examples where it fails completely if the cut cell method is not used. In Fig. 3 the density and velocity norm should only depend on the radial distance. Yet, we have unrealistic piling of material in the solver that does not utilize the cut-cell method.

## References

1. Ahmadi, M.: Modelling and Quantification of Structural Uncertainties in Petroleum Reservoirs Assisted by a Hybrid Cartesian Cut Cell / Enriched Multipoint Flux Approximation Approach. Ph.D. thesis (2012)
2. Bocquet, L., Losert, W., Schalk, D., Lubensky, T., Gollub, J.: Granular shear flow dynamics and forces: experiment and continuum theory. Phys. Rev. E **65**(1), 011,307 (2001)
3. Chung, M.H.: Cartesian cut cell approach for simulating incompressible flows with rigid bodies of arbitrary shape. Comput. Fluids **35**, 607–623 (2006)
4. Hartmann, D., Meinke, M., Schröder, W.: A strictly conservative Cartesian cut-cell method for compressible viscous flows on adaptive grids. Comput. Meth. Appl. Mech. Eng. **200**(9–12), 1038–1052 (2011)
5. Klein, R., Bates, K.R., Nikiforakis, N.: Well-balanced compressible cut-cell simulation of atmospheric flow. Phil. Trans. R. Soc. A **367**, 4559–4575 (2009)
6. Latz, A., Schmidt, S.: Hydrodynamic modeling of dilute and dense granular flow. Granular Matter **12**(4), 387–397 (2010). doi:10.1007/s10035-010-0187-6
7. Lorensen, W.E., Cline, H.E.: Marching Cubes: a high resolution 3D surface construction algorithm. Comput. Graph. **21**(4), 163–169 (1987)
8. Mittal, R.R., Iaccarino, G.: Immersed boundary methods. Annu. Rev. Fluid Mech. 37(1), 239–261 (2005). doi:10.1146/annurev.fluid.37.061903.175743
9. Neusius, D., Schmidt, S.: A Cartesian cut-cell method for the isothermal compressible viscous Navier-Stokes Equations. Berichte des Fraunhofer ITWM 231 (2013)
10. Savage, S.: Analysis of slow high-concentration flows of granular materials. J. Fluid Mech. **377**, 1–26 (1998)
11. Tucker, P.G., Pan, Z.: A Cartesian cut cell method for incompressible viscous flow. Appl. Math. Model. **24**, 591–606 (2000)
12. Zemerli, C.: Continuum Mechanical Modeling of Dry Granular Systems : From Dilute Flow to Solid-like behavior. Ph.D. thesis (2013)

# Comparison of Realizable Schemes for the Eulerian Simulation of Disperse Phase Flows

**Macole Sabat, Adam Larat, Aymeric Vié and Marc Massot**

**Abstract** In the framework of fully Eulerian simulation of disperse phase flows, the use of a monokinetic closure for the kinetic based moment method is of high importance since it accurately reproduces the physics of low inertia particles with a minimum number of moments. The free transport part of this model leads to a pressureless gas dynamics system which is weakly hyperbolic and can generate $\delta$-shocks. These singularities are difficult to handle numerically, especially without globally degenerating the order or disrespecting the realizability constraints. A comparison between three second order schemes is conducted in the present work. These schemes are: a realizable MUSCL/HLL finite volume scheme, a finite volume kinetic scheme, and a convex state preserving Runge-Kutta discontinuous Galerkin scheme. Even though numerical computations have already been led in 2D and 3D with this model and numerical methods, the present contribution focuses on 1D results for a full understanding of the trade off between robustness and accuracy and of the impact of the limitation procedures on the numerical dissipation. Advantages and drawbacks of each of these schemes are eventually discussed.

## 1 Introduction

The study of two-phase flows is needed for a wide range of applications such as fluidized beds, spray dynamics, atomization of fuel in combustion chamber, alumina particles in rocket engines, cosmology, etc. In the present contribution, we focus on

M. Sabat (✉) · A. Larat · M. Massot
CNRS UPR 288, Laboratoire d'Energétique Moléculaire et Macroscopique Combustion (EM2C), Ecole Centrale Paris and Fédération de Mathématiques de l'ECP-FR CNRS 3487, Grande Voie des Vignes, 92295 Châtenay-Malabry, France
e-mail: macole.sabat@ecp.fr

A. Vié
Center for Turbulence Research, Stanford University, Stanford, CA 94305-3024, USA

the Eulerian resolution of the disperse phase using the kinetic-based moment method (**KBMM**). It approximates the solution of the Williams-Boltzmann equation (WBE) [12] at a macroscopic level using a finite set of integrated quantities over the phase space, called **moments**. The closure of the KBMMs is based on the choice of a presumed shape in the velocity space for the number density function (NDF), having as many parameters as the number of moments one needs to control [8, 11]. The main advantages of KBMMs are the (weakly)-hyperbolic character of the resulting system and the close link between the transported moments and the physics contained in the underlying NDF. This is of critical importance for numerical scheme design.

In the present work, we consider the case of high Knudsen number where the particle-particle collisions are negligible. Moreover, since one of the most delicate steps in the Eulerian modeling is the velocity closure for the convective part, we will focus only on the transport term in the WBE. It is essential to note that this term is the building block for all the KBMMs. The model studied in this work is the monokinetic closure model (MK) obtained by assuming that the velocity distribution is a Dirac measure [2, 8]. In 1D, it leads to a two equation weakly hyperbolic pressureless gas dynamics (PGD) system. Since the velocity is locally uniquely defined, this model correctly reproduces the dynamics of low inertia particles when no trajectory crossing occurs. The main features of the solution of this model are stiff accumulations regions and large depletion zones, what justifies the search of accurate and robust numerical methods. Indeed, the numerical scheme used can highly influence the captured physics and should not degenerate in the presence of void regions and singularities. In addition, moment methods require that the numerical scheme satisfies the realizability condition (every set of moments has to be associated with a positive NDF) in each cell. This realizability condition translates into the positivity of density for PGD. In the literature, the resolution of the PDG system has already been studied among others by Bouchut et al. [3], Larat et al. [7] and Yang et al. [13].

The main contribution of this work is to compare, for a physical model that takes into account the key part of the transport of the disperse phase, one of the latest developments in the field of numerical methods, a **realizable** new class of **RKDG** with a convex projection strategy, to better known methods such as a realizable MUSCL/HLL finite volume scheme (**MUSCL/HLL**) [11] and a finite volume kinetic scheme (**FVKS**) [3]. For the FV schemes the limitation strategy is assessed by comparing the minmod and the monotonized central-difference (MC) limiters [9]. Given the challenging aspect of the simulation of Dirac solutions, and for the sake of simplicity of the qualitative, quantitative and individual properties comparisons, only 1D space test cases are presented hereafter. The methods generalize to higher space dimensions and higher order KBMM models, the research on which has been fostered by SAFRAN. However, the key features of the methods can be characterized already in 1D, which is the purpose of the present contribution.

## 2 KBMM and Mono-Kinetic Closure Based Model

The first step of modeling is the kinetic approach inspired by the kinetic theory of gases. A statistical description of the disperse phase is used through a NDF $f(t, x, \xi)$, where $t$ is the time, $x$ the position and $\xi$ the internal phase space. The sole choice of the phase space is strongly related to the physics one wants to describe. For example, if we consider that the particles are spherical, $\xi = (c, S, T)$ is the phase space composed of velocity, size and temperature. In this case, the statistical approach leads to a mesoscopic description given by the WBE [12]. This equation contains the free transport of the discrete phase, a term for the acceleration of the particles, a term relative to the evaporation rate, an expression of the rate of change of the particle temperature and the source terms due to breakup and coalescence. Since our focus is on the free transport we will only deal with this part of the WBE (see Eq. (1)) in 1D. Furthermore, for simplicity we will consider a monodisperse phase even though polydispersity could be included through a Multi-Fluid size phase space discretization [8], for example. The transport part of the WBE reads:

$$\partial_t f + \partial_x (cf) = 0 \tag{1}$$

After integrating Eq. (1) over the phase space, one gets a system of moment equations with $M_i = \int U^i f \, dU$ being the general $i$th order moment in velocity:

$$\partial_t M_i + \partial_x (M_{i+1}) = 0. \tag{2}$$

This system is not closed: for every set of $N + 1$ moments, the moment of order $N + 1$ is required as the $N$-th moment flux ($\mathscr{F}(M_i) = M_{i+1}$) [11]: a closure relation $M_{N+1} = f(M_0, \ldots, M_N)$ has to be provided to model the unknown flux. Such a closure depends on the physics one needs to describe. We focus here on the mono-kinetic closure [2, 8]. It correctly reproduces the formation of depletion zones and accumulations regions in the case of low inertia particles for which no particle trajectory crossing occurs [4, 11]. The NDF is assumed to write $f(t, x, c) = \rho(t, x)\delta(c - u(t, x))$, where $u(t, x)$ is the mean velocity of the dispersed phase. The system of moments closes at first order i.e. $N = 1$ and we get therefore the two equation **PGD** system:

$$\left\{ \partial_t \rho + \partial_x (\rho u) = 0 \quad ; \quad \partial_t (\rho u) + \partial_x (\rho u^2) = 0. \right\} \tag{3}$$

This system is weakly hyperbolic and can generate $\delta$-shocks [2]. These singularities are difficult to handle numerically, especially without globally degenerating the order of accuracy. In addition, the physical meaning of the numerical solution relies on the realizability condition: every pair of moments $(\rho, \rho u)$ is associated with a positive NDF. As a result, the positivity of the number density $\rho$ should be preserved. Moreover, the velocity $u$ has to respect a maximum principle [2, 3].

## 3 Numerical Schemes

Three second order schemes are tested here in order to check their ability to meet the accuracy, realizability and robustness requirements. The first scheme is a MUSCL/HLL finite volume scheme [1, 10, 11]. This scheme is obtained using the MUSCL strategy [10] with a linear conservative reconstruction of the primitive variables ($\mathscr{U} = (\rho, u)$) within each cell in order to calculate the interface values. The evaluation of the fluxes is then done with a first order HLL flux:

$$2\,\mathscr{F}^{HLL}(\mathbf{M}_L, \mathbf{M}_R) = \mathscr{F}(\mathbf{M}_L) + \mathscr{F}(\mathbf{M}_R) - |\lambda_m|\left(\mathbf{M}^* - \mathbf{M}_L\right) - |\lambda_M|\left(\mathbf{M}_R - \mathbf{M}^*\right) \tag{4}$$

where $\mathbf{M} = (\rho, \rho u)^T$ is the state of moments, $\mathbf{M}_L$ and $\mathbf{M}_R$ are the initial states at each side of the interface and $\lambda_M$ and $\lambda_m$ are respectively the maximum and minimum eigenvalues of the Jacobian over the cell interfaces: $\lambda_M = \max(\mathscr{F}'_L, \mathscr{F}'_R)$ and $\lambda_m = \min(\mathscr{F}'_L, \mathscr{F}'_R)$. For the integration in time, a strong stability preserving two-step Runge-Kutta (SSP2RK) method [6] is used. The resulting scheme is of second order in time and space and preserves the realizablity of the moments. It has already been used for example on 2D Taylor Green, and homogeneous isotropic turbulence (HIT) test cases for different KBMMs [11]. For more information on this scheme one may refer to the work of Vié et al. [11] and references therein.

The second scheme is the finite volume kinetic scheme (FVKS) [3]. It uses the exact solution in time of the underlying kinetic description and is therefore intrinsically realizable. For the second order scheme, piecewise linear reconstructions are considered for the density and velocity. This scheme was previously used to solve 3D HIT and other combustion applications ([4] and references therein).

The slope limiter used in the first two methods is either a minmod or a MC limiter [4, 9]. These limiters are obtained from Eq. (5) by respectively taking $\alpha = 1$ or $\alpha = 2$ with $\Delta^+\rho = \rho_{i+1}^n - \rho_i^n$, $\quad \Delta^-\rho = \rho_i^n - \rho_{i-1}^n$, $\quad \chi = \Delta x(1 + \Delta x D_{\rho_i}/6\rho_i^n)$

$$D_{\rho_i} = \frac{1}{2}(sgn(\Delta^+\rho) + sgn(\Delta^-\rho)) \times \min\left(\frac{|\Delta^+\rho + \Delta^-\rho|}{2\Delta x}, \frac{\alpha|\Delta^-\rho|}{\Delta x}, \frac{\alpha|\Delta^+\rho|}{\Delta x}\right), \tag{5}$$

$$D_{u_i} = \frac{1}{2}(sgn(\Delta^+u) + sgn(\Delta^-u)) \times \min\left(\frac{|\Delta^+u + \Delta^-u|}{2\chi}, \frac{\alpha|\Delta^+u|}{2\Delta x - \chi}, \frac{\alpha|\Delta^-u|}{\chi}, \frac{1}{\Delta t}\right)$$

The last scheme is a convex state preserving Runge-Kutta discontinuous Galerkin scheme (RKDG) [7, 14]. First, the variational formulation is conducted using $k + 1$ basis functions $\phi_i^j$, polynomials of order k in cell $\mathscr{C}_i$. Then, according to the classical DG formulation, $\mathbf{M}_h$ is the piecewise polynomial solution of the following differential system, where $\mathscr{M}$ is the mass matrix and $\mathscr{F}^*$ is the numerical flux: $\forall i = 1, ..., N; \ \forall j = 1, ..., k$

$$|\mathscr{C}_i|(\mathscr{M}_{jl})d_t\mathbf{M}_i^l + \left(\mathscr{F}_{i+\frac{1}{2}}^* \phi_i^j(x_{i+\frac{1}{2}}) - \mathscr{F}_{i-\frac{1}{2}}^* \phi_i^j(x_{i-\frac{1}{2}})\right) = \int_{\mathscr{C}_i} \mathscr{F}(\mathbf{M}_h)\partial_x\phi_i^j \, dx$$

(6)

By summing over all the degrees of freedom of the cell, one obtains the equation of evolution of the cell mean value. The time update is done using the same SSP2RK method used for MUSCL/HLL. Then, by using an appropriate Gauss-Lobatto quadrature rule the update of the cell mean value can be rewritten into a convex combination of abstract first order updates in the subcells of two neighbouring quadrature points. If the numerical flux is convex state preserving at first order (for example Rusanov flux) the updated set of mean moments is realizable if the solution at each Gauss-Lobatto quadrature points is realizable. This is obtained by reducing the deviation of the local polynomial just enough so that the realizability is met and the accuracy is not destroyed [14]. It is important to note that in the last test case an additional modified minmod slope limiter [5] had to be used to ensure the stability of the method when the solution becomes very singular. However this should not be the case theoretically. This remains an open question. This scheme was already tested on 2D Taylor Green, and homogeneous isotropic turbulence (HIT) test cases for different KBMMs on unstructured grid [7].

## 4 Results

We present here three test cases. We consider periodic boundary conditions for all the tests with CFL = 0.5 and a mesh of 100 cells (except for convergence study). First we want to assess the numerical method implemented with the most simplified version of the PGD system where the velocity is everywhere equal to unity. In this case, the linear advection equation is obtained. We consider a Gaussian-like initial condition:

$$\rho(x, 0) = [\cos(\pi(2x - 1))]^4 \text{ if } 0.25 < x < 0.75, \quad 0 \text{ otherwise; and } u(x) = 1$$

(7)

In Fig. 1 to the left, the solutions of the different schemes are represented after 10 cycles. The solutions of the schemes with the minmod limiter are clearly smeared out. We can also observe the leading phase error for the RKDG and MUSCL/HLL solutions which is a sign of numerical dispersion. For the FVKS MC solution we notice a minor flattening of the bump due to slope limitation. According to these results, the list of the schemes arranged in increasing order of numerical diffusion is: RKDG, FVKS MC, MUSCL/HLL MC, FVKS Minmod and MUSCL/HLL Minmod. We next perform a convergence study at $t = 1$ in the 2-norm, Fig. 1 right side. Correct second order is obtained for RKDG and FVKS with MC limiter, which is not the case for MUSCL/HLL particularly for coarse meshes. When using the minmod limiter instead of MC, the slopes are respectively reduced by 13 and 26 % for MUSCL/HLL and FVKS. Also an interesting feature is that RKDG maintains the exact second

**Fig. 1** Linear advection equation, to the left density at $t = 10$ and to the right convergence study for RKDG, FVKS MC and MUSCL/HLL MC (2-norm)

order in the 1-norm and in the $\infty$-norm, which is not the case of the two other FV schemes.

Remaining test cases solve for the PGD system with two different initial conditions. The **second test** is similar to numerical test I in Bouchut et al. [3]. The initial condition, for $0 \le x \le 2$, is $\rho(x, 0) = 0.5$ and:

$$u(x, 0) = \begin{cases} -0.4 & x < 0.5 \quad \text{or} \quad x > 1.8, \\ 0.4 & 0.5 < x < 1, \\ 1.4 - x & 1 < x < 1.8. \end{cases} \tag{8}$$

The density is plotted in Fig. 2 for the three schemes at $t = 0.5$. The RKDG solution is obtained by guaranteeing the positivity of the density and by limiting the absolute value of the velocity to 0.4. These two conditions define the convex state for this method. We can notice that all the schemes create small overshoots near the discontinuities (after $x = 1.2$ and before $x = 1.6$), these being already observed in Bouchut et al. [3]. These overshoots have the highest amplitude for the MUSCL/HLL (4.5 % above 1.0), are a little bit smaller for RKDG and reduce to only 1 % for the FVKS results. In addition, RKDG creates overshoots before $x = 0.3$ and after $x = 0.7$ since no limiter is used at these points and the scheme does not ensure local monotonicity. When FVKS gives the most satifactory solution compared to the exact, RKDG also shows the sharpest resolution of the discontinuities. Finally, velocity component is not shown, since it is nearly the same for all the schemes.

The **last test** is a more complex problem. It is a replicate, under the MK model, of two packets of particles approaching each other with opposite velocities. For $0 \le x \le 1$, the inital condition is:

**Fig. 2** Second test case: density and velocity results for the different schemes at $t = 0.5$



**Fig. 3** Last test case: density and velocity results for the different schemes at $t = 0.5$

$$\rho(x, 0) = [\sin(2\pi x)]^4 \quad ; \quad u(x, 0) = \begin{cases} -1 & \text{if } x > 0.5 \\ 1 & \text{otherwise} \end{cases} \tag{9}$$

At $t = 0.5$, the density exact solution is a Dirac measure at $x = 0.5$. Therefore, we should have all the matter concentrated in one cell at $x = 0.5$. RKDG result is obtained using an additional modified minmod limiter [5] and the convex constraint is defined as positive density and absolute velocity limited to 1.0. For this test case we consider a mesh of 101 cell in order to have a cell center at 0.5 to check if the schemes capture the right position of the Dirac. The numerical results are shown in Fig. 3.

   All the schemes are able to physically capture the singularity. The major part of the matter is concentrated in three cells for RKDG ([0.4802, 0.5198]: the mid-cell and its two neighboring cells), in two cells for MUSCL/HLL ([0.4901, 0.5198]: the

mid-cell and its right neighbor) but it is concentrated in only one cell for the FVKS result [0.4901, 0.5099]. The FVKS gives the highest density at $x = 0.5$ (37.87). At this point the density using MUSCL/HLL and RKDG is reduced respectively by 21 and 27 %. For a density less than $10^{-12}$ we consider void and set the velocity to zero. From the velocity results, it is seen that vacuum is not generated using FVKS. We have void outside the interval [0.2525, 0.7475] for MUSCL/HLL and [0.3515, 0.6485] for RKDG. The RKDG has the largest interval of vacuum but the FVKS gives the sharpest profile near the velocity discontinuity. According to the performance of RKDG in the previous problems, a better result was expected. For this reason, the same test case was repeated with a CFL number small enough to run RKDG without adding a slope limiter. In this case the results of the FVKS and MUSCL/HLL were not greatly affected, whereas the RKDG gave a sharper profile for the velocity discontinuity and therefore a localization of the density in two cells. The final RKDG result is however not totally satisfactory because mass accretion in the mid-cell is not as good as FVKS. Further investigations are needed. One possible reason being that we are comparing a vertex-centered scheme with transported polynomial values with cell-centered schemes with reconstructed slopes, by projecting the RKDG result as a cell-centered one.

For the presented test cases, the RKDG and the FVKS are competitive with each other and overpass the MUSCL/HLL. The FVKS provides slightly better results than RKDG and we believe this is due to the exact update in time for the former.

## 5 Conclusion

The comparison of the different numerical schemes presented in this work is an essential step toward the ultimate goal of finding an accurate, realizable, cost effective and parameter-free numerical scheme on unstructured grids that can be applied to the KBMM hierarchy. The RKDG and the MUSCL/HLL were applied to higher order models such as the isotropic Gaussian and Anisotropic Gaussian closures KBMM while the application of the FVKS is limited to the Dirac distributions for the velocity. Therefore, even though it is remarkable that the FVKS usage is attractive for this model, we are interested in a numerical scheme applicable to all the KBMMs and the new class of RKDG is a promising choice. To ensure the monotonicity of the RKDG results the best method should be found to make sure that the local maximum principle is respected without degenerating the accuracy. The modified minmod solves this problem but introduces a parameter depending on the initial condition. Further inverstigation of the PGD problem generating the $\delta$-shock should be carried out to better understand the RKDG result deterioration. Finally, 2D test cases on HIT were already conducted for the MUSCL/HLL and RKDG schemes (using the minmod limiter) and RKDG was found to be competitive from a quality/cost point of view. This study should be extended to include FVKS results, knowing that it is expected to be more efficient for the MK model.

# References

1. Berthon, C.: Stability of the MUSCL schemes for the Euler equations. Commun. Math. Sci. **3**(2), 133–157 (2005)
2. Bouchut, F.: On zero pressure gas dynamics. In: Perthame, B. (ed.) Advances in Kinetic Theory and Computing, Series on Advances in Mathematics for Applied Sciences, pp. 171–190. World Scientific Publishing, River Edge, NJ (1994)
3. Bouchut, F., Jin, S., Li, X.: Numerical approximations of pressureless and isothermal gas dynamics. SIAM J. Numer. Anal. **41**, 135–158 (2003)
4. de Chaisemartin, S.: Eulerian models and numerical simulation of turbulent dispersion for polydisperse evaporation sprays. Ph.D. thesis, Ecole Centrale Paris, France (2009)
5. Cockburn, B., Shu, C.W.: TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws II: General framework. Math. Comput. **52**(186), 411–435 (1989)
6. Gottlieb, S., Shu, C.W., Tadmor, E.: Strong stability-preserving high-order time discretization methods. SIAM Rev. **43**(1), 89–112 (2001)
7. Larat, A., Massot, M., Vié, A.: A stable, robust and high order accurate numerical method for Eulerian simulation of spray and particle transport on unstructured meshes. In: Annual Research Briefs 2012, pp. 205–216. Center for Turbulence Research, Stanford University, USA (2012)
8. Laurent, F., Massot, M.: Multi-fluid modeling of laminar poly-dispersed spray flames: origin, assumptions and comparison of the sectional and sampling methods. Combust. Theor. Model. **5**, 537–572 (2001)
9. LeVeque, R.J.: Finite volume methods for hyperbolic problems. Cambridge texts in applied mathematics. Cambridge University Press, Cambridge (2002)
10. van Leer, B.: Towards the ultimate conservative difference scheme V. A second order sequel to Godunov's method. J. Comput. Phys. **32**(1), 101–136 (1979)
11. Vié, A., Doisneau, F., Massot, M.: On the Anisotropic Gaussian closure for the prediction of inertial-particle laden flows. http://hal.archives-ouvertes.fr/hal-00912319 (2013) (Submitted)
12. Williams, F.: Spray combustion and atomization. Phys. Fluids **1**, 541–545 (1958)
13. Yang, Y., Wei, D., Shu, C.W.: Discontinuous Galerkin method for Krause consensus models and pressureless Euler equations. J. Comput. Phys. **252**, 109–127 (2013)
14. Zhang, X., Xia, Y., Shu, C.W.: Maximum-principle-satisfying and positivity-preserving high order discontinuous Galerkin schemes for conservation laws on triangular meshes. J. Sci. Comput. **50**(1), 29–62 (2012)

# Shock Capturing for Discontinuous Galerkin Methods using Finite Volume Subcells

**Matthias Sonntag and Claus-Dieter Munz**

**Abstract** We present a shock capturing procedure for high order discontinuous Galerkin methods, by which shock regions are refined and treated by the finite volume techniques. Hence, our approach combines the good properties of the discontinuous Galerkin method in smooth parts of the flow with the perfect properties of a total variation diminishing finite volume method for resolving shocks without spurious oscillations. Due to the subcell approach the interior resolution on the discontinuous Galerkin grid cell is preserved and the number of degrees of freedom remains the same. In this paper we focus on an implementation of this coupled method and show our first results.

## 1 Introduction

Discontinuous Galerkin methods of high order accuracy have the problem that shock waves travelling through grid cells introduce instabilities. The high order polynomial in the coarse grid cell generates spurious oscillations when such an inner element jump has to resolved. There exist different methods to circumvent these problems. One is the use of explicit artificial viscosity, which adds locally viscosity to the original equations to smear the discontinuities in such a way that it can be resolved by the numerical approximation. This was originally proposed by von Neumann and Richtmyer [7] for finite difference schemes. Persson and Peraire [4] adapted this to high order discontinuous Galerkin (DG) methods to eliminate the high frequencies without widening the shock over a couple of cells. They applied this artifical viscosity approach also on subcells. Another technique to capture shocks in a DG framework,

M. Sonntag (✉) · C.-D. Munz
Institute of Aerodynamics and Gas Dynamics, University of Stuttgart, Pfaffenwaldring 21,
70569 Stuttgart, Germany
e-mail: sonntag@iag.uni-stuttgart.de

C.-D. Munz
e-mail: munz@iag.uni-stuttgart.de

which is also inspired by the finite volume methodology, is the approach of refining the grid in shock regions, while reducing the degree of the polynomials [1, 2], often called $hp$-adaption. In general the reduction of the polynomial degree decreases the oscillations, while the resolution has to be preserved by $h$-refinement.

In this paper we investigate the latter approach. But, we use an inherent refinement of the discontinuous Galerkin elements into several finite volume subcells with a lower order approximation without changing the degree of freedoms or the general data structure. The outline of this paper is as follows. First we summarize the basic concepts of our DG method and define the degrees of freedom (DOF) of an element. In Sect. 3 we then derive a finite volume method on subcells associated with one degree of freedom within the DG grid cell. The interior FV method is capable to capture strong shocks due to its total variation diminishing character. The following section shows a numerical example to illustrate the effectiveness of our approach to handle shocks.

## 2 The Discontinuous Galerkin Spectral Element Method

In this section we recapitulate shortly the basic ideas of the discontinuous Galerkin method with use of spectral elements as implemented in our CFD code FLEXI. For a detailed description we refer the reader to Hindelang et al. [3]. The general system of conservation laws is given as

$$\mathbf{u}_t(x) + \nabla \cdot F(\mathbf{u}(x)) = 0 \quad \forall x \in \Omega, \tag{1}$$

where $\mathbf{u}$ is the vector of conservative variables, $F$ the physical fluxes and $\Omega$ the computational domain, which is subdivided into hexahedral elements. Mapping each of this elements onto the reference element $E = [-1, 1]^2$ yields

$$J(\xi)\mathbf{u}_t(t, \xi) + \nabla_\xi \cdot \mathscr{F}(\mathbf{u}(t, \xi)) = 0. \tag{2}$$

where $J(\xi)$ is the Jacobian of the mapping, $\mathscr{F}$ are the transformed fluxes and $\xi = (\xi^1, \xi^2)^\top$ are the coordinates in the reference space. We multiply the transformed conservation law (2) with a test function $\Phi$ and integrate over the reference element $E$ to obtain, after partial integration of the second integral, the weak formulation

$$\int_E J\mathbf{u}_t \, d\xi + \int_{\partial E} (\mathscr{F} \cdot n) \, \Phi \, dS_\xi - \int_E \mathscr{F} \cdot \nabla_\xi \Phi \, d\xi = 0, \tag{3}$$

where $n$ is the normal vector of the reference element $E$. We approximate the solution in the reference element by a polynomial tensor product basis of degree $N$ in each space direction

$$\mathbf{u}(\xi) = \sum_{i,j=0}^N \hat{\mathbf{u}}_{ij} \psi_{ij}(\xi) \quad \text{with } \psi_{ij} = l_i(\xi^1) l_j(\xi^2), \tag{4}$$

**Fig. 1** DG reference element $E$ with Gauss points ● and locations of the boundary fluxes □ at the DG interface for $N = 3$ in 2D

where $\hat{\mathbf{u}}_{ij}$ are the nodal degrees of freedom and $l_i(\xi)$ are the one-dimensional Lagrange interpolation polynomials defined by the Gauss nodes $\{\xi_i\}_{i=0}^{N}$. In our spectral element approach we use some sort of collocation technique. The integration in the discontinuous Galerkin framework is approximated by Gauss quadrature based on the same Gauss points with the Gauss weights $\{\omega_i\}_{i=0}^{N}$. Furthermore the Galerkin method uses the same ansatz and test functions $\Phi = \psi_{ij}$.

In the following we concentrate on the boundary integral of (3), because neighboring DG elements are coupled only by this term. Since the solution may be discontinuous at the interfaces the state is given twice, by the left and by the right element. Therefore, the flux $\mathscr{F}$ is approximated by a Riemann solver including the state of the actual element $\mathbf{u}$ and the state of the adjacent element $\mathbf{u}^+$. By the tensor product ansatz the boundary fluxes are also interpolated in Gauss points $\mathbf{u}_{\pm 1,j}$ or $\mathbf{u}_{i,\pm 1}$ at the left, right, top and bottom of the grid cell, respectively. The states in the Gauss points are computed by 1D-extrapolation of the inner nodal DOFs $\hat{\mathbf{u}}_{ij}$ along the $\xi^1$ or $\xi^2$ direction, see Fig. 1. Together with the extrapolated states $\mathbf{u}_{\pm 1,j^+}^+$ and $\mathbf{u}_{i^+,\pm 1}^+$ of the respective adjacent element the Riemann solver then computes the approximated fluxes $f_{\pm 1,j}$ and $f_{i,\pm 1}$ at the boundary, which are than used to calculate the boundary integral of (3) by use of a Gauss quadrature. All fluxes are computed, of course, only once for each edge and then added to both neighboring elements. Therefore, in our implementation at the beginning of every time stage the two states at each edge of the mesh are extrapolated from both sides of the edge. Later on they are then inserted into the Riemann solver to compute the fluxes at the boundaries.

Remember that DG elements are only coupled by the fluxes at the faces. Therefore a parallelization exchanges data only over the faces laying at the MPI borders. A short summary of the main steps of our implementation including the parallelization with MPI reads as follows:

**Algorithm 1 (DG method)**

1. For each element: Extrapolate the state $\mathbf{u}$ to the faces.
2. For each MPI face: Send extrapolated boundary state from master to slave.
3. For each element: Compute volume fluxes and the volume integral of (3).

**Fig. 2** DG reference element splitted into FV subcells $- - -$ with Gauss points •, Gauss weights $\omega_i$ and locations of the inner ■ and the interface □ boundary fluxes

4. For each face (excluding MPI faces on the master): Evaluate the fluxes.
5. For each MPI face: Send flux at the face back from slave to master.
6. For each element: Calculate the boundary integral of (3).

All these steps together calculate the time derivative of the DOFs

$$\frac{\partial \hat{\mathbf{u}}_{ij}}{\partial t} = \text{step } 1 \ldots 6 \quad \forall i, j \; \forall \text{elements.} \tag{5}$$

We integrate this derivative with an explicit Runge Kutta method in time.

## 3 Shock Capturing with Finite Volume Subcells

Numerical schemes of high order accuracy often have difficulties resolving shocks without generating new extrema or oscillations in the solution. Often seen in the region of the shock is a reduction of the polynomial degree to handle the problem that polynomials of higher order can't resolve discontinuities without oscillation. To avoid a loss in resolution this is then combined with a local mesh refinement. In this section we present a natural way of shock capturing by constructing a refinement of the high order discontinuous Galerkin element into several internal finite volume elements without introducing new degrees of freedom. The fixed number of DOFs helps us to keep the method as simple as possible and to reach a high computational performance in the actual implementation.

For a discontinuous Galerkin element of polynomial order $N$ we use, as described in Sect. 2, $N+1$ Gauss points for interpolation and integration in each space direction. Each of the $(N + 1)^2$ Gauss points $\{x_{ij}\}_{i,j=0}^N$ of the DG reference element in 2D is used as a node of a finite volume subcell $\kappa_{ij}$. In the DG reference element $E$ the $ij$-th finite volume subcell has the size $\omega_i \times \omega_j$, where $\omega_i, \omega_j$ are the Gauss weights corresponding to the point $x_{ij}$, see Fig. 2.

We now formulate the finite volume method for the transformed conservation law (2) on the reference element of the discontinuous Galerkin method $E$. Each subcell $\kappa_{ij}$ of $E$ is now a control volume in the finite volume context and the corresponding equation reads after applying the divergence theorem to the second integral as

$$\int_{\kappa_{ij}} J\mathbf{u}_t \, \mathrm{d}\xi + \int_{\partial\kappa_{ij}} \mathscr{F}(\mathbf{u}) \cdot n \, \mathrm{d}S_\xi = 0. \tag{6}$$

Since we do not introduce new DOFs, every finite volume subcell $\kappa_{ij}$ contains exactly one Gauss point $x_{ij}$ of the discontinuous Galerkin discretization. We use this Gauss point as "center" of the respective finite volume subcell and therefore take the nodal value $\hat{\mathbf{u}}_{ij}$ of the DG discretization (4) as the approximative value in the FV subcell. With the volume $\omega_i\omega_j$ of the subcell $\kappa_{ij}$ the volume integral of (6) becomes

$$\int_{\kappa_{ij}} J\mathbf{u}_t \, \mathrm{d}\xi = \omega_i\omega_j J_{ij} \frac{\partial \hat{\mathbf{u}}_{ij}}{\partial t}, \tag{7}$$

where $J_{ij}$ is the average of the Jacobian $J$ in the $ij$-th subelement.

Because the state in a FV subcell is constant we can replace the boundary integral of (6) by the midpoint rule point as

$$\int_{\partial\kappa_{ij}} \mathscr{F}(\mathbf{u}) \cdot n \, \mathrm{d}S_\xi = \omega_j \left( f_{i-\frac{1}{2},j}(\mathbf{u}, \mathbf{u}^+, n) + f_{i+\frac{1}{2},j}(\mathbf{u}, \mathbf{u}^+, n) \right)$$
$$+ \omega_i \left( f_{i,j-\frac{1}{2}}(\mathbf{u}, \mathbf{u}^+, n) + f_{i,j+\frac{1}{2}}(\mathbf{u}, \mathbf{u}^+, n) \right), \tag{8}$$

where $f_{i-\frac{1}{2},j}(\mathbf{u}, \mathbf{u}^+, n)$ denotes the flux at the left edge of the $ij$-th subcell and $\omega_i$ and $\omega_j$ denote and the lengths of the edges. The numerical flux is computed by a Riemann solver involving the state $\mathbf{u}^+$ of the neighboring subcells. In total the finite volume method for the $ij$-th subcell reads after division by the volume $\omega_i\omega_j$ and the Jacobian $J_{ij}$ as

$$\frac{\partial \hat{\mathbf{u}}_{ij}}{\partial t} = -\frac{1}{J_{ij}\omega_i} \left( f_{i-\frac{1}{2},j}(\mathbf{u}, \mathbf{u}^+, n) + f_{i+\frac{1}{2},j}(\mathbf{u}, \mathbf{u}^+, n) \right)$$
$$- \frac{1}{J_{ij}\omega_j} \left( f_{i,j-\frac{1}{2}}(\mathbf{u}, \mathbf{u}^+, n) + f_{i,j+\frac{1}{2}}(\mathbf{u}, \mathbf{u}^+, n) \right). \tag{9}$$

Therewith we have another expression than (5) for the time derivative of the DOFs in one DG cell, which can be directly interchanged within a time step of the explicit time integration. Since we use the same nodal DOFs for the DG and the FV method, the approximation can be interpreted either as DG polynomial or as set of FV values in every stage of the Runge Kutta method.

Comparing Fig. 1 with Fig. 2 it is clear that the fluxes $f_{\pm 1,j}$ and $f_{i,\pm 1}$ at the faces of the DG element are calculated in the same points as the fluxes $f_{0-\frac{1}{2},j}$, $f_{N+\frac{1}{2},j}$,

$f_{i,0-\frac{1}{2}}$ and $f_{i,N+\frac{1}{2}}$ of the FV subcells lying directly at the DG boundary (white squares). In both cases the boundary fluxes are evaluated by a Riemann solver using the master and the slave state at this face, irrespective of wether they are extrapolated from a DG element or a FV subcell.

Computing the inner fluxes of the FV subcells (gray squares in Fig. 2) is comparable to the volume integral of the DG method as an absolutely local operation, not involving any data from neighboring elements. We modify the DG method (Algorithm 1) to a coupled DG / FV subcell algorithm, where the new steps are italic printed, as follows

**Algorithm 2  (Modified steps of coupled DG / FV-subelement Method)**

0.  *For each element: Indicator based switching between DG and FV.*
3a. For each DG element: Compute volume fluxes and the volume integral of (3).
3b. *For each FV element: Compute inner fluxes (gray squares) and divide by weights $\omega_i$ or $\omega_j$ and Jacobian $J_{ij}$ (cf. equation (9)).*
6a. For each DG element: Calculate the boundary integral of (3).
6b. *For each FV element: Divide fluxes at the DG boundaries (white squares) by weight $\omega_i$ or $\omega_j$ and Jacobian $J_{ij}$ (cf. equation (9)).*

As we see there are only three main differences to the previously stated DG algorithm. The first one is the change of the volume integral into multiple inner surface integrals evaluated with Riemann solvers (step 3). Secondly the FV subcell method uses the same fluxes at the outer boundaries as the DG method (step 6). A difference is that in the DG method a boundary flux affects all DOFs in that considered cell, while in the FV subcell method the subcell adjacent to the DG boundary is influenced only. The third change is hidden in step 1. The extrapolation to the DG faces must be modified for FV subcells. Of course, the algorithm must be also extended by an indicator (step 0), which decides where to use DG or FV.

*Remark 1* (*Block unstructured FV method*) This DG method with the finite volume subcell framework in every DG element may be interpreted also as blockwise finite volume method on unstructured curved blocks It may be also interpreted as a heterogeneous domain decomposition appraoch with a weak coupling of the sub-domains. Since the DOFs of both algorithms are the same, there is no difference in data formats and a FV subcell solution can be directly compared to the DG solution. This gives us also the ability to investigate the advantages and disadvantages of the discontinuous Galerkin method by comparison with the finite volume method on the same mesh. Another useful application of the blockwise FV subcell method is also the stabilization at the beginning of a computation when initializing with freestream.

*Remark 2* (*Higher order reconstruction*) In this paper we only consider the refinement of DG elements with FV subcells without using a higher order reconstruction within the finite volume method to keep the derivation of the general method as simple as possible. Of course our implementation includes a FV reconstruction coupled with different types of slope limiters.

**Fig. 3** Density of DMR on a grid with $480 \times 120$ DG elements for $N = 4$ at $t = 0.2$



**Fig. 4** DMR: Refined cells (*black*) and density in DG cells

## 4 Numerical Examples

The double mach reflection sends a Mach 10 shock diagonally into a reflecting wall and was originially introduced by Woodward and Colella [8]. This problem has been widely used as a test case for high resolution schemes in the literature. With the MPI parallelized version of our code we computed this example on 16 cores until the time $t = 0.2$ with a polynomial degree of $N = 4$. In Fig. 3 one can see 30 equally spaced contour lines from $\rho = 1.5$ to $\rho = 22.9705$; Fig. 4 shows the refined regions. For this results we used a local Lax-Friedrichs Riemann solver and in the finite volume subcells a second order reconstruction with the Sweby slope limiter ($\beta = 1.4$) [6]. The indicator which switches between DG and FV subcells was chosen as the famous Persson indicator [4]. The shown results are in good agreement to the results of other groups, for example Shi et al. [5].

The second example is the forward facing step, also described by Woodward and Colella [8]. In Fig. 5 the density at time $t = 4.0$ is plotted. This example was also computed on 16 cores by using the same Riemann solver and slope limiting as above. In this case the polynomial degree of the DG solution is $N = 6$ and the grid is equidistant with $h = 1/50$, in total 6300 DG cells. A detail view of the FFS briefly compares, in Fig. 6, our method with a full finite volume scheme.

**Fig. 5** Grid and density of forward facing step at $t = 4.0$ with $N = 6$. Grid lines of FV subcells let refined cells look like *black squares*



**Fig. 6** Comparison of hybrid DG/FV subcell method (*left*) against full finite volume scheme

## 5 Conclusion

We have presented a shock-capturing strategy for discontinuous Galerkin schemes, which uses a natural subcell decomposition and a total variation diminishing finite volume method on the subcells. This procedure preserves the whole data structure of the underlying DG scheme and can be used in an adaptive way in grid cells by a simple switch. Our discontinuous Galerkin scheme was based on spectral elements and used the same nodal DOFs for both numerical schemes. This approach may be considered as a combination of a DG scheme with a finite volume scheme on an $h$-refined grid. In smooth parts of the flow large grid cells are used and high order of accuracy, which is very efficient on massively parallel systems, while in troubled cells with strong gradients we switch to a total variation diminishing finite volume solver on subcells. In this sense the DG approach may be considered as a general framework of a heterogeneous domain decomposition.

## References

1. Baumann, C.E., Oden, J.T.: A discontinuous hp finite element method for the euler and navier-stokes equations. Int. J. Numer. Meth. Fluids **31**, 7995 (1999)
2. Burbeau, A., Sagaut, P., Bruneau, C.H.: A problem-independent limiter for high-order rungekutta discontinuous galerkin methods. J. Comput. Phys. **169**(1), 111–150 (2001)

3. Hindenlang, F., Gassner, G., Altmann, C., Beck, A., Staudenmaier, M., Munz, C.D.: Explicit discontinuous galerkin methods for unsteady problems. Comput. Fluids **61**, 86–93 (2012)
4. Persson, P.O., Peraire, J.: Sub-cell shock capturing for discontinuous galerkin methods. In: Proceedings of the 44th AIAA Aerospace Sciences Meeting and Exhibit. American Institute of Aeronautics and Astronautics (2006)
5. Shi, J., Zhang, Y.T., Shu, C.W.: Resolution of high order weno schemes for complicated flow structures. J. Comput. Phys. **186**(2), 690–696 (2003)
6. Sweby, P.K.: High resolution schemes using flux limiters for hyperbolic conservation laws. SIAM J. Numer. Anal. **21**(5), 995–1011 (1984)
7. VonNeumann, J., Richtmyer, R.D.: A method for the numerical calculation of hydrodynamic shocks. J. Appl. Phys. **21**(3), 232–237 (1950)
8. Woodward, P., Colella, P.: The numerical simulation of two-dimensional fluid flow with strong shocks. J. Comput. Phys. **54**(1), 115–173 (1984)

# A Simple Well-Balanced, Non-negative and Entropy-Satisfying Finite Volume Scheme for the Shallow-Water System

**Emmanuel Audusse, Christophe Chalons and Philippe Ung**

**Abstract** This work considers the numerical approximation of the shallow-water equations. In this context, one faces three important issues related to the well-balanced, non-negativity and entropy-preserving properties, as well as the ability to consider vacuum states. We propose a Godunov-type method based on the design of a three-wave Approximate Riemann Solver (ARS) which satisfies all these properties together.

## 1 Introduction

In this work, we look for a numerical scheme for the shallow-water equations given by:

$$\begin{cases} \partial_t h + \partial_x \left( hu \right) = 0, \\ \partial_t \left( hu \right) + \partial_x \left( hu^2 + \dfrac{gh^2}{2} \right) = -gh\partial_x b(x), \end{cases} \tag{1}$$

where $b(x)$ is a sufficiently smooth topography, $g$ refers to the gravitational acceleration, and the water height $h$ and the velocity $u$ depend on time $t$ and space $x$; $h$ and $u$ are the primitive variables and $b$ is given. In addition, the associated entropy inequality is written as:

E. Audusse (✉)
University Paris 13, 99 Avenue Jean Baptiste Clément, 93430 Villetaneuse, France
e-mail: audusse@math.univ-paris13.fr

C. Chalons
University of Versailles-Saint-Quentin-en-Yvelines, 55 Avenue de Paris, 78000 Versailles, France
e-mail: christophe.chalons@uvsq.fr

P. Ung
University of Orléans, 6 Avenue du Parc Floral, 45100 Orléans, France
e-mail: ung@math.cnrs.fr

$$\begin{cases} \partial_t \mathscr{U}(w) + \partial_x \mathscr{F}(w) \leqslant -ghu\partial_x b, \\ \mathscr{U}(w) = \dfrac{hu^2}{2} + \dfrac{gh^2}{2}, \quad \mathscr{F}(w) = \left(\dfrac{u^2}{2} + gh\right)hu, \end{cases} \qquad (2)$$

where $w = (h, hu)^T \in \mathbb{R}^+ \times \mathbb{R}$. The scheme should preserve three important properties that are satisfied by the exact solution of the shallow-water equations: the non-negativity of water heights, a discrete entropy inequality, and the steady states of the lake at rest defined by

$$h_L + b_L = h_R + b_R, \quad \text{and} \quad u_L = u_R = 0, \qquad (3)$$

where the indices $L$ and $R$ refer to the left and right states in the Riemann problem detailed later. Furthermore, it should be able to handle vacuum, in particular, the steady state of the wet-dry transition

$$h_L + b_L \leqslant b_R, \quad h_R = 0, \quad \text{and} \quad u_L = u_R = 0, \qquad (4)$$

There is a huge amount of work about this topic but most of the schemes fail to satisfy these three properties at once. Up to our knowledge, four methods [2, 4, 5, 10] are proved to fulfill the three requirements but they are costly in terms of computing runtime and/or based on quite complex algorithms. In this work, we propose a numerical scheme adapted to vacuum that endows the three properties and that is very cheap and simple to implement. Numerical experiments are proposed to compare the new method with some popular non-negative and well-balanced schemes for which no fully discrete entropy property is proved [1, 3, 7, 8].

## 2 Numerical Scheme

In the following, we describe a Godunov-type finite volume scheme for (1) and (2). Let us first introduce some notations. We consider a sequence of points $x_{i+1/2}$ such that

$$x_{i-1/2} < x_{i+1/2}, \quad \forall i \in \mathbb{Z}$$

and we define the cells $C_i$ and space steps $\Delta x_i = \Delta x$, such that

$$C_i = ]x_{i-1/2}, x_{i+1/2}[, \quad \Delta x = x_{i+1/2} - x_{i-1/2}.$$

In addition, we set $x_i = (x_{i-1/2} + x_{i+1/2})/2$.

We also introduce a time step $\Delta t > 0$ that allows to define a sequence of intermediate times $t^n$ by

$$t^{n+1} = t^n + \Delta t.$$

Starting from a given piecewise constant approximate solution at time $t^n$, we construct the solution at time $t^{n+1}$ in two steps:

- we build an approximate solution of the Riemann problem at each interface $x_{i+1/2}$,
- we obtain the new solution by calculating the average value of the juxtaposition of these solutions in each cell $C_i$ at time $t^{n+1}$.

As an approximate Riemann solution associated with initial data

$$(w(x, 0), \tilde{b}(x)) = \begin{cases} (w_L, b_L) & x < 0, \\ (w_R, b_R) & x > 0, \end{cases} \tag{5}$$

we consider a simple approximate Riemann solver composed by three waves propagating with velocities $\lambda_L$, $\lambda_0 = 0$ and $\lambda_R$. Note that the most simple approximate Riemann solvers contain two waves (as the well-known HLL flux [4, 9] but are not able to preserve steady states. The choice of a three waves solver is then a compromise between simplicity and accuracy that was also adopted in [2, 7, 8]. Note also that the quantities $b_L$ and $b_R$ have to be related to the given bottom topography $b(x)$ to ensure the consistency with the source term in (1)

$$b_L = \frac{1}{\Delta x} \int_{-\Delta x}^{0} b(x)dx, \quad b_R = \frac{1}{\Delta x} \int_{0}^{\Delta x} b(x)dx,$$

From [7–9], it is known that such an approximate Riemann solver is consistent in the integral sense with (1) provided that the intermediate states satisfy the following consistency relations:

$$f(w_R) - f(w_L) - \Delta x \, s \, (\Delta x; \, w_L, \, w_R, \, b_L, \, b_R) = \lambda_L(w_L^* - w_L) + \lambda_R(w_R - w_R^*), \tag{6}$$

with $f(w) = (f^h(w), f^q(w))^T = (hu, \, hu^2 + gh^2/2)^T$ and $s(\Delta x; w_L, w_R, b_L, b_R)$ is defined as an approximation of the source term in (1), since it has to satisfy:

$$\lim_{\substack{w_L, w_R \to w \\ \Delta x \to 0}} s \, (\Delta x; \, w_L, \, w_R, \, b_L, \, b_R) = \begin{pmatrix} 0 \\ -gh\partial_x b \end{pmatrix}. \tag{7}$$

Recall also that the scheme satisfies a discrete version of the entropy inequality (2) provided that the following conditions on the intermediate states is fulfilled

$$\mathscr{F}(w_R) - \mathscr{F}(w_L) - \Delta x \sigma_v(\Delta x; \, w_L, \, w_R, \, b_L, \, b_R)$$
$$\leqslant \lambda_L(\mathscr{U}(w_L^*) - \mathscr{U}(w_L)) + \lambda_R(\mathscr{U}(w_R) - \mathscr{U}(w_R^*)), \tag{8}$$

where $\sigma_v \, (\Delta x; \, w_L, \, w_R, \, b_L, \, b_R)$ has to be defined such that

$$\lim_{\substack{w_L, w_R \to w \\ \Delta x \to 0}} \sigma_v \left( \Delta x; w_L, w_R, b_L, b_R \right) = -ghu\partial_x b. \tag{9}$$

We refer again to [7–9] for more details and we recall that properties (7) and (9) are in accordance with the fact that we consider only smooth topographies.

## 2.1 Expression of the Solution in the Intermediate States

We propose to define the two intermediate states by imposing the consistency relations in the integral sense resulting from the Eq. (1), see (6):

$$\begin{cases} h_R u_R - h_L u_L = \lambda_L \left( h_L^* - h_L \right) + \lambda_R \left( h_R - h_R^* \right), \\ \left( h_R u_R^2 + \dfrac{g h_R^2}{2} \right) - \left( h_L u_L^2 + \dfrac{g h_L^2}{2} \right) + g \Delta x \{h\partial_x b\} \\ \qquad = \lambda_L \left( h_L^* u_L^* - h_L u_L \right) + \lambda_R \left( h_R u_R - h_R^* u_R^* \right). \end{cases} \tag{10}$$

The definition of the approximation of the source term $\{h\partial_x b\}$ will be related to the well-balanced property and is given hereafter.

In order to close this system, two relations are missing and we suggest to impose two relations across the standing waves

$$\begin{cases} h_L^* + b_L = h_R^* + b_R, \\ h_L^* u_L^* = h_R^* u_R^*, \end{cases} \tag{11}$$

which are consistent with the steady states of the system (1).

By solving the Eqs. (10) and (11), we define the water heights in the intermediate states

$$h_L^* = h_{HLL} + \frac{\lambda_R}{\lambda_R - \lambda_L} \Delta b \quad \text{and} \quad h_R^* = h_{HLL} + \frac{\lambda_L}{\lambda_R - \lambda_L} \Delta b, \tag{12}$$

where $\Delta b = b_R - b_L$ and

$$h_{HLL} = \frac{\lambda_R h_R - \lambda_L h_L}{\lambda_R - \lambda_L} - \frac{1}{\lambda_R - \lambda_L} \left( h_R u_R - h_L u_L \right) \tag{13}$$

is the intermediate water height associated to the HLL solver [4, 9].

Then, from the Eqs. (10) and (11), we deduce the intermediate discharge $q^* := h_L^* u_L^* = h_R^* u_R^*$,

$$q^* = q_{HLL} - \frac{g}{\lambda_R - \lambda_L} \Delta x \{h\partial_x b\}, \tag{14}$$

with

$$q_{HLL} = \frac{\lambda_R h_R u_R - \lambda_L h_L u_L}{\lambda_R - \lambda_L} - \frac{\left(h_R u_R^2 + \frac{g h_R^2}{2}\right) - \left(h_L u_L^2 + \frac{g h_L^2}{2}\right)}{\lambda_R - \lambda_L},$$

the intermediate discharge involved in the HLL scheme [4, 9]. From now on $\lambda_L$ and $\lambda_R$ refer to the values that ensure the stability properties of the classical HLL scheme (non-negativity of the water heights but also validity of an entropy inequality as we will need later on). In practice, we apply the following expression from [4]

$$\lambda_L = \min_{w=w_L,w_R} (u - \sqrt{gh}, 0) \quad \text{and} \quad \lambda_R = \max_{w=w_L,w_R} (u + \sqrt{gh}, 0). \quad (15)$$

## 2.2 Properties of the Scheme: Non-negativity, Well-Balancing and Entropy

We first study the non-negativity of the water height. In regard to the expression of the intermediate water heights (12), it is not possible to ensure the non-negativity of these quantities. That is why, we suggest to modify these intermediate values depending on the sign of $\Delta b$. In the case $\Delta b \geqslant 0$, we clearly have

$$h_R^* \leqslant h_{HLL} \leqslant h_L^*.$$

In order to ensure the non-negativity of $h_R^*$, we introduce a cut in its definition and we modify the definition of $h_L^*$ in order to still fulfill the consistency relation (10)

$$\tilde{h}_R^* = \max\left(h_R^*, 0\right) \quad \text{and} \quad \tilde{h}_L^* = h_L^* - \frac{\lambda_R}{\lambda_L}\left(h_R^* - \tilde{h}_R^*\right). \quad (16)$$

Note that the expression of the discharge $q^*$ is unchanged and then consistency relation (10) remains also valid.

In practice and in order to avoid threshold values near 0 for the wave velocities, we will exclusively work with the quantities $\lambda_R \tilde{h}_R^*$ and $\lambda_L \tilde{h}_L^*$,

$$\lambda_R \tilde{h}_R^* = \max\left(\lambda_R h_R^*, 0\right) \quad \text{and} \quad \lambda_L \tilde{h}_L^* = \lambda_L h_L^* - \lambda_R\left(h_R^* - \tilde{h}_R^*\right). \quad (17)$$

The case $\Delta b < 0$ can be treated by applying the same method.

We now turn to the well-balancing property. Preserving the steady states of a lake at rest in the context of the proposed simple Riemann solver is an easy task. We introduce the following natural discretization of the source term

$$\{h\partial_x b\} = \frac{h_L + h_R}{2\Delta x}\,\Delta b. \tag{18}$$

Simple computations show that this discretization preserves the steady state of the lake at rest (3).

Note that we can also preserve the lake at rest in the case of a wet-dry transition (4) with a slight modification of the source term.

We now turn to the study of the entropy property. The scheme is entropy preserving if it satisfies the discrete entropy inequality (8). Inspired by [2], we will in fact prove a variant of this entropy inequality and prove that there exists some term $\varepsilon(\Delta x)$ with property

$$\lim_{\Delta x \to 0} \varepsilon(\Delta x) = 0. \tag{19}$$

such that the following relation holds

$$\mathscr{F}(w_R) - \mathscr{F}(w_L) - \Delta x \sigma_v(\Delta x; w_L, w_R, b_L, b_R) + \Delta x\,\varepsilon(\Delta x)$$
$$\leqslant \lambda_L(\mathscr{U}(w_L^*) - \mathscr{U}(w_L)) + \lambda_R(\mathscr{U}(w_R) - \mathscr{U}(w_R^*)) \tag{20}$$

Indeed this correction term does not affect the validity of the Lax-Wendroff theorem, see [2]. Note that this correction term is not related to a modification of the proposed scheme but to the fact that, with the proposed scheme, we can not prove the classical entropy inequality (8).

We do not have the place to detail the proof of entropy inequality (20) here. We just mention that it is based on the fact that our scheme can be seen as a modification of the HLL solver for which a similar entropy inequality is known to be valid. Starting from the relation satisfied by the HLL solver, some calculation lead to the desired property with

$$\sigma_v(\Delta x; w_L, w_R, b_L, b_R) = -g\bar{h}\frac{q_{HLL}}{h_{HLL}}\frac{\Delta b}{\Delta x}, \tag{21}$$

which is clearly consistent with (2). We do not give the precise form of the error term $\varepsilon(\Delta x)$ but we mention that it involves the jump of bottom topography $\Delta b$. It follows that relation (19) is fulfilled only if the topography is at least continuous.

Finally we insist on the fact that, unlike [2], the main idea of the proposed scheme focuses on the proofs of the non-negativity and entropy-preserving properties which are obtained for $\lambda_L$ and $\lambda_R$ defined exactly as in the HLL scheme, and not defined asymptotically large according to specific behaviors like $-\lambda_L/\lambda_R \gg 1$ or $-\lambda_R/\lambda_L \gg 1$.

**Fig. 1** Fluvial flow: Comparison of orders of error for the water height $h$ (*left*) and the discharge $q$ (*right*) for different schemes

## 3 Numerical Results

We are interested in the behaviour of our scheme for fluvial flow regime and the way it handles the wet-dry transition. In this aim, we propose the well-known test case of a flow over a bump in the fluvial regime and the Thacker test case for the wet-dry transition.

In the following, we compare the L1-errors committed by the present scheme with the results obtained by using the HLL flux with different adaptation to the source term such as a centered discretization, the hydrostatic reconstruction [1] and the hydrostatic upwind scheme [3], together with the scheme proposed by Gallice [7, 8].

For the two test cases, the number of points goes from 100 to 1,600.

In the fluvial flow test case, the steady states are governed by the following equations

$$hu = K_1 \quad \text{and} \quad \frac{u^2}{2} + g(h + b) = K_2. \tag{22}$$

where $K_1$ and $K_2$ are two constants. Here we set $K_1 = 1$ and $K_2 = 25$. The domain is reduced to the interval $[-2; 2]$ and the bottom topography is defined by

$$\begin{cases} b(x) = \dfrac{\cos(10\pi(x+1)) + 1}{4}, & \text{if} -0.1 < x \leqslant 0.1, \\ 0, & \text{elsewhere.} \end{cases}$$

The initial datas correspond to this steady state.

The error curves (Fig. 1) emphasize the accuracy of the proposed scheme—called *simple solver* on the figures. Indeed, it gives a better approximation of the exact solution than other existing schemes with a gain of several orders.

**Fig. 2** Thacker test case: Comparison of orders of error for the water height $h$ (*left*) and the discharge $q$ (*right*) for different schemes at time $T = 16$

The Thacker test case brings out the ability of the scheme to handle vacuum, especially in the case of a wet-dry and dry-wet transition. The details of this test case are given in [6]. We can wisely precise that in this test case, the discharges are very low which can explain the large values of the relative errors one can observe for coarse mesh (Fig. 2).

## 4 Conclusion

In this paper, we have proposed what is up to our knowledge, the first simple to implement, non-negative, entropic and well-balanced scheme for the shallow water equations. The scheme proved to be very accurate on several typical test cases. The very motivation of this work comes from the numerical approximation of the solutions of the Saint-Venant–Exner equations for the problem of sediment bedload transport, to which we would like to adapt the proposed scheme therein. This is the matter of a work currently in progress.

## References

1. Audusse, E., Bouchut, F., Bristeau, M.O., Klein, R., Perthame, B.: A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. J. Sci. Comp. **25**, 2050–2065 (2004)
2. Berthon, C., Chalons, C.: A fully well-balanced, positive and entropy-satisfying Godunov-type method for the shallow-water equations. hal-00956799 (2014). http://hal.archives-ouvertes.fr/hal-00956799
3. Berthon, C., Foucher, F.: Efficient well-balanced hydrostatic upwind schemes for shallow-water equations. J. Comput. Phys. **231**, 4993–5015 (2012)

4. Bouchut, F.: Nonlinear stability of finite volume methods for hyperbolic conservation laws and well-balanced schemes for sources. Birkhäuser Verlag (2004)
5. Chinnayya, A., Leroux, A.Y., Seguin, N.: A well-balanced numerical scheme for the approximation of the shallow-water equations with topography: the resonance phenomenon. Int. J. Finite **1**, 1–33 (2004)
6. Delestre, O.: Simulation du ruissellement d'eau de pluie sur des surfaces agricoles. Ph.D. thesis, Université d'Orléans (2010)
7. Gallice, G.: Solveurs simples positifs et entropiques pour les systèmes hyperboliques avec terme source. C. R. Math. Acad. Sci. Paris **334**(8), 713–716 (2002)
8. Gallice, G.: Positive and entropy stable godunov-type schemes for gas dynamics and mhd equations in lagrangian or eulerian coordinates. Num. Math. **94**(4), 673–713 (2003)
9. Harten, A., Lax, P.D., van Leer, B.: On upstream differencing and godunov-type schemes for hyperbolic conservation laws. SIAM Rev. **25**(1), 53–61 (1983)
10. Perthame, B., Simeoni, C.: A kinetic scheme for the saint-venant system with source term. Calcolo **38**, 201–231 (2001)

# Well-Balanced Inundation Modeling for Shallow-Water Flows with Discontinuous Galerkin Schemes

**Stefan Vater and Jörn Behrens**

**Abstract** Modeling coastal inundation for tsunami and storm surge hazard mitigation is an important application of geoscientific numerical modeling. While the complex topography demands for robust and locally accurate schemes, computational parallel efficiency and discrete conservation properties of the scheme are required. In order to meet these requirements, Runge-Kutta discontinuous Galerkin numerical methods are attractive. However, maintaining conservation and well-balancedness of these schemes with wetting/drying boundary conditions poses a challenge. We address this issue by a local nondestructive modification of the flux computation at boundary cells, which maintains accuracy, conservation and well-balancedness. The development can be viewed as a specialized flux limiter, which proves its usefulness with three different test cases for inundation simulation.

## 1 Introduction

Coastal ocean modeling becomes a more and more important field in geoscientific research, as recent natural disasters like the 2004 Indian Ocean Tsunami, the 2011 Japan Tsunami, or the 2013 Super-Taifun Haiyan hitting the Philippines demonstrated. Therefore, planning for hazard mitigation and achieving early warning capacity heavily relies on coastal modeling [3]. The requirements for such models pose challenging demands on the numerical schemes used therein: while high computational efficiency is required, it is paramount that the schemes are robust to spatially and temporally changing irregular boundary conditions and domain shapes

S. Vater (✉) · J. Behrens
Center for Earth System Research and Sustainability (CEN), Universität Hamburg,
Hamburg, Germany
e-mail: stefan.vater@zmaw.de

J. Behrens
e-mail: joern.behrens@zmaw.de

(inundation). Additionally, since the equations used in such geoscientific applications represent sensitive geophysical balances, the numerical schemes must adhere to the conservation of quantities and the structure of the given problem.

In this study we propose a well-balanced inundation scheme for a Runge-Kutta discontinuous Galerkin (DG) scheme solving the 1D shallow water equations with bathymetry. This serves as a test bed for converging robust and conservative inundation schemes in more complex near realistic models (see e.g. [11]). The system of equations is given by

$$\mathbf{U}_t + \mathbf{F}(\mathbf{U})_x = \mathbf{S}(\mathbf{U}) , \tag{1}$$

where the vector of unknowns is given by $\mathbf{U} = (h, hu)^T$. The quantity $h = h(x, t)$ denotes the water height of a uniform density water layer and $u = u(x, t)$ is the particle velocity. The flux function is defined by $\mathbf{F}(\mathbf{U}) = (hu, hu^2 + \frac{g}{2}h^2)^T$, where $g$ is the gravitational constant. Furthermore, the bathymetry or bottom topography $b = b(x)$ is represented by the source term $\mathbf{S}(\mathbf{U}) = (0, -ghb_x)^T$.

The DG method is an attractive numerical scheme for wave and fluid dynamics modeling due to its discrete conservation property, its potential high order of accuracy, its parallel scalability, and not least its geometrical flexibility. However, while in finite volume methods inundation modeling has reached some kind of maturity with the introduction of hydrostatic reconstruction [1], it is still an ongoing research topic for DG methods [4]. A common approach is to enforce a minimum water level globally and define a threshold for the water elevation $h_{\min}$ below which a cell is declared dry [5, 7]. Another approach is based on modifying the bathymetry to achieve a globally wet simulation [13]. Our approach strives to maintain the bathymetry and only locally corrects the sea surface height $h$ in order to avoid unphysical behavior. Special care needs to be taken in order to maintain conservation and well-balancedness of the scheme.

## 2  Runge-Kutta Discontinuous Galerkin Method

We briefly introduce the numerical scheme with a focus on the wetting and drying treatment. For a more complete presentations of Runge-Kutta DG methods the reader is referred to e.g. [9, 12].

The governing equations are solved on the one-dimensional domain $[x_{\min}, x_{\max}]$, which is divided into intervals (cells) $I_i = (x_{i-1/2}, x_{i+1/2})$. On each interval, the Eq. (1) are multiplied by a test function $\varphi$ and integrated. Integration by parts of the flux term leads to the weak DG formulation

$$\int_{I_i} \mathbf{U}_t \varphi \, \mathrm{d}x - \int_{I_i} \varphi_x \mathbf{F}(\mathbf{U}) \, \mathrm{d}x + \left[ \mathbf{F}^*(\mathbf{U}) \varphi \right]_{x_{i-1/2}}^{x_{i+1/2}} = \int_{I_i} \mathbf{S}(\mathbf{U}) \varphi \, \mathrm{d}x .$$

A second integration by parts on the inner part of the interval leads to

$$\int_{I_i} \mathbf{U}_t \varphi \, dx + \int_{I_i} \mathbf{F}(\mathbf{U})_x \varphi \, dx + \left[ \left( \mathbf{F}^*(\mathbf{U}) - \mathbf{F}(\mathbf{U}) \right) \varphi \right]_{x_{i-1/2}}^{x_{i+1/2}} = \int_{I_i} \mathbf{S}(\mathbf{U}) \varphi \, dx , \quad (2)$$

which is the so-called strong DG formulation and will be used throughout this paper. Note that the interface flux $\mathbf{F}^*$ is not defined in general, since the solution can have different values in the adjacent cells. This problem is circumvented in the discretization by using the (approximate) solution of the corresponding Riemann problem.

System (2) is discretized using a semi-discretization in space with a piecewise polynomial ansatz for the discrete solution components and test functions. To obtain second-order accuracy, we use piecewise linear functions, which are represented by nodal Lagrange basis functions [12]. In view of a two-dimensional extension of the scheme, $n$-point Gauß-Legendre quadrature is applied to obtain an (exact) discretization of the integral terms. The remaining system of ordinary differential equations is then solved using Heun's method, which is the second-order representative of a standard Runge-Kutta total-variation diminishing (TVD) scheme [10, 16]. Well-balancing is achieved by using the same discretization for the inner cell pressure flux term $ghh_x$ and source $ghb_x$ in the momentum equation. At the interfaces no problems occur, since a continuous representation is used for the bottom topography. For stabilization and to avoid spurious oscillations near discontinuities a minmod slope limiter is applied after each internal Runge-Kutta stage [8, 15] in the hydrostatic variables $(h + b, hu)$.

## 2.1 Inundation Algorithm

In the coastal zone, where the water inundates and recedes due to wave and tidal activity, cells repeatedly become wet and dry. At the wet/dry interface, semi-dry cells occur, which have to be approximated by piecewise smooth functions in a DG discretization (see Fig. 1). This may introduce an artificial height gradient which destroys the well-balancedness of the scheme.

In the present scheme, such semi-dry cells are further distinguished into two subclasses (cf. [2, 5]). In the first case, the highest surface elevation within the cell is attained in a wet node. This is the so-called dambreak-type, in which the element may undergo a rapid wetting from above, and such cells are treated in the same manner as fully wet cells. On the other hand, flooding-type cells, where the highest surface elevation is attained in a dry node, are treated specially. Here, the (possibly unphysical) surface elevation gradient, which enters the momentum balance, must be neglected to ensure well-balancedness for nearly still-water states. Furthermore, all tendencies, which occur in an originally dry node, are redistributed to the wet node in order to obtain mass-conservation. In the current implementation a node is considered dry if it has a fluid height smaller than $10^{-8}$. Since slope limiting might conflict with the treatment in a flooding-type semi-dry cell, limiting is disabled in such cells.

**Fig. 1** Discretization of a semi-dry cell using the discontinuous Galerkin scheme with piecewise linear elements. Note the artificial height gradient introduced by the discretization

## 3 Numerical Results

We present three different test cases, which show the behavior of our scheme with respect to wetting and drying and its well-balanced property. The first test case is the classical "lake at rest" with an island in the middle of the domain. Additionally, the model is validated for two transient test cases with wetting and drying, for which the analytical solution is known. In all simulations the gravitational constant is set to $g = 9.81$. Here and below we omit the dimensions of the physical quantities, which should be thought in the standard SI system with $m$ (meter), $s$ (seconds) etc. as basic units. The discrete initial conditions and the bottom topography are derived from the analytical ones by interpolation at the nodal (cell interface) points.

### 3.1 Lake at Rest

This test is usually conducted to illustrate the effectiveness of the discrete balance between the flux term and the source term due to bottom topography in the momentum equation. However, it is not only crucial to maintain the exact balance, but also to show that small perturbations do not lead to an unphysical behavior of the numerical scheme. The test is conducted on the domain [0, 1] with periodic boundary conditions. Given $r = |x - 0.5|$, the bottom topography is set to

$$b(x) = \tilde{b}(r) = \begin{cases} a \cdot \frac{\exp(-0.5/(r_m^2 - r^2))}{\exp(-0.5/r_m^2)} & \text{if } r < r_m, \\ 0 & \text{otherwise,} \end{cases}$$

where the parameters are given by $a = 1.2$ and $r_m = 0.4$. The initial height is set to $h(x, 0) = \min(0, 1 - b(x))$, such that the bathymetry forms an island in the middle of the domain. In a first setup the initial momentum is set to $(hu)(x, 0) \equiv 0$. In a second simulation it is perturbed by a random disturbance of the order $10^{-8}$. The domain is discretized into 50 cells, and the timestep is set to 0.002. This corresponds to a CFL number of 0.3. The solution is integrated over 10000 timesteps until $t_{\max} = 20$.

**Fig. 2** Deviation of the surface elevation (*left*) and momentum (*right*) over time for the Lake at rest test case measured in the $L^\infty$ norm. Test initialized with zero momentum field (*top*) and small deviations in the momentum field of the order $10^{-8}$ (*bottom*)

In Fig. 2 the deviation of the surface elevation and momentum measured in the $L^\infty$ norm is plotted over time. As one can see, in case of initial still water (first setup) the deviations remain at machine accuracy. In the case of the perturbed initial state deviations can be observed in both variables, but they gradually vanish over time.

### 3.2 Oscillatory Flow in a Parabolic Bowl

In this problem an oscillatory flow in a domain with parabolic bottom topography is considered. The analytical solution to this numerically challenging test was originally derived by Thacker [17] and has been applied to several schemes (e.g. [14, 19]). The solution involves a periodical movement of the wet/dry interface at both sides of the basin. On the domain $[-5000, 5000]$ the bottom topography is defined by $b(x) = h_0(x/a)^2$, where $a = 3000$ and $h_0 = 10$ define the shape of the parabolic basin. Note that the boundary conditions for the domain should not matter since the boundary is in the dry part of the solution. An analytic solution of the water surface is then given by

$$h(x, t) + b(x) = h_0 - \frac{B^2}{4g}(1 + \cos(2\omega t)) - \frac{Bx}{2a}\sqrt{\frac{8h_0}{g}}\cos\omega t,$$

where we set $\omega = \sqrt{2gh_0}/a$ and $B = 5$. The initial momentum at $t = 0$ is set to zero over the whole domain. The solution is discretized using 200 cells and a timestep 0.05, which corresponds to an approximate CFL number of 0.01.

**Fig. 3** Free surface elevation of an oscillatory flow in a parabolic bowl. Initial condition (*top left*), Numerical (*black*) and exact solution (*gray dashed*) at times $t = 1000$, $t = 2000$, $t = 3000$

The simulation is executed until $t_{max} = 3000$, when the flow has oscillated a bit more than two periods. The initial surface elevation and the numerical solution compared to the analytical at times 1000, 2000, 3000 is shown in Fig. 3. Only small deviations can be observed. The largest differences arise at the wet/dry interface, where the numerical solution lags a bit behind. This behavior probably has to do with the limiting within the wet domain. At the end of the simulation, the total energy $E = \int_\Omega hu^2/2 + gh(h/2 + b)$ has only decreased by 1.2 % of its initial value. Furthermore, we note that a relatively small time step must be chosen to get a stable solution in this case. These problems slightly degrade the efficiency of the scheme and are currently being investigated.

### 3.3 Tsunami Runup onto a Sloping Beach

In a final test case the propagation of a tsunami wave onto a uniformly sloping beach is simulated. This benchmark problem was originally defined in [18]. Besides the slope of the beach the initial surface elevation and momentum with $(hu)(x, 0) \equiv 0$ is given. The solution is sought on the domain $[-500, 50000]$ and the bottom topography is set to $b(x) = 5000 - 0.1x$. On the right boundary of the domain a simple transparent boundary condition is implemented. However, the crucial task is

**Fig. 4** Tsunami runup onto a sloping beach. Initial surface elevation at $t = 0$ (*top left*), computed height from 1D DG method (*black*) compared to analytical solution (*gray dashed*) at times $t = 160$ (*top right*), $t = 175$ (*bottom left*), $t = 220$ (*bottom right*)

to correctly simulate the inundation process on the interval $[-200, 800]$. The analytic solution at times $t = 160$, $175$ and $220$ can be derived by the initial-value-problem technique introduced in [6] and is given in [18].

In the presented simulation, the domain is discretized into 1010 cells and the timestep is 0.05, which approximately corresponds to a CFL number of 0.22 at the deepest point (right side) of the domain. The results compared to the analytical solution are displayed in Fig. 4. Considering the relatively coarse discretization, the inundation process is well approximated with the numerical scheme. Also in this case, a small lag at the wet/dry interface can be observed.

## 4 Conclusion

We introduced an efficient treatment of semi-dry cells, which occur in the coastal area at the wet/dry interface, in a DG inundation scheme of the shallow water equations. By neglecting possibly unphysical surface gradients, and redistributing the tendencies of the local quantities to the wet part of the cell, we obtain an efficient, stable and robust inundation scheme in one space dimension. Furthermore, the method conserves mass and is well-balanced for nearly still water states. This is demonstrated by the three test cases shown in Sect. 3, which also illustrate the scheme's ability for long term integrations like in the parabolic bowl test case.

Since our method does not rely on specific one-dimensional features of the discretization, we expect to generalize the findings to our triangular two-dimensional

near-realistic setting [4]. In this respect, the combination with discretizations of other near-coast processes like friction and non-hydrostatic effects through proper extensions of the shallow water equations will be also investigated.

# References

1. Audusse, E., Bouchut, F., Bristeau, M.O., Klein, R., Perthame, B.: A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. SIAM J. Sci. Comput. **25**(6), 2050–2065 (2004). doi:10.1137/S1064827503431090
2. Bates, P.D., Hervouet, J.M.: A new method for moving-boundary hydrodynamic problems in shallow water. Proc. R. Soc. Lond. A **455**(1988), 3107–3128 (1999). doi:10.1098/rspa.1999.0442
3. Behrens, J.: Numerical methods in support of advanced Tsunami early warning, Chap. 14. In: Freeden, W. (ed.) Handbook of Geomathematics, pp. 399–416. Springer, Heidelberg (2010)
4. Beisiegel, N., Behrens, J., Castro, C.E.: Development of an adaptive discontinuous Galerkin inundation model. J. Comput. Phys. (Submitted) (2013)
5. Bunya, S., Kubatko, E.J., Westerink, J.J., Dawson, C.: A wetting and drying treatment for the Runge-Kutta discontinuous Galerkin solution to the shallow water equations. Comput. Methods Appl. Mech. Eng. **198**, 1548–1562 (2009). doi:10.1016/j.cma.2009.01.008
6. Carrier, G.F., Wu, T.T., Yeh, H.: Tsunami run-up and draw-down on a plane beach. J. Fluid Mech. **475**, 79–99 (2003). doi:10.1017/S0022112002002653
7. Chen, C., Qi, J., Li, C., Beardsley, R.C., Lin, H., Walker, R., Gates, K.: Complexity of the flooding/drying process in an estuarine tidal-creek salt-marsh system: an application of FVCOM. J. Geophys. Res. **113**, C07, 052 (2008). doi:10.1029/2007jc004328
8. Cockburn, B., Shu, C.W.: The Runge-Kutta discontinuous Galerkin method for conservation laws V: multidimensional systems. J. Comput. Phy. **141**(2), 199–224 (1998). doi:10.1006/jcph.1998.5892
9. Giraldo, F.X., Warburton, T.: A high-order triangular discontinuous Galerkin oceanic shallow water model. Int. J. Numer. Meth. Fluids **56**(7), 899–925 (2008). doi:10.1002/fld.1562
10. Gottlieb, S., Shu, C.W., Tadmor, E.: Strong stability-preserving high-order time discretization methods. SIAM Rev. **43**(1), 89–112 (2001). doi:10.1137/S003614450036757X
11. Harig, S., Chaeroni, C., Pranowo, W.S., Behrens, J.: Tsunami simulations on several scales: comparison of approaches with unstructured meshes and nested grids. Ocean Dyn. **58**, 429–440 (2008). doi:10.1007/s10236-008-0162-5
12. Hesthaven, J.S., Warburton, T.: Nodal discontinuous Galerkin methods: algorithms, analysis, and applications. Springer, New York (2008)
13. Kärnä, T., de Brye, B., Gourgue, O., Lambrechts, J., Comblen, R., Legat, V., Deleersnijder, E.: A fully implicit wetting-drying method for DG-FEM shallow water models, with an application to the scheldt estuary. Comput. Methods Appl. Mech. Eng. **200**(5–8), 509–524 (2011). doi:10.1016/j.cma.2010.07.001
14. Kesserwani, G., Liang, Q.: Well-balanced RKDG2 solutions to the shallow water equations over irregular domains with wetting and drying. Comput. Fluids **39**, 2040–2050 (2010). doi:10.1016/j.compfluid.2010.07.008
15. Shu, C.W.: TVB uniformly high-order schemes for conservation laws. Math. Comput. **49**, 105–121 (1987). doi:10.1090/S0025-5718-1987-0890256-5

16. Shu, C.W., Osher, S.: Efficient implementation of essentially non-oscillatory shock-capturing schemes. J. Comput. Phys. **77**(2), 439–471 (1988). doi:10.1016/0021-9991(88)90177-5
17. Thacker, W.C.: Some exact solutions to the nonlinear shallow-water wave equations. J. Fluid Mech. **107**, 499–508 (1981). doi:10.1017/S0022112081001882
18. The Third International Workshop on Long-Wave Runup Models: Benchmark problem #1: Tsunami runup onto a plane beach (2004). http://isec.nacse.org/workshop/2004_cornell/bmark1.html
19. Xing, Y., Zhang, X., Shu, C.W.: Positivity-preserving high order well-balanced discontinuous Galerkin methods for the shallow water equations. Adv. Water Resour. **33**(12), 1476–1493 (2010). doi:10.1016/j.advwatres.2010.08.005

# Comparison of Cell-Centered and Staggered Pressure-Correction Schemes for All-Mach Flows

**Nicolas Therme and Chady Zaza**

**Abstract** Defining a robust scheme for solving the compressible Euler equations at all-Mach number is a challenging issue. We consider here an original pressure-correction scheme which solves the internal energy using a specific corrective term, ensuring the positivity of the internal energy and the global consistency of the scheme. The scheme has already proved its effectiveness on several Riemann problems with both staggered and cell-centered discretizations. We test here these two discretizations against the incompressible limit of the Euler equations.

## 1 Introduction

We address in this paper the compressible Euler equations written with the internal energy as energy variable:

$$\partial_t \rho + \mathrm{div}(\rho \mathbf{u}) = 0, \tag{1a}$$

$$\partial_t(\rho \mathbf{u}) + \mathrm{div}(\rho \mathbf{u} \otimes \mathbf{u}) + \boldsymbol{\nabla} p = 0, \tag{1b}$$

$$\partial_t(\rho e) + \mathrm{div}(\rho e \mathbf{u}) + p\,\mathrm{div}\,\mathbf{u} = 0, \tag{1c}$$

$$p = (\gamma - 1)\rho e, \tag{1d}$$

where $t$ stands for the time ; $\rho$, $\mathbf{u}$, $p$ and $e$ are the density, velocity, pressure and internal energy respectively, and $\gamma > 1$ is a coefficient specific to the fluid. The problem is

N. Therme
IRSN/PSN-RES/SA2I/LIE, CEA Cadarache, Bat. 288, 13115 St Paul-lez-Durance, France
e-mail: nicolas.therme@irsn.fr

C. Zaza (✉)
CEA/DEN/DANS/DM2S/STMF/LMEC, CEA Cadarache, Bat. 238,
13108 St Paul-lez-Durance, France
e-mail: chady.zaza@cea.fr; chady.zaza@latp.univ-mrs.fr

975

defined over $\Omega \times (0, T)$, where $\Omega$ is an open bounded connected subset of $\mathbb{R}^d$, $1 \leq d \leq 3$, and $(0, T)$ is a finite time interval.

Defining a robust scheme for the numerical solution of the compressible Euler equations at all Mach number is a challenging issue. Indeed, in the zero Mach limit, the pressure gradient has a singular limit and the acoustic time scale vanishes [1]. As a result approximate Riemann solvers face severe limitations, among which the loss of accuracy of the pressure gradient approximation and the time step limitation. Pressure-correction methods may be relevant for addressing this issue, in particular because of their built-in stability properties.

While pressure-correction schemes were originally introduced for the incompressible Navier-Stokes equations [3, 12] many extensions to compressible flows have been attempted [4, 9]. In this work we compare two finite volume discretizations—staggered and cell-centered—of an original pressure-correction scheme first introduced in [8, 10].

The use of the internal energy as energy variable is motivated by our will to control its positivity through the numerical scheme. The internal energy balance must be discretized carefully in order to force the scheme to be consistent with the total energy equation. Indeed, similarly to the continuous case, we obtain a (discrete) kinetic energy equation from the (discrete) momentum balance and the (discrete) mass balance in which there is a numerical diffusion term. This term must be compensated in the discrete internal energy balance so that the sum of the internal and kinetic energy equations yields the correct total energy equation.

## 2 Pressure Correction Scheme

We first introduce the pressure correction method in a semi-discrete time setting. Let $\delta t$ be a time discretization step. We define the discrete time $t^n = n\delta t$ with $t^N = T$ and $N = \lfloor T/\delta t \rfloor$. The pressure-correction scheme reads:

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

$$\frac{1}{\delta t}\left(\rho^n \tilde{\mathbf{u}}^{n+1} - \rho^{n-1}\tilde{\mathbf{u}}^n\right) + \mathrm{div}\left(\rho^n \tilde{\mathbf{u}}^{n+1} \otimes \mathbf{u}^n\right) + \sqrt{\frac{\rho^n}{\rho^{n-1}}}\nabla p^n = 0.$$

- Solve for $p^{n+1}$, $\mathbf{u}^{n+1}$, $\rho^{n+1}$ and $e^{n+1}$ the non-linear system:

$$\frac{\rho^n}{\delta t}\left(\mathbf{u}^{n+1} - \tilde{\mathbf{u}}^{n+1}\right) + \nabla p^{n+1} - \sqrt{\frac{\rho^n}{\rho^{n-1}}}\nabla p^n = 0,$$

$$\frac{1}{\delta t}\left(\rho^{n+1} - \rho^n\right) + \mathrm{div}\left(\rho^{n+1}\mathbf{u}^{n+1}\right) = 0,$$

$$\frac{1}{\delta t} \left( \rho^{n+1} e^{n+1} - \rho^n e^n \right) + \mathrm{div} \left( \rho^{n+1} e^{n+1} \mathbf{u}^{n+1} \right) + p^{n+1} \mathrm{div} \left( \mathbf{u}^{n+1} \right) = 0,$$
$$p^{n+1} = (\gamma - 1) \rho^{n+1} e^{n+1}.$$

The first step is a classical semi-implicit discretization of the momentum balance to obtain a predicted velocity. The second step is a non-linear pressure correction step which combines the mass balance and the internal energy balance. This non-linear coupling is important to ensure the positivity of the energy. It is solved using Newton's method.

## 3 Spatial Discretization

We suppose that the boundaries of the domain are sections of hyperplanes normal to a coordinate axis. Let $\mathcal{M}$ be a decomposition of $\Omega$. The cells are either rectangles ($d = 2$) or rectangular parallelepipeds ($d = 3$). By $\mathcal{E}$ and $\mathcal{E}(K)$ we denote the set of all $(d - 1)$-faces $\sigma$ of the mesh and of the element $K \in \mathcal{M}$ respectively. The set of faces included in the boundary of $\Omega$ is denoted by $\mathcal{E}_{\mathrm{ext}}$ and the set of internal faces (i.e. $\mathcal{E} \setminus \mathcal{E}_{\mathrm{ext}}$) is denoted by $\mathcal{E}_{\mathrm{int}}$; a face $\sigma \in \mathcal{E}_{\mathrm{int}}$ separating the cells $K$ and $L$ is denoted by $\sigma = K|L$. The outward normal vector to a face $\sigma$ of $K$ is denoted by $\mathbf{n}_{K,\sigma}$. For $K \in \mathcal{M}$ and $\sigma \in \mathcal{E}$, we denote by $|K|$ the measure of $K$ and by $|\sigma|$ the $(d-1)$-measure of the face $\sigma$. For $1 \leq i \leq d$, we denote by $\mathcal{E}^{(i)} \subset \mathcal{E}$ and $\mathcal{E}_{\mathrm{ext}}^{(i)} \subset \mathcal{E}_{\mathrm{ext}}$ the subset of the faces of $\mathcal{E}$ and $\mathcal{E}_{\mathrm{ext}}$ respectively which are perpendicular to the $i$th unit vector of the canonical basis of $\mathbb{R}^d$. The definition of the divergence operator is similar in both the cell-centered and the staggered scheme. For $(\mathbf{u}_\sigma^n)_{\sigma \in \mathcal{E}}$, we set:

$$\text{for } K \in \mathcal{M}, \quad (\mathrm{div}\,\mathbf{u})_K^n = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma|\, u_{K,\sigma}^n, \qquad (2)$$

with $u_{K,\sigma}^n = \mathbf{u}_\sigma^n \cdot \mathbf{n}_{K,\sigma}$ the advecting velocity.

### 3.1 Cell-Centered Scheme

The unknowns are associated to the cells of the mesh $\mathcal{M}$ and are denoted by:

$$\{\rho_K,\ e_K,\ p_K,\ \mathbf{u}_K,\ K \in \mathcal{M}\}.$$

We first explain the initial discrete conditions: $\rho^0$, $p^0$ and $\mathbf{u}^0$ are given; then we set for $K \in \mathcal{M}$ and $1 \leq i \leq d$:

$$\rho_K^0 = \frac{1}{|K|} \int_K \rho_0(\mathbf{x}) \, d\mathbf{x}, \quad e_K^0 = \frac{1}{|K|} \int_K e_0(\mathbf{x}) \, d\mathbf{x}, \quad \text{and} \quad u_{K,i}^0 = \frac{1}{|K|} \int_K (\mathbf{u}_0(\mathbf{x}))_i \, d\mathbf{x}.$$

The fully discrete scheme then reads, for $n = 0, 1, \ldots, N - 1$:

- *Velocity prediction step*:

$$\frac{|K|}{dt} (\rho_K^n \tilde{\mathbf{u}}_K^{n+1} - \rho_K^{n-1} \mathbf{u}_K^n) + \sum_{\sigma \in \mathcal{E}(K)} \tilde{\mathbf{u}}_\sigma^{n+1} F_{K,\sigma}^n + \sqrt{\frac{\rho_K^n}{\rho_K^{n-1}}} |K| (\nabla p)_K^n = 0. \quad (3)$$

- *Projection step*: solve the non-linear system

$$\mathbf{u}_K^{n+1} = \tilde{\mathbf{u}}_K^{n+1} - \frac{dt}{\rho_K^n} \left( (\nabla p)_K^{n+1} - \sqrt{\frac{\rho_K^n}{\rho_K^{n-1}}} (\nabla p)_K^n \right), \quad (4a)$$

$$\frac{|K|}{dt} (\rho_K^{n+1} - \rho_K^n) + \sum_{\sigma \in \mathcal{E}(K)} F_{K,\sigma}^{n+1} = 0, \quad (4b)$$

$$\frac{|K|}{dt} (\rho_K^{n+1} e_K^{n+1} - \rho_K^n e_K^n) + \sum_{\sigma \in \mathcal{E}(K)} e_\sigma^{n+1} F_{K,\sigma}^{n+1} + p_K^{n+1} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| u_{K,\sigma}^{n+1} - S_K^n = 0,$$
$$(4c)$$

$$p_K^{n+1} = (\gamma - 1) \rho_K^{n+1} e_K^{n+1}, \quad (4d)$$

where $\tilde{\mathbf{u}}_\sigma^{n+1}$ in (3) is a centered interpolation of the velocity, $F_{K,\sigma}^{n+1} = |\sigma| \rho_\sigma^{n+1} u_{K,\sigma}^{n+1}$ is the mass flux and $\rho_\sigma^{n+1}$, $e_\sigma^{n+1}$ are upwind interpolations with respect to the sign of $u_{K,\sigma}^{n+1}$ and $F_{K,\sigma}^{n+1}$ respectively. In the expression of the advecting velocity, we use a centered interpolation of the velocity at the face $\sigma$. In order to ensure the consistency of the scheme, the pressure gradient is constructed by duality with the discrete divergence of the velocity and reads:

$$(\nabla p)_K^n = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}(K)} |\sigma| p_\sigma^n \mathbf{n}_{K,\sigma}, \quad (5)$$

with $p_\sigma^n$ a centered interpolation of the pressure at face $\sigma$.

The corrective term $S_K^n$ is defined as:

$$S_K^n = \frac{|K|}{2 \, dt} \rho_K^{n-1} (\tilde{\mathbf{u}}_K^{n+1} - \mathbf{u}_K^n)^2. \quad (6)$$

## *3.2 Staggered Scheme*

The space discretization is staggered, using the Marker-And Cell (MAC) scheme. The degrees of freedom for scalar variables are still associated to the cells of the mesh, but the degrees of freedom for the $i$th component of the velocity are defined at the center of the faces $\sigma \in \mathcal{E}^{(i)}$, so the whole set of discrete velocity unknowns reads:

$$\{u_{\sigma,i}, \ \sigma \in \mathcal{E}^{(i)}, \ 1 \le i \le d\}.$$

We introduce dual meshes for each direction $i$ centered on $\sigma \in \mathcal{E}^{(i)}$, which are used for the finite volume approximation of the time derivative and convection terms in the momentum balance. For $\sigma = K|L \in \mathcal{E}^{(i)}$, we build a dual cell $D_\sigma$ made of two half cells $D_{K,\sigma}$ and $D_{L,\sigma}$ included in $K$ and $L$ respectively. Each cell $D_{K,\sigma}$ is a rectangle or a rectangular parallelepiped of basis $\sigma$ and of measure $|K|/2$. We denote by $|D_\sigma|$ the measure of $D_\sigma$ and by $\epsilon = D_\sigma | D_{\sigma'}$ the face separating $D_\sigma$ and $D_{\sigma'}$. We denote by $\tilde{\mathcal{E}}$ the set of dual faces, $\tilde{\mathcal{E}}_{\text{int}}^{(i)}$ the internal faces in the direction $i$ and $\tilde{\mathcal{E}}(D_\sigma)$ those belonging to $D_\sigma$.

We will only point out the major changes with respect to the cell-centered scheme. Initial conditions for the velocity differ from the cell-centered scheme only for the velocities, which are now defined on the dual cells:

$$\forall \sigma \in \mathcal{E}_{\text{int}}, \quad u_{\sigma,i}^0 = \frac{1}{|D_\sigma|} \int_{D_\sigma} (\mathbf{u}_0(\mathbf{x}))_i \, d\mathbf{x}. \tag{7}$$

The definition of the divergence operator is the same as before but the discrete gradient is now defined on the dual mesh:

$$\forall \sigma = K|L \in \mathcal{E}_{\text{int}}, \quad (\nabla p)_\sigma^n = \frac{|\sigma|}{|D_\sigma|}(p_L - p_K)\mathbf{n}_{K,\sigma}. \tag{8}$$

Equations for scalar variables have just minor changes. Unlike the cell-centered discretization the convective fluxes $u_{K,\sigma}^{n+1}$ are obtained without interpolation as the velocity unknowns are defined on the edges. We still use an upwind interpolation for $\rho_\sigma$ and $e_\sigma$ in (4b) and (4c) respectively. We need to rewrite the velocity updates (3) and (4a) on the dual mesh, which read for all $i \in [1, d]$, for all $\sigma \in \mathcal{E}_{\text{int}}^{(i)}$:

$$\frac{|D_\sigma|}{dt}(\rho_{D_\sigma}^n \tilde{u}_{\sigma,i}^{n+1} - \rho_{D_\sigma}^{n-1} u_{\sigma,i}^n) + \sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} \tilde{u}_{\epsilon,i}^{n+1} F_{\sigma,\epsilon}^n + \sqrt{\frac{\rho_{D_\sigma}^n}{\rho_{D_\sigma}^{n-1}}}|D_\sigma|(\nabla p)_\sigma^n = 0, \tag{9}$$

$$u_{\sigma,i}^{n+1} = \tilde{u}_{\sigma,i}^{n+1} - \frac{dt}{\rho_{D_\sigma}^n}\left((\nabla p)_\sigma^{n+1} - \sqrt{\frac{\rho_{D_\sigma}^n}{\rho_{D_\sigma}^{n-1}}}(\nabla p)_\sigma^n\right). \tag{10}$$

The dual fluxes $F_{\sigma,\epsilon}^n$ and densities $\rho_{D_\sigma}$ are defined such that we recover a discrete mass balance over the dual cells. As we mentioned in the introduction this is crucial for obtaining a discrete kinetic balance. The corrective term $S_K^n$ in the internal energy equation reads for all $K \in \mathcal{M}$:

$$S_K^n = \sum_{i=1}^{d} S_{K,i}^n, \quad \text{with} \quad S_{K,i}^n = \frac{1}{2}\rho_K^{n-1} \sum_{\sigma \in \mathcal{E}(K) \cap \mathcal{E}^{(i)}} \frac{|D_{K,\sigma}|}{\delta t}(\tilde{u}_{\sigma,i}^{n+1} - u_{\sigma,i}^n)^2. \quad (11)$$

## 4 Discrete Properties

Thanks to the upwind choice for the density in the mass balance both schemes preserve the positivity of the density, see [6, Lemma 2.2] for further details. With either discretization a kinetic energy balance can be derived from the momentum prediction equation:

**Proposition 1** *(Discrete kinetic energy balance for the cell-centered discretization) A solution to the cell-centered (resp. staggered) scheme satisfies (12) (resp. (13)):*

$$\frac{|K|}{2\delta t}\left[\rho_K^n(\mathbf{u}_K^{n+1})^2 - \rho_K^{n-1}(\mathbf{u}_K^n)^2\right] + \frac{1}{2}\sum_{\sigma \in \mathcal{E}(K)} \tilde{\mathbf{u}}_K^{n+1}\tilde{\mathbf{u}}_L^{n+1}F_{K,\sigma}^n,$$

$$+ \mathbf{u}_K^{n+1} \cdot \sum_{\sigma \in \mathcal{E}(K)} |\sigma|p_\sigma^{n+1}\mathbf{n}_{K,\sigma} + P_K^{n+1} - R_K^{n+1} = 0. \quad (12)$$

$$\frac{|D_\sigma|}{2\delta t}\left[\rho_{D_\sigma}^n(u_{\sigma,i}^{n+1})^2 - \rho_{D_\sigma}^{n-1}(u_{\sigma,i}^n)^2\right] + \frac{1}{2}\sum_{\epsilon \in \tilde{\mathcal{E}}(D_\sigma)} \tilde{u}_{\sigma,i}^{n+1}\tilde{u}_{\sigma',i}^{n+1}F_{\sigma,\epsilon}^{n+1}$$

$$+ \tilde{u}_{\sigma,i}^{n+1}|D_\sigma|(\nabla p^{n+1})_\sigma^{(i)} + P_\sigma^{n+1} - R_{\sigma,i}^{n+1} = 0, \quad (13)$$

*with the following source terms:*

$$R_K^{n+1} = -\frac{|K|}{2\delta t}\rho_K^{n-1}(\tilde{\mathbf{u}}_K^{n+1} - \mathbf{u}_K^n)^2, \quad R_{\sigma,i}^{n+1} = -\frac{|D_\sigma|}{2\delta t}\rho_{D_\sigma}^{n-1}(\tilde{u}_{\sigma,i}^{n+1} - u_{\sigma,i}^n)^2,$$

$$P_K^{n+1} = \frac{\delta t}{2}\left[\frac{1}{\rho_K^n}\left((\nabla p)_K^{n+1}\right)^2 - \frac{1}{\rho_K^{n-1}}\left((\nabla p)_K^n\right)^2\right],$$

$$P_\sigma^{n+1} = \frac{\delta t}{2}\frac{|\sigma|^2}{|D_\sigma|^2}\left[\frac{1}{\rho_{D_\sigma}^n}\left((\nabla p)_\sigma^{n+1}\right)^2 - \frac{1}{\rho_{D_\sigma}^{n-1}}\left((\nabla p)_\sigma^n\right)^2\right].$$

For both schemes, the corrective term $S_K^n$ in the internal energy balance is intended to compensate the terms $R_K^{n+1}$ and $R_\sigma^{n+1}$ which tend to zero: hence the expression of the corrective term $S_K$ given in (6) and (11). Note that $S_K$ is always positive, which

ensures the positivity of the internal energy thanks to the following proposition proved in [11].

**Proposition 2** *(Positivity of the internal energy) If* $\forall K \in \mathcal{M}$, $e_K^n \geq 0$, $S_K^n \geq 0$ *and* $\rho_K > 0$ *then* $\forall K \in \mathcal{M}$, $e_K^{n+1} \geq 0$.

## 5 Numerical Results

The two discretizations are tested with a recent benchmark introduced in [2]. This benchmark aims at testing numerical schemes for the compressible Euler equations against their incompressible limit when the Mach number $M$ tends to zero. It consists in a Taylor vortex in a unit square cavity $\Omega = [0, 1] \times [0, 1]$. The initial solution verifies the incompressible Euler equations and reads in non-dimensional variables:

$$\rho_0(\mathbf{x}) = 1, \quad \mathbf{u}_0(\mathbf{x}) = \begin{pmatrix} \sin(\pi x)\cos(\pi y) \\ \cos(\pi x)\sin(\pi y)) \end{pmatrix}, \quad p_0(\mathbf{x}) = \frac{1}{\gamma M^2} + \frac{1}{4}(\cos(2\pi x) + \cos(2\pi y))$$

However, it does not lead to a steady flow with the compressible Euler equations, as the homogeneous density induces variations of the entropy. The main idea is to study the behaviour of the scheme at two scales: the macroscopic scale (slow variations associated with time variable $t$) and the acoustic scale (fast variations associated with time variable $\tau = t/M$). Each flow variable is decomposed as $X(\mathbf{x}, \tau, t) = \bar{X}(\mathbf{x}, t) + \delta X(\mathbf{x}, \tau, t)$ with $\bar{X}(\mathbf{x}, t)$ its time average over the acoustic scale and $\delta X(\mathbf{x}, \tau, t)$ the fast time fluctuations. The asymptotic expansion of the non-dimensional flow variables with respect to the Mach number yields [2]:

$$p(\mathbf{x}, t) = p_0(\mathbf{x}) + M\delta P_3(\mathbf{x}, \tau, 0) + M^2(\bar{P}_4(\mathbf{x}, t) + \delta P_4(\mathbf{x}, \tau, t)) + o(M^2),$$

$$\rho(\mathbf{x}, t) = \rho_0(\mathbf{x}) + M^2\bar{\rho}_2(\mathbf{x}, t) + M^3\delta\rho_3(\mathbf{x}, \tau, 0) + M^4(\bar{\rho}_4(\mathbf{x}, t) + \delta\rho_4(\mathbf{x}, \tau, t)) + o(M^4).$$

The particular field chosen for initialization allows the derivation of an analytic solution, well suited for spectral analysis. We focus on two terms of the asymptotic expansion: $\bar{\rho}_2$, associated with the slow variations and $\mathscr{P}_3 = \delta P_3(\mathbf{x}, \tau, 0) + M(\bar{P}_4(\mathbf{x}, t) + \delta P_4(\mathbf{x}, \tau, t))$ associated with the fast variations. In practice these two terms are computed as $\bar{\rho}_2 = (\rho - \rho_0)/M^2$ and $\mathscr{P}_3 = (p - p_0)/M$.

Our numerical simulations are carried out on a $400 \times 400$ grid with $M = 0.1$ and $M = 0.01$. We observe very similar results for both cell-centered and staggered discretizations. At the macroscopic scale, the upwind diffusion damps the main modes of the density, which looks smooth at $T = 8.8$ (Fig. 1). As for the term $\bar{\rho}_2$, the oscillations of the solution are completely damped after $t = 4$ (Fig. 2, left). The Mach number does not appear to have any influence on this term. At the acoustic scale, the fluctuations of the pressure $\mathscr{P}_3$ on the short time interval $(0, 5)$ are also close with both discretizations (Fig. 2, center). After $t = 0.5$, the amplitude of the main mode of $\mathscr{P}_3$ (frequency $f = \sqrt{10}/2$) is decreased by two orders of magnitude.

**Fig. 1** Density field for the staggered discretization at $t = 0.5$, $t = 2$, $t = 4$ and $t = 8.8$ for $M = 0.1$. The density fields obtained with $M = 0.01$ and with the cell-centered discretization are the same



**Fig. 2** Evolution of $\bar{\rho}_2$ (*left*) at position $(0.5, 0.05)$ and $\mathscr{P}_3$ (*middle* and *right*) at position $(0.66, 0.05)$ for $M = 0.1$

The results of this benchmark do not feature spurious pressure modes for the cell-centered discretization as we might have expected. Indeed the internal energy balance (4c) can be reformulated as a non-linear equation on the pressure using the velocity update (4a) and the equation of state (4d):

$$
M^2 \left\{ \frac{|K|}{\delta t}(P_K^{n+1} - P_K^n) + \sum_{\sigma \in \mathscr{E}(K)} |\sigma|(P_\sigma^{n+1} - (\gamma - 1)P_K^{n+1}) \left[ \frac{\delta t}{2}(\mathbf{g}_{K,\sigma}^{n+1} + \mathbf{g}_{L,\sigma}^{n+1}) \cdot \mathbf{n}_{K,\sigma} \right. \right.
$$
$$
\left. \left. + \tilde{u}_{K,\sigma}^{n+1} \right] - (\gamma - 1)S_K^n \right\} + \sum_{\sigma \in \mathscr{E}(K)} |\sigma| \left[ \frac{\delta t}{2}(\mathbf{g}_{K,\sigma}^{n+1} + \mathbf{g}_{L,\sigma}^{n+1}) \cdot \mathbf{n}_{K,\sigma} + \tilde{u}_{K,\sigma}^{n+1} \right] = 0
$$

with the change of variables $P = p - 1/(\gamma M^2)$, $P_\sigma^{n+1}$ the upwind interpolation with respect to $F_{K,\sigma}^{n+1}$ and $\mathbf{g}_{K,\sigma}^{n+1} = (\rho_K^{n-1}\rho_K^n)^{-1/2}(\nabla P)_K^n - (\rho_K^n)^{-1}(\nabla P)_K^{n+1}$ for the cell-centered discretization. In the zero Mach limit this equation degenerates to the classical Poisson equation of the projection step of incremental pressure-correction schemes for incompressible flows. For the cell-centered discretization the resulting discrete Laplace operator introduces a decoupling between neighboring pressure unknowns, which is not the case with the staggered discretization. We managed to introduce spurious pressure modes for the cell-centered discretization by adding artificially a Dirac to the right hand side of this pressure equation at $t = 0$. However, these oscillations are quickly damped by the boundary conditions. We expect sustained spurious pressure modes in the case of an open boundary.

# References

1. Alazard, T.: A minicourse on the low Mach number limit. Discrete and Continuous Dyn. Syst. Ser. S **1**, 365–404 (2008)
2. Cadiou, A., Le Penven, L., Buffat, M.: Asymptotic and numerical analysis of an inviscid bounded vortex flow at low Mach number. J. Comput. Phys. **227**(18), 8268–8289 (2008)
3. Chorin, A.: Numerical Solution of the Navier-Stokes equations. Math. Comput. **22**, 745–762 (1968)
4. Degond, P., Tang, M.: All speed scheme for the low Mach number limit of the Isentropic Euler equation. Commun. Comput. Phys. **10**, 1–31 (2011)
5. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Ciarlet, P.G., Lions, J.L. (eds.) Handbook of Numerical Analysis. North-Holland, Amsterdam (2000)
6. Gastaldo, L., Herbin, R., Latché, J.C.: A discretization of the phase mass balance in fractional step algorithms for the drift-flux model. IMA J. Numer. Anal. **31**, 116–146 (2011)
7. Gastaldo, L., Herbin, R., Latché, J.C., Therme, N.: Consistent explicit staggered schemes with muscle and artificial viscosity techniques for the Euler equations. (in preparation) (2014)
8. Grapsas, D., Kheriji, W., Herbin, R., Latché, J.C.: An unconditionally stable finite-element-finite volume pressure correction scheme for the compressible Navier-Stokes equations. (in preparation) (2013)
9. Harlow, F., Amsden, A.: Numerical calculation of almost incompressible flow. J. Comput. Phys. **3**, 80–93 (1968)
10. Herbin, R., Kheriji, W., Latché, J.C.: Consistent pressure correction staggered schemes for the shallow water and Euler equations. M2AN (submitted) (2013)
11. Herbin, R., Latché, J.C., Zaza, C.: A cell-centered pressure-correction scheme for compressible flows. (in preparation) (2014)
12. Temam, R.: Sur l'approximation de la solution des équations de Navier-Stokes par la méthode des pas fractionnaires I. Arch. Ration. Mech. Anal. **32**, 135–153 (1969)

# Author Index