

Asymptotic preserving Deferred Correction Residual Distribution schemes

Rémi Abgrall* and Davide Torlo*[†]

Abstract

This work aims to extend the residual distribution (RD) framework to stiff relaxation problems. The RD is a class of schemes which is used to solve hyperbolic system of partial differential equations. Up to our knowledge, it was used only for systems with mild source terms, such as gravitation problems or shallow water equations. What we propose is an IMEX (implicit–explicit) version of the residual distribution schemes, that can resolve stiff source terms, without refining the discretization up to the stiffness scale. This can be particularly useful in various models, where the stiffness is given by topological or physical quantities, e.g. multiphase flows, kinetic models, viscoelasticity problems. Moreover, the provided scheme is able to catch different relaxation scales automatically, without losing accuracy. The scheme is asymptotic preserving and this guarantees that in the relaxation limit, we recast the expected macroscopic behaviour. To get a high order accuracy, we use an IMEX time discretization combined with a Deferred Correction (DeC) procedure, while naturally RD provides high order space discretization. Finally, we show some numerical tests in 1D and 2D for stiff systems of equations.

Keywords: Residual distribution, IMEX, relaxation, deferred correction, asymptotic preserving, kinetic model.

AMS subject classification: 65M12, 65L04, 65M60

1 Introduction

In many models, such as kinetic models, multiphase flows, viscoelasticity or relaxing gas flows, we have to deal with hyperbolic systems with relaxation terms. The relaxation term is often led by a parameter ε , the relaxation parameter, that can represent the mean free path, the average distance between two collisions of particles, the time needed to reach the equilibrium between two phases, etc. Expanding these equations asymptotically with respect to ε , one can find the limit equations that describe the average, effective or macroscopic physical behaviour [8, 17, 19].

In particular, we focus on the kinetic model proposed by Aregba-Driollet and Natalini in [8, 9]. This model is able to solve any hyperbolic system of equation, through an artificial relaxation, which leads to a linear advection system with a relaxation source term. It can be used to test classical hyperbolic systems in the relaxation limit case. This model must be subjected to a generalization of Whitham’s subcharacteristic condition [8, 17], which assures that we are adding numerical viscosity to the limit equations. We use this model to approximate transport linear equation, Burgers’ equation and Euler equation in 1D and 2D. There are various other models and physical problems which behave similarly to this kinetic model. The perspective is, in future, to extend the method to multiphase flows, viscoelasticity problems, and so on.

We use the residual distribution (RD) framework [3, 6, 13, 20] to discretize our space. This class of schemes is a generalization of finite element schemes and allows to recast different well

*Institut für Mathematik, Winterthurststrasse 190, CH 8057 Zürich, Switzerland.

[†]Corresponding author, (davide.torlo@math.uzh.ch).

known finite element, finite volume and discontinuous Galerkin schemes [5]. The main ingredients of the scheme are three: we have to compute total residuals for each cell of the discretized domain, then, we have to distribute each residual to degrees of freedom of the cell, finally, we sum all contributions at each node. In order to get a high order scheme, the RD is coupled with a Deferred Correction (DeC) iterative method to have computationally explicit schemes [4, 14, 18]. It needs two operators: the first one is a low order method, but easy to be inverted, while the second one, must be higher order, but we do not need to solve it directly. The coupling of these two allows to reach the high order through a few iterative intermediate steps. Thanks to this, we can produce a scheme which is fast, high order and stable. Up to our knowledge, RD was utilised only for hyperbolic equations with mild source terms, such as in gravitation problems or shallow water equations, but never on strongly stiff source terms.

To deal with the stiffness of the relaxation term, we have to introduce some special treatments. An explicit scheme with CFL conditions tuned on the macroscopic regime would, indeed, present instabilities. To properly catch the small scale of the microscopic model, one must classically recur to very fine time and space discretizations that are not always feasible in terms of computational time. The natural alternative is to use an implicit or semi-implicit formulation, which guarantees the stability of the scheme. We use an IMEX (implicit-explicit) scheme to treat implicitly the relaxation term and explicitly the advection part [17, 19]. Nevertheless, we propose a computationally explicit scheme, thanks to some properties of the model. Then, we introduce an IMEX discretization for the DeC RD schemes with the details of its implementation. Furthermore, we prove that the new DeC RD IMEX scheme is asymptotic preserving (AP). The AP property of a numerical method allows to preserve the asymptotic behaviour of the model from the microscopic to the macroscopic case. These schemes solve the microscopic equations, avoiding coupling of different models, and automatically are able to solve the asymptotic macroscopic limit in a robust way. In the appendix, we also provide a proof of the accuracy of the total scheme.

We show the performance of the high order scheme on some tests. In particular, we simulated different examples in 1D and 2D for linear transport equation and Euler equation. Thanks to these results, we validate the accuracy of our method and the capability of shock limiting along discontinuities.

The outline of the manuscript is as follows. In section 2 we present the kinetic model we want to solve and the conditions under which it is stable. In section 3 we describe the RD schemes for the spatial discretization with the DeC high order time discretization. In section 4, we need to adjust the time discretization according to an IMEX scheme, to deal with stiff source terms and we prove the asymptotic preserving property of the scheme. We show numerical results for 1D and 2D problems in section 5. Finally, in section 6, we describe the conclusions and some future investigations that may be done.

2 Kinetic relaxation model for hyperbolic systems

In this section, we introduce the kinetic relaxation model presented by D. Aregba-Driollet and R. Natalini in [8, 9]. This is a first step to solve general hyperbolic systems of conservation laws via a relaxed system.

Let $u : \Omega \subset \mathbb{R}^D \times [0, T] \rightarrow \mathbb{R}^K$ be a weak solution of the following system of equations

$$u_t + \sum_{d=1}^D \partial_{x_d} A_d(u) = 0 \quad (1)$$

with initial conditions $u(x, 0) = u_0(x)$. Here, $A_d : \mathbb{R}^K \rightarrow \mathbb{R}^K$ are locally Lipschitz continuous on

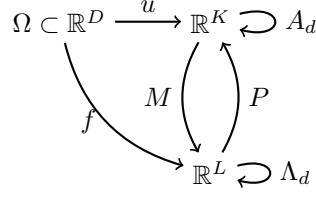


Figure 1: Relaxation functions

\mathbb{R}^K with values in \mathbb{R}^K . We approximate the problem with a relaxed system

$$f_t^\varepsilon + \sum_{d=1}^D \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} (M(Pf^\varepsilon) - f^\varepsilon), \quad f^\varepsilon(x, 0) = f_0^\varepsilon(x) \quad (2)$$

where $f^\varepsilon : \Omega \subset \mathbb{R}^D \times [0, T] \rightarrow \mathbb{R}^L$ with $M : \mathbb{R}^K \rightarrow \mathbb{R}^L$ Lipschitz continuous Maxwellian function, $P : \mathbb{R}^L \rightarrow \mathbb{R}^K$ a constant projection matrix ($L > K$) and Λ_d diagonal $L \times L$ matrices as sketched in figure 1.

Moreover, we require that for all u in a certain manifold of interest of \mathbb{R}^K the relations

$$\begin{cases} P(M(u)) = u \\ P\Lambda_d M(u) = A_d(u) \end{cases} \quad (3)$$

hold. If f^ε converges in some strong topology to a limit f and $Pf_0^\varepsilon \rightarrow u_0$, then Pf is a solution of the first system (1).

To show this, we define $u^\varepsilon := Pf^\varepsilon$, $v_j^\varepsilon := P\Lambda_j f^\varepsilon$ for $j = 1, \dots, D$. Then we have from (2) that

$$\begin{cases} \partial_t u^\varepsilon + \sum_{j=1}^D \partial_{x_j} v_j^\varepsilon = 0 \\ \partial_t v_d^\varepsilon + \sum_{j=1}^D \partial_{x_j} (P\Lambda_j \Lambda_d f^\varepsilon) = \frac{1}{\varepsilon} (A_d(u^\varepsilon) - v_d^\varepsilon), \quad \forall d \in \{1, \dots, D\} \end{cases} \quad (4)$$

Again, thanks to (2), we consider a formal expansion of f^ε in Taylor series with respect to ε in the form of

$$f^\varepsilon = M(u^\varepsilon) + \varepsilon g^\varepsilon + \mathcal{O}(\varepsilon^2), \quad (5)$$

from the second equation of (4) we can write $\forall d = 1, \dots, D$

$$v_d^\varepsilon = A_d(u^\varepsilon) - \varepsilon \left(\partial_t v_d^\varepsilon + \sum_{j=1}^D \partial_{x_j} (P\Lambda_d \Lambda_j f^\varepsilon) \right) + \mathcal{O}(\varepsilon^2) \quad (6)$$

$$= A_d(u^\varepsilon) - \varepsilon \left(\partial_t v_d^\varepsilon + \sum_{j=1}^D \partial_{x_j} (P\Lambda_d \Lambda_j M(u^\varepsilon)) \right) + \mathcal{O}(\varepsilon^2). \quad (7)$$

If we substitute this result in (4), we get

$$\partial_t u^\varepsilon + \sum_{d=1}^D \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^D \partial_{x_d} \left(\partial_t v_d^\varepsilon + \sum_{j=1}^D \partial_{x_j} (P\Lambda_d \Lambda_j M(u^\varepsilon)) \right) + \mathcal{O}(\varepsilon^2). \quad (8)$$

Now, we have that

$$\partial_t v_d^\varepsilon = \partial_t A_d(u^\varepsilon) + \mathcal{O}(\varepsilon) = A'_d(u^\varepsilon) \partial_t u^\varepsilon + \mathcal{O}(\varepsilon) = - \sum_{j=1}^D A'_d(u^\varepsilon) A'_j(u^\varepsilon) \partial_{x_j} u^\varepsilon + \mathcal{O}(\varepsilon). \quad (9)$$

Then, we eventually obtain up to second order in ε

$$\partial_t u^\varepsilon + \sum_{d=1}^D \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^D \partial_{x_d} \left(\sum_{j=1}^D B_{dj}(u^\varepsilon) \partial_{x_j} u^\varepsilon \right) \quad (10)$$

where

$$B_{dj}(u) := P \Lambda_d \Lambda_j M'(u) - A'_d(u) A'_j(u) \quad (11)$$

is a $K \times K$ matrix.

This limit equation is stable if the following condition holds:

$$\sum_{j,d=1}^D (B_{dj} \xi_j, \xi_d) \geq 0, \quad \forall \xi_1, \dots, \xi_D \in \mathbb{R}^K. \quad (12)$$

This property is a generalization of the *Whitham's subcharacteristic condition* [8, 17, 9].

We have to choose M, P, Λ that respect conditions (3) to completely define the kinetic model. First of all, let us take in consideration $L = N \times K$ with $P = (I_K, \dots, I_K)$ the juxtaposition of N blocks of identity matrices $I_K \in \mathbb{R}^{K \times K}$. Here, we can consider several $f_n^\varepsilon \in \mathbb{R}^K$ with $n = 1, \dots, N$ instead of a single vector $f^\varepsilon \in \mathbb{R}^{N \times K}$, several Maxwellians $M_n : \mathbb{R}^K \rightarrow \mathbb{R}^K$ and a block diagonal matrix $\forall d = 1, \dots, D$

$$\Lambda_d = \text{diag}(C_1^{(d)}, \dots, C_N^{(d)}) \quad C_n^{(d)} = \lambda_n^{(d)} I_K, \quad \text{with } \lambda_n^{(d)} \in \mathbb{R}, \forall n = 1, \dots, N.$$

With this formalism we can rewrite (2) as

$$\begin{cases} \partial_t f_n^\varepsilon + \sum_{d=1}^D \lambda_n^{(d)} \partial_{x_d} f_n^\varepsilon = \frac{1}{\varepsilon} (M_n(u^\varepsilon) - f_n^\varepsilon), & \forall n = 1, \dots, N \\ u^\varepsilon = \sum_{n=1}^N f_n^\varepsilon \end{cases} \quad (13)$$

Let us present the *diagonal relaxation method (DRM)*. Here $N = D + 1$. Then we have to define Maxwellians M_n and matrices $C_j^{(d)}$. Take $\lambda > 0$, that will be chosen according to *Whitham's subcharacteristic condition* (12), and

$$C_j^{(d)} = \begin{cases} -\lambda I_K & j = d \\ \lambda I_K & j = D + 1 \\ 0 & \text{else} \end{cases} \quad (14)$$

The Maxwellians can be defined as follows:

$$\begin{cases} M_{D+1}(u) = \left(u + \frac{1}{\lambda} \sum_{d=1}^D A_d(u) \right) / (D + 1) \\ M_j(u) = -\frac{1}{\lambda} A_j(u) + M_{D+1}(u) \end{cases} \quad (15)$$

For one-dimensional system of conservation laws this formulation coincides with Jin–Xin relaxation model [17], the simplest example that we can think of in this context. Indeed, if we set $u^\varepsilon := P f^\varepsilon$ and $v^\varepsilon := P \Lambda f^\varepsilon$, we get

$$\begin{cases} \partial_t u^\varepsilon + \partial_x v^\varepsilon = 0 \\ \partial_t v^\varepsilon + \partial_x u^\varepsilon = \frac{1}{\varepsilon} (A(u^\varepsilon) - v^\varepsilon). \end{cases} \quad (16)$$

3 Residual distribution schemes

Let us now introduce the spatial and time discretization given by RD schemes [1, 13] and DeC approach [4, 14].

3.1 Notation

Let us start introducing the notation of RD schemes. For sake of simplicity, we explain the RD approach for steady equations, the time derivative part will be discussed in section 3.3. So, we can focus on the following equation

$$\nabla \cdot A(U) - S(U) = 0.$$

We define a triangulation Ω_h on our domain Ω and denote by K the generic element of the mesh and by h the characteristic mesh size (implicitly supposing some regularity on the mesh). Following

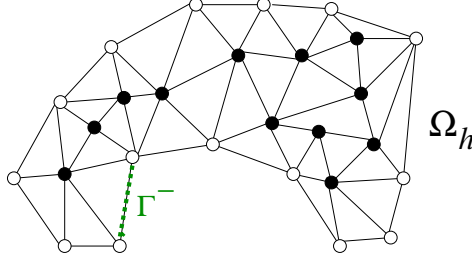


Figure 2: Triangulation of the domain Ω

the ideas of the Galerkin finite element method (FEM), we use a solution approximation space V_h given by globally continuous piecewise polynomials of degree k :

$$V_h = \{U \in \mathcal{C}^0(\Omega_h), U|_K \in \mathbb{P}^d, \forall K \in \Omega_h\}. \quad (17)$$

Now we can rewrite the numerical solution $U_h(x) \approx U(x)$ as a linear combination of basis functions $\varphi_\sigma \in V_h$:

$$U_h(x) = \sum_{\sigma \in \mathcal{D}_h} U_\sigma \varphi_\sigma(x) = \sum_{K \in \Omega_h} \sum_{\sigma \in K} U_\sigma \varphi_\sigma|_K(x), \quad \forall x \in \Omega \quad (18)$$

where \mathcal{D}_h is the set of all the degrees of freedom of Ω_h , so that $\{\varphi_\sigma : \sigma \in \mathcal{D}_h\}$ is a basis for V_h , and the coefficient U_σ must be found by a numerical method.

3.2 Residual distribution scheme

RD schemes can be summarized as follows.

1. Define $\forall K \in \Omega_h$ a fluctuation term (total residual)

$$\phi^K = \int_K (\nabla \cdot A(U_h) - S(U_h)) dx \quad (19)$$

2. Define a nodal residual ϕ_σ^K as a contribution to fluctuation term ϕ^K for each degree of freedom σ within the element K , so that the sum of all the contributions over an element is the fluctuation itself, i.e.,

$$\phi^K = \sum_{\sigma \in K} \phi_\sigma^K, \quad \forall K \in \Omega_h. \quad (20)$$

In appendix A or [2, 7] one can find more details on possible definitions of the nodal residuals.

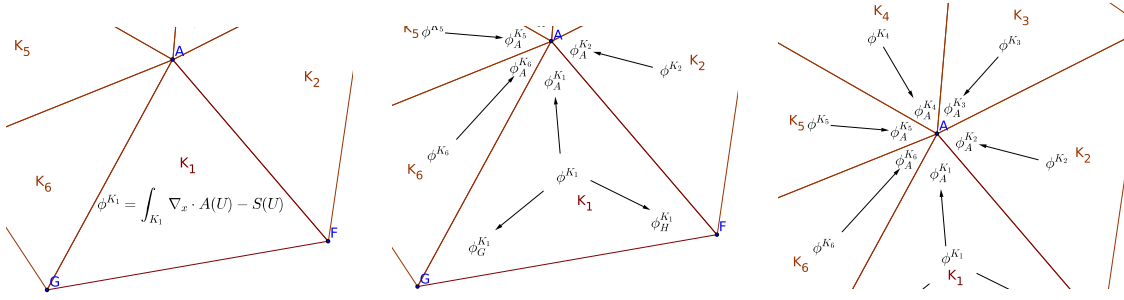


Figure 3: Defining total residual, nodal residuals and building the RD scheme

3. The resulting scheme is obtained by summing all the nodal residual contributions of one degree of freedom from different elements K , that is

$$\sum_{K|\sigma \in K} \phi_\sigma^K = 0, \quad \forall \sigma \in D_h. \quad (21)$$

This is a RD scheme.

The main sketch of the scheme is done in picture 3. The key of the scheme is the definition of nodal residuals. This choice is leading the whole spatial discretization. The scheme can be highly accurate in space, just choosing higher order polynomial basis functions and consistent nodal residuals. In [1, 4, 5] it has been shown that well known finite element or finite volume schemes (such as SUPG, DG, FV-WENO, etc.) can be rewritten in terms of RD, just choosing the proper nodal residuals.

Details and some examples of the schemes can be found in the appendix A.

3.3 Time discretization

For time discretization, we want to get a high order accurate approximation. To do so, we discretize the timestep $[t^n, t^{n+1}]$ into M subimesteps $[t^{n,0}, t^{n,1}], \dots, [t^{n,M-1}, t^{n,M}]$ and the variable U_h in time at each subimestep $U_h^{n,m}$ as in picture 4.

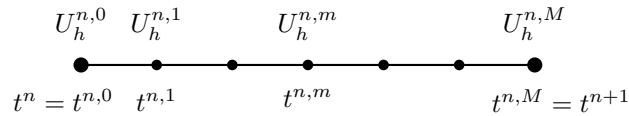


Figure 4: Subimesteps

Using the Picard–Lindelöf theorem, we can write for $m = 1, \dots, M$

$$U_h^{n,m} - U_h^n + \int_{t^n}^{t^{n,m}} (\nabla \cdot A(U_h(x, s)) - S(U_h(x, s))) ds = 0. \quad (22)$$

For sake of simplicity, we drop from now on the subscript h . More precisely, the scheme that we want to solve is a system of equations, where each entry is the discretization of (22) for a different

$m = 1, \dots, M$. In practice, we can write it as

$$\begin{aligned} \mathcal{L}_\sigma^2(U^{n,0}, \dots, U^{n,M}) &= \\ &= \begin{pmatrix} \sum_{K|\sigma \in K} \int_K \varphi_\sigma(U^{n,1} - U^{n,0}) dx + \sum_{K|\sigma \in K} \int_{t^{n,0}}^{t^{n,1}} \mathcal{I}_M(\phi_\sigma^K(U^{n,0}), \dots, \phi_\sigma^K(U^{n,M}), s) ds \\ \vdots \\ \sum_{K|\sigma \in K} \int_K \varphi_\sigma(U^{n,M} - U^{n,0}) dx + \sum_{K|\sigma \in K} \int_{t^{n,0}}^{t^{n,M}} \mathcal{I}_M(\phi_\sigma^K(U^{n,0}), \dots, \phi_\sigma^K(U^{n,M}), s) ds \end{pmatrix}. \end{aligned} \quad (23)$$

Here, we have M equations with M unknowns $U^{n,1}, \dots, U^{n,M}$, \mathcal{I}_M is an interpolation polynomial in nodes $\{t^{n,m}\}_{m=0}^M$ and the time integration is computed using quadrature formulas in the interpolation points. Of course, this system may contain a lot of nonlinear terms as functions of U , so we would like not to solve it directly. Nevertheless, the solution to (23) is what we are interested in. It is an approximation of the real solution with an accuracy of order $M + 1$ in time and $d + 1$ in space, where d is the degree of utilised polynomials.

The spirit of the DeC algorithm is to use two schemes, one high order and another one explicit or easy to solve. So, we introduce a first order approximation of the scheme \mathcal{L}^2 , that we will call \mathcal{L}^1 :

$$\begin{aligned} \mathcal{L}_\sigma^1(U^{n,0}, \dots, U^{n,M}) &= \\ &= \begin{pmatrix} (U_\sigma^{n,1} - U_\sigma^{n,0}) \sum_{K|\sigma \in K} \int_K \varphi_\sigma dx + \sum_{K|\sigma \in K} \int_{t^{n,0}}^{t^{n,1}} \mathcal{I}_0(\phi_\sigma^K(U^{n,0}), \dots, \phi_\sigma^K(U^{n,M}), s) ds \\ \vdots \\ (U_\sigma^{n,M} - U_\sigma^{n,0}) \sum_{K|\sigma \in K} \int_K \varphi_\sigma dx + \sum_{K|\sigma \in K} \int_{t^{n,0}}^{t^{n,M}} \mathcal{I}_0(\phi_\sigma^K(U^{n,0}), \dots, \phi_\sigma^K(U^{n,M}), s) ds \end{pmatrix}. \end{aligned} \quad (24)$$

The first simplification we applied is a mass lumping on the derivative in time, substituting U with U_σ . This is only possible if $|\mathcal{C}_\sigma| = \sum_K \int_K \varphi_\sigma(x) dx > 0$. For this reason, we will always consider Bernstein polynomials \mathbb{B}^d , which are nonnegative everywhere, instead of Lagrange polynomial \mathbb{P}^d .

The second one is in the residual part, where we substituted the high order interpolant \mathcal{I}_M with a piecewise constant interpolant \mathcal{I}_0 , which is explicit or easy to solve. An example of interpolant polynomial can be $\mathcal{I}_0(\phi_\sigma^K(U^{n,0}), \dots, \phi_\sigma^K(U^{n,M}), s) \equiv \phi_\sigma^K(U^{n,0})$. The detail of the interpolant will be given in section 4. The approximation error brought from these two approximations is a $\mathcal{O}(\Delta t + \Delta x)$.

3.4 Deferred Correction algorithm

Now, we present the deferred correction (DeC) algorithm to couple the two formulations. It was introduced by A. Dutt in [14] and we can see another approach in [18], but we follow the formulation by Abgrall in [4]. The aim of DeC schemes is to avoid implicit methods, without losing the high order of accuracy of a scheme. In our case, the high order method that we want to approximate is \mathcal{L}^2 of (23). To use the DeC procedure, we also need another method, which is easy and fast to be solved, we use diagonal mass matrix explicit methods, with low order of accuracy \mathcal{L}^1 , as in (24). The DeC algorithm is providing an iterative procedure that wants to approximate the solution of the \mathcal{L}^2 scheme U^* in the following way.

$$\begin{aligned} \mathcal{L}^1(U^{(1)}) &= 0, \\ \mathcal{L}^1(U^{(k)}) &= \mathcal{L}^1(U^{(k-1)}) - \mathcal{L}^2(U^{(k-1)}) \text{ with } k = 2, \dots, K, \end{aligned} \quad (25)$$

where K is the number of iterations that we compute. In particular, we need as many iteration as the order of accuracy that we want to reach: $K = d + 1 = M + 1$. Notice that, in every step, we solve the equations for the unknown variable $U^{(k)}$ which appears only in the \mathcal{L}^1 formulation, the one that can be solved easily. While \mathcal{L}^2 is only applied to already computed predictions of the solution $U^{(k-1)}$. Thus, we can state the following proposition as in [4].

Proposition 3.1. *Let \mathcal{L}^1 and \mathcal{L}^2 be two operators defined on \mathbb{R}^m , which depend on the discretization scale $\Delta \sim \Delta x \sim \Delta t$, such that*

- \mathcal{L}^1 is coercive with respect to a norm, i.e., $\exists \alpha_1 > 0$ independent of Δ , such that for any U, V we have that

$$\alpha_1 \|U - V\| \leq \|\mathcal{L}^1(U) - \mathcal{L}^1(V)\|,$$

- $\mathcal{L}^1 - \mathcal{L}^2$ is Lipschitz with constant $\alpha_2 > 0$ uniformly with respect to Δ , i.e., for any U, V

$$\|(\mathcal{L}^1(U) - \mathcal{L}^2(U)) - (\mathcal{L}^1(V) - \mathcal{L}^2(V))\| \leq \alpha_2 \Delta \|U - V\|.$$

We also assume that there exists a unique U_Δ^* such that $\mathcal{L}^2(U_\Delta^*) = 0$. Then, if $\eta := \frac{\alpha_2}{\alpha_1} \Delta < 1$, the DeC is converging to U^* and after k iterations the error $\|U^{(k)} - U^*\|$ is smaller than $\eta^k \|U^{(0)} - U^*\|$.

The proof of the proposition can be found in appendix B.1, while the proof of the properties of \mathcal{L}^1 and \mathcal{L}^2 , which depend on their definitions, can be found for our specific case in appendix B.2.

The theorem tells us that, if the method \mathcal{L}^2 is accurate with order of accuracy r , then we should perform r iterations for every timestep of the method and that we need only $r - 1$ sub-time steps. For example, if we use \mathbb{B}^1 basis functions, we will have 2 iterations of the DeC method (1 prediction and 1 correction) with 1 sub-time steps ($t^{n,0} = t^n$, $t^{n,1} = t^{n+1}$): this amounts to one version of the second order Runge Kutta method, see [20]. For \mathbb{B}^2 , we need 3 iterations (1 prediction, 2 corrections) and 2 sub-time steps ($t^{n,0} = t^n$, $t^{n,1} = \frac{1}{2}(t^n + t^{n+1})$, $t^{n,2} = t^{n+1}$) and so on. If not specified, in all our test cases we will use the same number of degree of polynomial, corrections-1 and subimesteps, i.e., $d = K - 1 = M$.

4 IMEX asymptotic preserving kinetic scheme

Before introducing an IMEX scheme, let us explain what is the problem concerning the kinetic model that we are considering. Solving equation (2), we have to be careful in treating the source term. If we discretize it in an explicit way, it would produce strongly stiff terms as $\varepsilon \rightarrow 0$. To classically solve this problem, one should take very small Δt values of the order of $\Delta t \sim \varepsilon$. On the other hand, the solution of the system would induce very long computational time. That is why, this can not always be a feasible way. The alternative is to treat implicitly the source term. Namely, we can use this type of time discretization:

$$\frac{f^{n+1,\varepsilon} - f^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^D \Lambda_d \partial_{x_d} f^{n,\varepsilon} = \frac{1}{\varepsilon} (M(P f^{n+1,\varepsilon}) - f^{n+1,\varepsilon}), \quad (26)$$

$$f^{0,\varepsilon}(x) = f_0^\varepsilon(x), \quad (27)$$

where the superscript index in f^n indicates the n -th timestep. This type of discretization is called implicit-explicit (IMEX), since the advection term is explicit, while the source term is implicit. This approach guarantees stability to the time discretization and we can relax the constraint on Δt until the usual CFL conditions proportional to the eigenvalue of the jacobian of the flux, which is λ in DRM model. Overall, the time-step can be chosen such that $\Delta t \leq \text{CFL} \lambda \Delta x$, where the CFL depends on the degree of the used polynomial basis functions.

As it is written, the time discretization (26) presents some nonlinear implicit terms. We can get rid of this technical problem, so that the scheme turns out to be computationally explicit. What we can notice is that the source is depending nonlinearly on $Pf^{n+1,\varepsilon} = u^{n+1,\varepsilon}$ and linearly on $f^{n+1,\varepsilon}$. To reach our goal, we can solve the following auxiliary equation for $u^{n+1,\varepsilon}$, which is the results of the multiplication of (26) by P and properties (3):

$$\frac{u^{n+1,\varepsilon} - u^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^D P\Lambda_d \partial_{x_d} f^{n,\varepsilon} = 0. \quad (28)$$

We can see that for this equation we are simply applying forward Euler method, which is explicit, since the source term turns out to be zero. So, we can solve it and then substitute $u^{n+1,\varepsilon}$ in equation (26) and solve it without recurring to implicit methods nor inversion of mass matrices. Indeed, the equation (26) can be rewritten in the following form, where the right-hand-side is explicit:

$$f^{n+1,\varepsilon} \left(\frac{1}{\Delta t} + \frac{1}{\varepsilon} \right) = \frac{f^{n,\varepsilon}}{\Delta t} - \sum_{d=1}^D \Lambda_d \partial_{x_d} f^{n,\varepsilon} + \frac{1}{\varepsilon} M(u^{n+1,\varepsilon}). \quad (29)$$

One can, indeed, express the variable f^{n+1} in the following way

$$f^{n+1,\varepsilon} = \frac{\varepsilon}{\Delta t + \varepsilon} f^{n,\varepsilon} - \frac{\varepsilon \Delta t}{\Delta t + \varepsilon} \sum_{d=1}^D \Lambda_d \partial_{x_d} f^{n,\varepsilon} + \frac{\Delta t}{\Delta t + \varepsilon} M(u^{n+1,\varepsilon}). \quad (30)$$

We can see that, in this formulation, ε does not appear alone in any denominator, so, for $\varepsilon \rightarrow 0$, $f^{n+1,\varepsilon}$ is well defined and tends to the Maxwellian $M(u^{n+1,\varepsilon})$.

4.1 Residual distribution IMEX operators

What we need to do now, is to apply the IMEX time discretization to the DeC and RD frameworks. This implies the change of the time discretization only of the operator \mathcal{L}^1 . Indeed, that is the only operator that we actually need to invert to get solutions of the DeC algorithm. While, we can not modify \mathcal{L}^2 because we do not want to drop the order of accuracy and because it will be anyway computed on previously computed solutions.

To do so, we want to choose the zero order interpolant \mathcal{I}_0 in a way that the source term is evaluated constantly on the end of the subimestep, namely in $t^{n,m}$, while the advection term is evaluated on the beginning of the timestep $t^{n,0}$, i.e.,

$$\mathcal{I}_0(\phi_\sigma^K(f^{n,0}), \dots, \phi_\sigma^K(f^{n,M}), s) \equiv \phi_{ad,\sigma}^K(f^{n,0}) + \phi_{source,\sigma}^K(f^{n,m}). \quad (31)$$

This requires a further definition of the nodal residuals that splits the source term and the advection part. The choice of the source residual is done accordingly to IMEX discretization. Indeed, what we require is its implicitness, the linear dependence on $f_\sigma^{n,m}$ and that it does not depend on other degrees of freedom. To reach these goals, we will perform a mass lumping on the whole source term and we evaluate everything in $t^{n,m}$. This results in

$$\phi_{source,\sigma}^K = \int_K \varphi_\sigma(x) \frac{M(Pf_\sigma^{n,m,\varepsilon}) - f_\sigma^{n,m,\varepsilon}}{\varepsilon} dx. \quad (32)$$

This allows us to collect $f_\sigma^{n,m,\varepsilon}$ on the left hand side of the equation and solve it explicitly. The advection part $\phi_{ad,\sigma}^K$ can be defined in different ways [13, 6, 1]. We give some examples in appendix A. Anyway, in this time discretization, it will be always explicit.

From now on we will drop the index n that indicates the timestep we are referring to and the index ε which refers to relaxation variables. They will be used only when necessary.

Overall, if we define $|\mathcal{C}_\sigma| := \int_\Omega \varphi_\sigma(x) dx$, the \mathcal{L}^1 operator will be at the m -th component

$$\mathcal{L}_{\sigma,u}^{1,m}(f^0, u^m) = |\mathcal{C}_\sigma|(u_\sigma^m - P f_\sigma^0) + \Delta t^m \sum_{K|\sigma \in K} P \phi_{ad,\sigma}^K(f^0); \quad (33a)$$

$$\begin{aligned} \mathcal{L}_\sigma^{1,m}(f^0, f^m) = & |\mathcal{C}_\sigma| \left(1 + \frac{\Delta t^m}{\varepsilon} \right) f_\sigma^m - |\mathcal{C}_\sigma| f_\sigma^0 + \\ & + \Delta t^m \sum_{K|\sigma \in K} \phi_{ad,\sigma}^K(f^0) - |\mathcal{C}_\sigma| \frac{\Delta t^m}{\varepsilon} M(u_\sigma^m). \end{aligned} \quad (33b)$$

We can see that both the equations of the \mathcal{L}^1 with the IMEX discretization are computationally explicit. Moreover, as before, we can see that, as $\varepsilon \rightarrow 0$, equation (33b) does not lead to terms with ε alone at the denominator. Indeed, it can be rewritten as

$$\begin{aligned} \mathcal{L}_\sigma^{1,m}(f^0, f^m) = & f_\sigma^m - \frac{\varepsilon}{\varepsilon + \Delta t^m} f_\sigma^0 + \\ & + \frac{\varepsilon \Delta t^m}{|\mathcal{C}_\sigma|(\varepsilon + \Delta t^m)} \sum_{K|\sigma \in K} \phi_{ad,\sigma}^K(f^0) - \frac{\Delta t^m}{\varepsilon + \Delta t^m} M(u_\sigma^m). \end{aligned} \quad (33c)$$

Finally, we can write a general term of the correction DeC procedure for the $(k+1)$ th correction and the m th subimestep. First, we have the u auxiliary equation

$$\begin{aligned} \mathcal{L}_{\sigma,u}^{1,m,(k+1)} - \mathcal{L}_{\sigma,u}^{1,m,(k)} + \mathcal{L}_{\sigma,u}^{2,m,(k)} = & |\mathcal{C}_\sigma|(u_\sigma^{m,(k+1)} - u_\sigma^{m,(k)}) + \\ & + \sum_{K|\sigma \in K} \left[\int_K \varphi_\sigma(x) (u^{m,(k)}(x) - u^{0,(k)}(x)) dx + \right. \\ & \left. + \int_{t^{n,0}}^{t^{n,m}} \mathcal{I}_M(P \phi_{\sigma,ad}^K(f^{0,(k)}), \dots, P \phi_{\sigma,ad}^K(f^{M,(k)}), s) ds \right]; \end{aligned} \quad (34a)$$

and, then, the f equation

$$\begin{aligned} \mathcal{L}_\sigma^{1,m,(k+1)} - \mathcal{L}_\sigma^{1,m,(k)} + \mathcal{L}_\sigma^{2,m,(k)} = & \\ |\mathcal{C}_\sigma| \left(1 + \frac{\Delta t^m}{\varepsilon} \right) (f_\sigma^{m,(k+1)} - f_\sigma^{m,(k)}) - & |\mathcal{C}_\sigma| \frac{\Delta t^m}{\varepsilon} \left(M(u_\sigma^{m,(k+1)}) - M(u_\sigma^{m,(k)}) \right) + \\ + \sum_{K|\sigma \in K} \left[\int_K \varphi_\sigma(x) (f^{m,(k)}(x) - f^{0,(k)}(x)) dx + \int_{t^0}^{t^m} \mathcal{I}_M(\phi_\sigma^K(f^{0,(k)}), \dots, \phi_\sigma^K(f^{M,(k)}), s) ds \right]. \end{aligned} \quad (34b)$$

Again, thanks to the factor $(1 + \frac{\Delta t^m}{\varepsilon})$ in front of the unknown, we are sure not to have any stiff term, even in the source of \mathcal{L}^2 .

4.2 AP property

An *asymptotic preserving (AP)* scheme preserves the asymptotic behaviour of the model from the microscopic to the macroscopic case. It solves the microscopic equations, avoiding coupling of different models, and, automatically, it is able to solve the asymptotic macroscopic limit as the relaxation parameter ε tends to its limit.

The behaviour of an AP scheme is sketched in figure 5. Let us call \mathcal{F}^ε the microscopic model which depends on ε and its asymptotic macroscopic limit $\mathcal{F}^0 := \lim_{\varepsilon \rightarrow 0} \mathcal{F}^\varepsilon$. We denote the numerical discretization of \mathcal{F}^ε as $\mathcal{F}_\Delta^\varepsilon$, where Δ is the mesh size and/or the time step length (in our case they are always linked by some CFL conditions). Then, we call the asymptotic limit as $\varepsilon \rightarrow 0$ of this scheme $\mathcal{F}_\Delta^0 := \lim_{\varepsilon \rightarrow 0} \mathcal{F}_\Delta^\varepsilon$ (for fixed Δ), if it exists. We can say that the scheme $\mathcal{F}_\Delta^\varepsilon$ is an AP scheme,

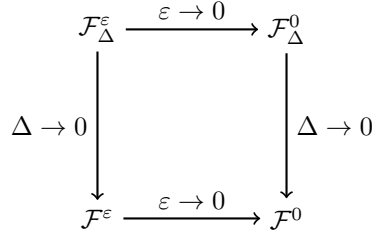


Figure 5: Asymptotic preserving schemes

if \mathcal{F}_Δ^0 is a consistent and stable approximation of \mathcal{F}^0 , i.e., $\mathcal{F}_\Delta^0 = \mathcal{F}^0 + \mathcal{O}(\Delta)$. In our model, the limit model \mathcal{F}^0 is the equation (1) and the relaxed model \mathcal{F}^ε is the equation (2). What we need to check is that the discrete model, namely, the scheme that we proposed, is asymptotic preserving. This implies that, first, we let $\varepsilon \rightarrow 0$, then the discretization scale $\Delta \rightarrow 0$. In other words, we can consider $\frac{\varepsilon}{\Delta} = o(1)$.

To start, let us suppose that the initial conditions u_0^ε and f_0^ε verify $f_0^\varepsilon = M(u_0^\varepsilon)$. We can prove that this is also true for the beginning of every timestep by induction. We want to show that at the end of each timestep we maintain the following relation for $u^\varepsilon = Pf^\varepsilon$ at the discrete level:

$$\frac{u^{n+1} - u^n}{\Delta t} + \sum_{d=1}^D \partial_{x_d} A_d(u^{n+1}) + \mathcal{O}(\varepsilon) + \mathcal{O}(\Delta) = 0. \quad (35)$$

To prove it by induction, we want to add the following relation in the induction hypothesis

$$\frac{f^{n+1} - f^n}{\Delta t} + \sum_{d=1}^D \partial_{x_d} \Lambda_d(f^{n+1}) - \frac{M(u^{n+1}) - f^{n+1}}{\varepsilon} + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta) = 0, \quad (36)$$

that implies

$$f^{n+1} = M(u^{n+1}) + \mathcal{O}(\varepsilon) + \mathcal{O}(\Delta). \quad (37)$$

The initial conditions verify the hypothesis (37) for $n = 0$. So we check the $n + 1$ th timestep given the relations (35) and (37) for the n th timestep. We start from the prediction $\mathcal{L}^1 = 0$, in forms (33a) and (33b). Since the scheme begins at each step with (33a), we can write $\forall m \in [1, \dots, M]$:

$$\frac{u_\sigma^{m,(1)} - u_\sigma^0}{\Delta t^m} + \frac{1}{|\mathcal{C}_\sigma|} \sum_{K|\sigma \in K} P\phi_{ad,\sigma}^K(f^0) = 0, \quad (38a)$$

and, if we use the fact that the sum of nodal residual is a consistent discretization of space derivatives as shown in [4], we get

$$\frac{u_\sigma^{m,(1)} - u_\sigma^0}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} P\Lambda_d f_\sigma^0 + \mathcal{O}(\Delta) = 0. \quad (38b)$$

Using the induction hypothesis on property (37), we obtain

$$\frac{u_\sigma^{m,(1)} - u_\sigma^{0,\varepsilon}}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} P\Lambda_d M(u_\sigma^0) + \mathcal{O}(\Delta) + \mathcal{O}(\varepsilon) = 0, \quad (38c)$$

while, using the properties in (3), the equation (38a) itself and the fact that A_d are Lipschitz continuous, we reach

$$\frac{u_\sigma^{m,(1)} - u_\sigma^0}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} A_d(u_\sigma^{m,(1)}) + \mathcal{O}(\Delta) + \mathcal{O}(\varepsilon) = 0. \quad (38d)$$

Then, from (33b) we can recast the second property (36), with similar reasoning:

$$0 = \frac{f_\sigma^{m,(1)} - f_\sigma^0}{\Delta t^m} + \frac{1}{|\mathcal{C}_\sigma|} \sum_{K|\sigma \in K} \phi_{ad,\sigma}^K(f^0) - \frac{M(u_\sigma^{m,(1)}) - f_\sigma^{m,(1)}}{\varepsilon}, \quad (39a)$$

then, using consistency of residuals, we can say that

$$0 = \frac{f_\sigma^{m,(1)} - f_\sigma^0}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} \Lambda_d f_\sigma^0 - \frac{M(u_\sigma^{m,(1)}) - f_\sigma^{m,(1)}}{\varepsilon} + \mathcal{O}(\Delta), \quad (39b)$$

and, finally, substituting (39a) in (39b), we get

$$0 = \frac{f_\sigma^{m,(1)} - f_\sigma^0}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} \Lambda_d f_\sigma^{m,(1)} - \frac{M(u_\sigma^{m,(1)}) - f_\sigma^{m,(1)}}{\varepsilon} + \mathcal{O}(\Delta) + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right). \quad (39c)$$

We proved that the prediction is asymptotic preserving, since it recast the limit equation (1). A more rigorous proof of a similar property for the norm convergence of the kinetic scheme is in [8, 9].

What is left to prove are the same properties for every correction $(k+1)$, using induction hypothesis on the previous correction (k) . For prediction $(k) = (1)$, we have already given the proof. Now, let us consider u equation in (34a) to prove property (35).

$$\begin{aligned} & \mathcal{L}_{\sigma,u}^{1,m} - \mathcal{L}_{\sigma,u}^{1,m} + \mathcal{L}_{\sigma,u}^{2,m} \\ &= \frac{u_\sigma^{m,(k+1)} - u_\sigma^{m,(k)}}{\Delta t^m} + \sum_{K|\sigma \in K} \int_K \varphi_\sigma(x) \frac{u^{m,(k)}(x) - u^{0,(k)}(x)}{|\mathcal{C}_\sigma| \Delta t^m} dx + \\ &+ \sum_{K|\sigma \in K} \frac{1}{|\mathcal{C}_\sigma| \Delta t^m} \int_{t^0}^{t^m} \mathcal{I}_M(P\phi_{\sigma,ad}^K(f^{0,(k)}), \dots, P\phi_{\sigma,ad}^K(f^{M,(k)}), s) ds = 0, \end{aligned} \quad (40a)$$

Then, we apply a mass lumping of time derivative term in \mathcal{L}^2 , moreover, we know that the quadrature of the interpolant is a first order approximation of any of its points, which are a consistent approximation of the flux. So,

$$\begin{aligned} & \frac{u_\sigma^{m,(k+1)} - u_\sigma^{m,(k)}}{\Delta t^m} + \frac{u_\sigma^{m,(k)} - u_\sigma^{0,(k)}}{\Delta t^m} + \mathcal{O}(\Delta) + \\ &+ \sum_{d=1}^D \partial_{x_d} P \Lambda_d f_\sigma^{m,(k)} + \mathcal{O}(\Delta) = 0, \end{aligned} \quad (40b)$$

we can now apply property (37) in the induction hypothesis on correction (k) and properties (3) to get

$$\frac{u_\sigma^{m,(k+1)} - u_\sigma^{0,(k)}}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} A_d(u_\sigma^{m,(k)}) + \mathcal{O}(\Delta) + \mathcal{O}(\varepsilon) = 0, \quad (40c)$$

and then we can substitute (40a) in (40c) to gain another $\mathcal{O}(\Delta)$ using also the Lipschitz continuity of fluxes A_d :

$$\frac{u_\sigma^{m,(k+1)} - u_\sigma^{0,(k+1)}}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} A_d(u_\sigma^{m,(k+1)}) + \mathcal{O}(\Delta) + \mathcal{O}(\varepsilon) = 0. \quad (40d)$$

Then, to prove property (37), we can proceed from (34b). We can split the three terms of the sum $\mathcal{L}_\sigma^{1,m,(k+1)} - \mathcal{L}_\sigma^{1,m,(k)} + \mathcal{L}_\sigma^{2,m,(k)}$. Let us start from $\mathcal{L}_\sigma^{2,m,(k)}$:

$$\begin{aligned} \mathcal{L}_\sigma^{2,m,(k)} &= \frac{1}{|\mathcal{C}_\sigma|} \sum_{K|\sigma \in K} \int_K \varphi_\sigma(x) \frac{f^{m,(k)}(x) - f^{0,(k)}(x)}{\Delta t^m} dx + \\ &+ \frac{1}{|\mathcal{C}_\sigma| \Delta t^m} \sum_{K|\sigma \in K} \int_{t^0}^{t^m} \mathcal{I}_M(\phi_\sigma^K(f^{0,(k)}), \dots, \phi_\sigma^K(f^{M,(k)}), s) ds \end{aligned} \quad (41a)$$

Then, we use a mass lumping on the time derivative, which brings an error of the order of Δ , the fact that the interpolant is a first order approximation of any of the interpolation points and that the residuals are consistent approximation of the flux and the source. So, we obtain

$$\mathcal{L}_\sigma^{2,m,(k)} = \frac{f_\sigma^{m,(k)} - f_\sigma^{0,(k)}}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} \Lambda_d f_\sigma^{m,(k)} + \frac{M(u_\sigma^{m,(k)}) - f_\sigma^{m,(k)}}{\varepsilon} + \mathcal{O}(\Delta). \quad (41b)$$

If we then use the induction hypothesis on (k) correction, we get

$$\mathcal{L}_\sigma^{2,m,(k)} = \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta). \quad (41c)$$

Analogously, for $\mathcal{L}_\sigma^{1,m,(k)}$ we can prove that it is an $\mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta)$ using the induction hypothesis. Finally, using what we just proved, we have that

$$\mathcal{L}_\sigma^{1,m,(k+1)} - \mathcal{L}_\sigma^{1,m,(k)} + \mathcal{L}_\sigma^{2,m,(k)} = \mathcal{L}_\sigma^{1,m,(k+1)} + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta) = 0. \quad (42a)$$

If we express explicitly the formula, we get

$$\begin{aligned} &\frac{f_\sigma^{m,(k+1)} - f_\sigma^{0,(k+1)}}{\Delta t^m} + \sum_{K|\sigma \in K} \frac{\phi_{ad,\sigma}^K(f^{0,(k+1)})}{|\mathcal{C}_\sigma|} \\ &+ \frac{M(u_\sigma^{m,(k+1)}) - f_\sigma^{m,(k+1)}}{\varepsilon} + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta) = 0, \end{aligned} \quad (42b)$$

Using the fact that the residuals are a consistent approximation of the fluxes and that the term at the m th subimestep is an approximation of the term at the 0th time step, up to an $\mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta)$ from (42b), we finally reach

$$\frac{f_\sigma^{m,(k+1)} - f_\sigma^{0,(k+1)}}{\Delta t^m} + \sum_{d=1}^D \partial_{x_d} \Lambda_d f_\sigma^{m,(k+1)} + \frac{M(u_\sigma^{m,(k+1)}) - f_\sigma^{m,(k+1)}}{\varepsilon} + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta) = 0. \quad (42c)$$

So, we proved property (36) for all subimesteps and corrections. This implies that the scheme is AP as $\varepsilon \rightarrow 0$ for any discretization scale Δ .

5 Numerical simulations

To validate the scheme we presented, we test the method on different problems. We will show 1D and 2D test cases for scalar equations and systems of equations. Generally, we will start from the asymptotic limit u and we will draw from that the whole kinetic systems for variable f . The shown results are related to the variable u . We have some parameter to choose in order to perform our tests. First of all, the convection coefficient λ , which should satisfy the Whitham's subcharacteristic conditions (12) and the relaxation parameter ε that will be often very small to get the asymptotic behaviour. Then, the CFL conditions, namely a bound on the size of Δt . Thanks to the scheme presented, we do not need CFL conditions linked to the source term, so, we can just choose them such that

$$\Delta t \leq \frac{\text{CFL } \Delta x}{\lambda}, \quad (43)$$

where λ is the convection parameter. While, with a standard RD DeC method without IMEX technique, the Δt should scale as

$$\Delta t \leq \text{CFL} \min \left\{ \frac{\Delta x}{\lambda}, \frac{\varepsilon}{\lambda} \right\},$$

which would require very small timesteps that lead to a huge computational demand. The CFL number depends on the degree of the polynomial chosen, and it scales as $\frac{1}{d}$, but for a comparison of the methods, we will choose it uniformly through different polynomial degrees. In all our computations we will also specify the θ_k parameter, which are leading the stabilization of the jump of the derivative, that we are using in the definition of the nodal residual. More details about the used nodal residual and the jump stabilization can be found in appendix A and in [6].

5.1 1D numerical tests

5.1.1 Burgers' equation

First of all, we start with 1D scalar equations. We want to approximate the Burgers' equation, i.e.,

$$\partial_t u(x, t) + \partial_x \left(\frac{u(x, t)^2}{2} \right) = 0, \quad x \in [0, 1], \quad t \in [0, T] \quad (44)$$

using the relaxation system (2). As initial condition, we take $u_0(x) = \sin(2\pi x)$ and $f_0(x) = M(u_0(x))$ and the boundary conditions are periodic. To satisfy Whitham's condition, we choose $\lambda = 2$, so that $|A'(u)| = |u| \leq \lambda$ in an area of interest.

In following figures some approximated solutions for different number of elements are shown. To solve the equation we used the scheme (58) in appendix A with $\theta_1 = 1$ and, only for \mathbb{B}^3 , we used $\theta_2 = 0.5$. The relaxation parameter is set to $\varepsilon = 10^{-9}$ and $\text{CFL} = 0.1$. Final time is $T = 0.5$. We can see in picture 6 that the scheme is well catching the shock position and, as the order of the polynomials increases, we can see improvements in the sharpness of the solution.

5.1.2 Convergence for linear transport equation

Then, we test our scheme with different orders to check the convergence rate. For all the smooth test cases, where we want to study the order of convergence, we use the scheme which involves only Galerkin residuals and stabilizations of jumps in derivative, as presented in [11] and in the scheme (52) in appendix A. We use a linear scalar transport equation $u_t + u_x = 0$ as limit equation with the relaxation system presented above, on domain $[0, 1]$. The initial condition is $u_0(x) = e^{-80(x-0.4)^2}$ and $f_0 = M(u_0)$, until final time $T = 0.12$ with periodic boundary conditions.

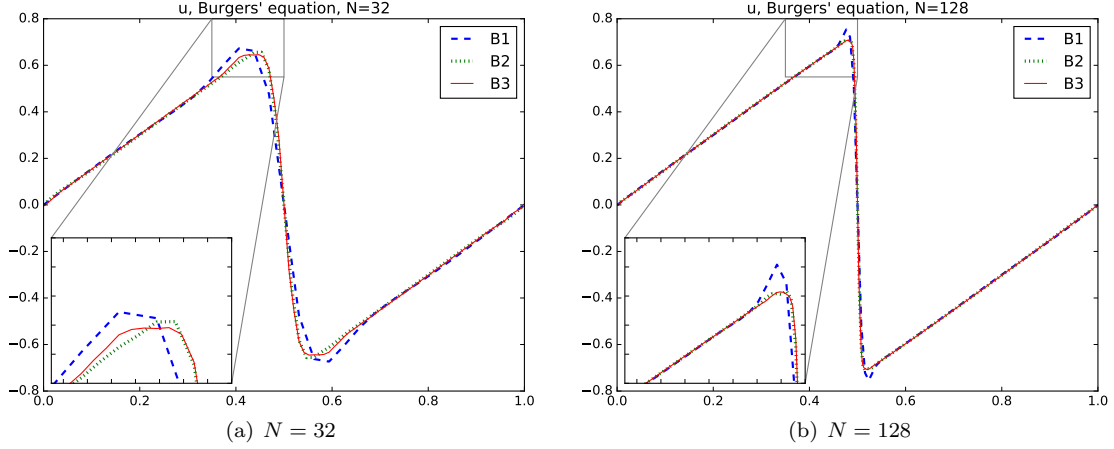


Figure 6: Burgers' equations

We use the relaxation coefficient $\varepsilon = 10^{-9}$, convection $\lambda = 1.5$ and CFL=0.1. In particular, for \mathbb{B}^1 we used $\theta_1 = 1$, for \mathbb{B}^2 we used $\theta_1 = 1, \theta_2 = 0$ and for \mathbb{B}^3 we used $\theta_1 = 1, \theta_2 = 5$. Final time of the solution is $T = 0.12$. For \mathbb{B}^3 we see that only increasing a bit the number of corrections with respect to the theoretical ones we achieve the correct slope for the error convergence, i.e., $K \gtrsim 7$. The reason of this behaviour is still under investigation. As we can see in figure 7(a), the convergence of the scheme is what we expected from theory.

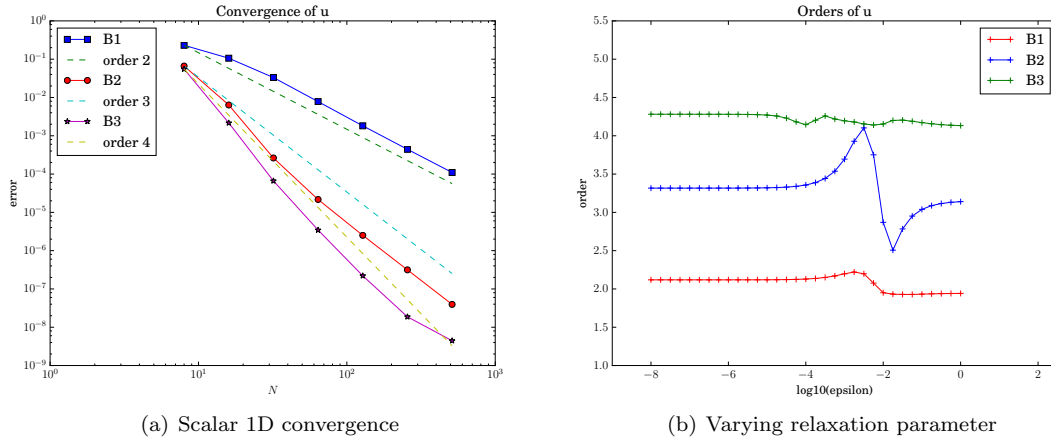


Figure 7: Scalar linear 1D test

Moreover, we can see in figure 7(b) that, also varying the relaxation parameter ε , the order of accuracy is the expected one. There are slight oscillations in particular for \mathbb{B}^2 solutions. This is a well known problem of order reduction as ε is approaching the magnitude of Δ , which affects lots of schemes, including some RK methods, as stated in [10]. Anyway, we can say that the scheme is getting an order of accuracy bigger or equal than the expected one, except for few mid-range values of ε . Moreover, we can state that the scheme is stable, for any value of ε we use.

5.1.3 Euler equation – Isentropic flow

Now, we can pass to systems of equations. In particular, we will focus on Euler equation

$$\begin{pmatrix} \rho \\ \rho v \\ E \end{pmatrix}_t + \begin{pmatrix} \rho v \\ \rho v^2 + p \\ (E + p)v \end{pmatrix}_x = 0 \quad (45)$$

on domain $[-1, 1]$, where ρ is the density, v the speed, p the pressure and E the total energy. The quantities are linked by the equation of state (EOS)

$$E = \frac{p}{\gamma - 1} + \frac{1}{2}\rho v^2. \quad (46)$$

To test the convergence of the scheme on 1D Euler equations, we use the case of isentropic flow, when $\gamma = 3$ and $p = \rho^\gamma$. With following initial conditions

$$\begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 + 0.5 \cdot \sin(\pi x) \\ 0 \\ \rho_0^\gamma \end{pmatrix} \text{ for } x \in [-1, 1],$$

final time $T = 0.1$ and periodic boundary conditions.

Now, we use $\varepsilon = 10^{-9}$, convection coefficient $\lambda = 3$ and $\text{CFL} = 0.2$. The θ parameter used for this convergence test, are the same of the scalar one: for \mathbb{B}^1 we used $\theta_1 = 1$, for \mathbb{B}^2 we used $\theta_1 = 1, \theta_2 = 0$ and for \mathbb{B}^3 we used $\theta_1 = 1, \theta_2 = 5$. Also here, we need a bit more of corrections for \mathbb{B}^3 to reach the 4th order of accuracy ($K \approx 7$). As we can see in picture 8, the order of convergence

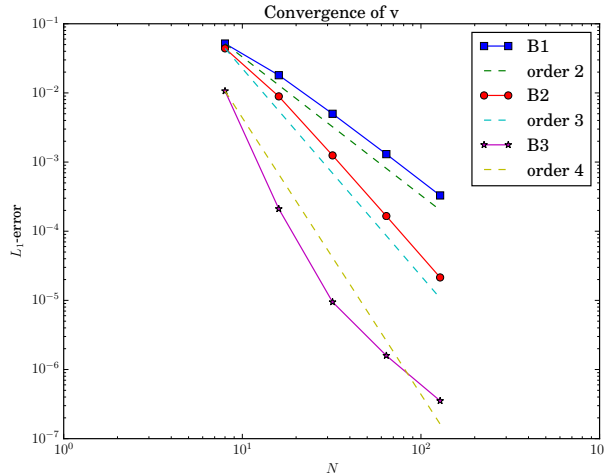


Figure 8: Convergence of Euler system

is what we expected.

5.1.4 Euler equation – Sod shock test

Now we can start testing the scheme on not smooth solutions. Let us begin with the Euler Sod test case. The Sod test case is solving equation (45) on domain $[0, 1]$, with EOS $E = \frac{p}{\gamma - 1} + \frac{1}{2}\rho v^2$, where $\gamma = 1.4$. The initial conditions are the following

$$\begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \text{ for } x \leq 0.5 \quad \begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 0.125 \\ 0 \\ 0.1 \end{pmatrix} \text{ for } x > 0.5,$$

final time is $T = 0.16$ and we have outflow boundary conditions.

In figure 9, we can see what we obtained for $\varepsilon = 10^{-9}$ in the formulation of IMEX Kinetic scheme (58). We used convection coefficient $\lambda = 2$, CFL = 0.2. For \mathbb{B}^1 $\theta_1 = 1$, for \mathbb{B}^2 $\theta_1 = 1, \theta_2 = 0.5$, for \mathbb{B}^3 $\theta_1 = 2.5, \theta_2 = 4$. In picture 9 we show the density plots for different mesh sizes $N = 64, 256$.

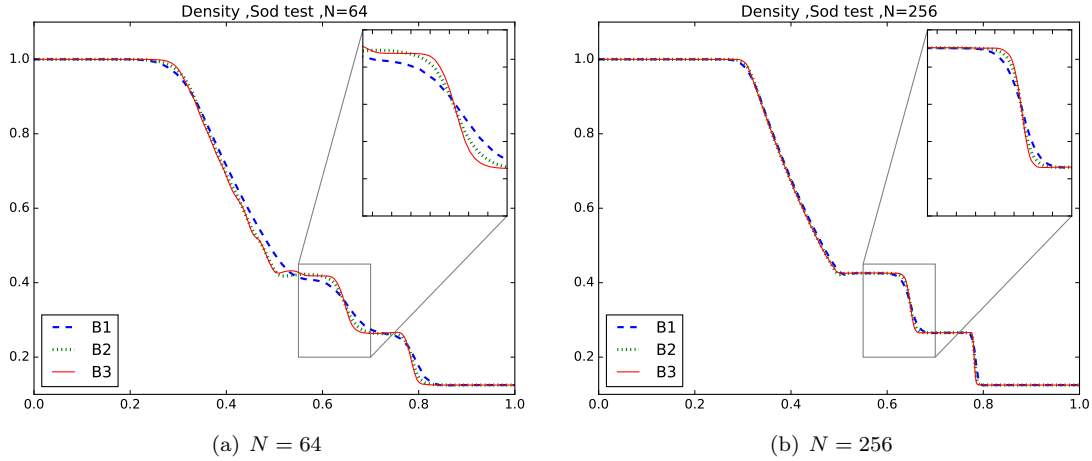


Figure 9: Density of Sod test case 1D

As we can see, even with few points the \mathbb{B}^3 solution is outperforming the other solutions, catching in a better way the edges of the discontinuities.

5.1.5 Euler equation – Woodward Colella

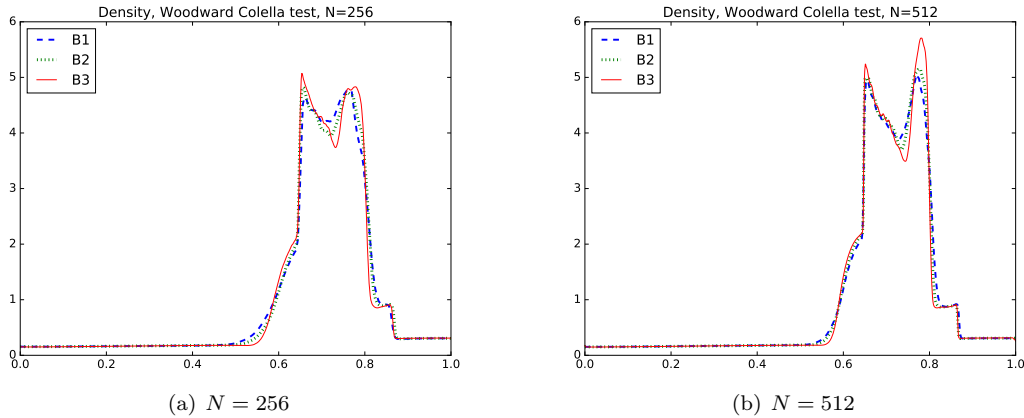


Figure 10: Density of Woodward Colella test

We can see even better the advantages of using a high order scheme in the following examples. First, we present the one proposed by Woodward and Colella [12]. It solves again Euler equation (45) on domain $[0, 1]$ with EOS (46) with $\gamma = 1.4$. The final time is 0.038, the initial conditions

are

$$\rho_0 = 1, \quad v_0 = 0, \quad p_0 = \begin{cases} 10^3 & \text{for } x \in [0, 0.1], \\ 10^{-2} & \text{for } x \in [0.1, 0.9], \\ 10^2 & \text{for } x \in [0.9, 1] \end{cases}$$

and we use outflow boundary conditions. For \mathbb{B}^1 , $\theta_1 = 0.5$. For \mathbb{B}^2 , $\theta_1 = 0.8$, $\theta_2 = 1$. For \mathbb{B}^3 , $\theta_1 = 5$, $\theta_2 = 1$. In figure 10 there is the result for $\varepsilon = 10^{-9}$, convection coefficient = 20, CFL = 0.1, $N = 256, 512$.

We can notice that in this case, only \mathbb{B}^3 is able to catch the shape of the second peak (with 512 elements).

5.1.6 Euler equation – Shu Osher test

Last test we performed in 1D was proposed by Shu and Osher [21]. Again we have Euler equation (45) on domain $[-5, 5]$ with EOS (46) with $\gamma = 1.4$. Here initial conditions are

$$\begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 3.857143 \\ 2.629369 \\ 10.333333 \end{pmatrix} \text{ if } x \in [-5, -4], \quad \begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 + 0.2 \sin(5x) \\ 0 \\ 1 \end{pmatrix} \text{ if } x \in [-4, 5].$$

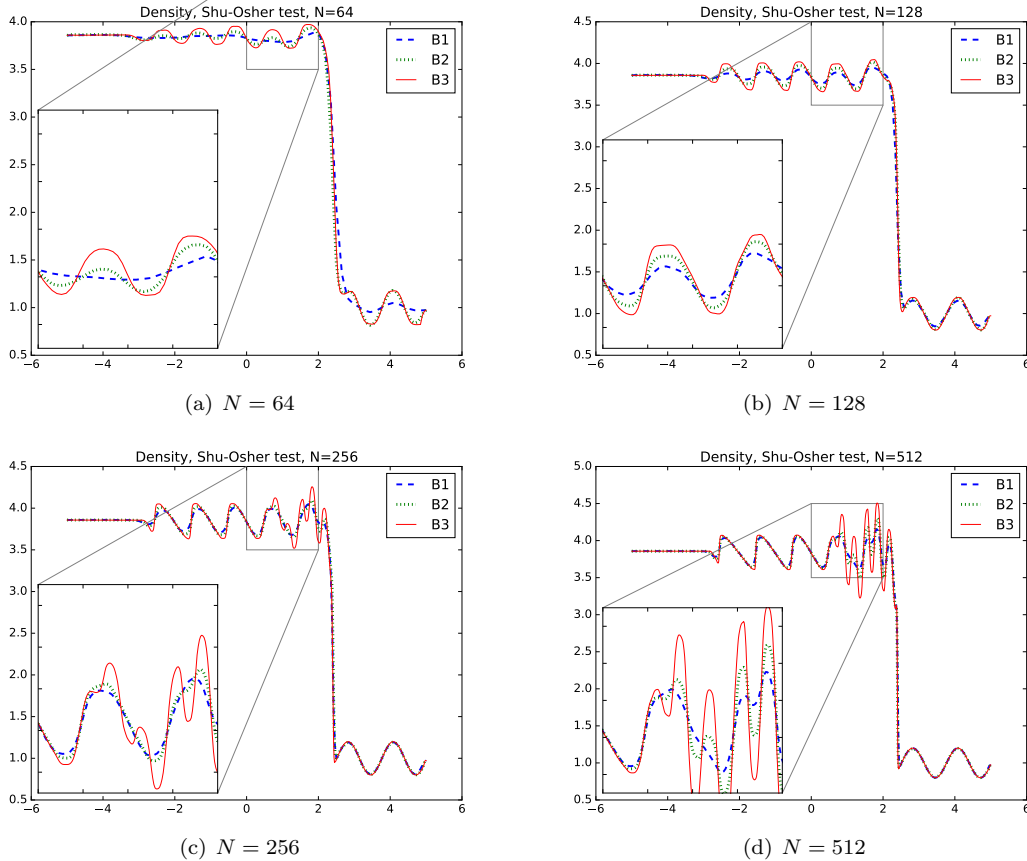


Figure 11: Density of Shu–Osher’s test

Final time is $T = 1.8$, we use outflow boundary conditions, $\varepsilon = 10^{-9}$, convection coefficient $\lambda = 3$, CFL = 0.1. For \mathbb{B}^1 $\theta_1 = 0.5$, for \mathbb{B}^2 $\theta_1 = 0.8$, $\theta_2 = 1$, for \mathbb{B}^3 $\theta_1 = 3$, $\theta_2 = 1$. In figure 11, we can see results for several N s. Even here, we can see that the second and third order polynomials perform better with respect to the first order one. In particular, we can see how the oscillations are already captured with few points and how the precision increases quickly if the order is greater.

In all these cases, we have seen that our method performs nicely and capture the correct behaviours of the equations solutions. Moreover, we see that it can be convenient to switch to higher order to better get the solution of our test cases with less mesh elements.

5.2 2D numerical tests

Let us present some numerical test defined on a 2D domain. We use again the DRM model proposed by [8] and the scheme we presented. We see only examples of Euler equation in 2D:

$$\partial_t U(\mathbf{x}, t) + \partial_x A_1(U(\mathbf{x}, t)) + \partial_y A_2(U(\mathbf{x}, t)) = 0, \quad \mathbf{x} = (x, y) \in \Omega \subset \mathbb{R}^2,$$

$$U = \begin{pmatrix} \rho \\ \rho u \\ \rho v \\ E \end{pmatrix}, \quad A_1(U) = \begin{pmatrix} \rho u \\ \rho u^2 + p \\ \rho uv \\ u(E + p) \end{pmatrix}, \quad A_2(U) = \begin{pmatrix} \rho v \\ \rho uv \\ \rho v^2 + p \\ v(E + p) \end{pmatrix} \quad (47)$$

where ρ is the density, u is the speed in x direction, v is the speed in y direction, E the total energy and p the pressure. They are linked by the following EOS:

$$p = (\gamma - 1) \left(E - \frac{1}{2} \rho (u^2 + v^2) \right). \quad (48)$$

5.2.1 Euler equation – Smooth vortex test case

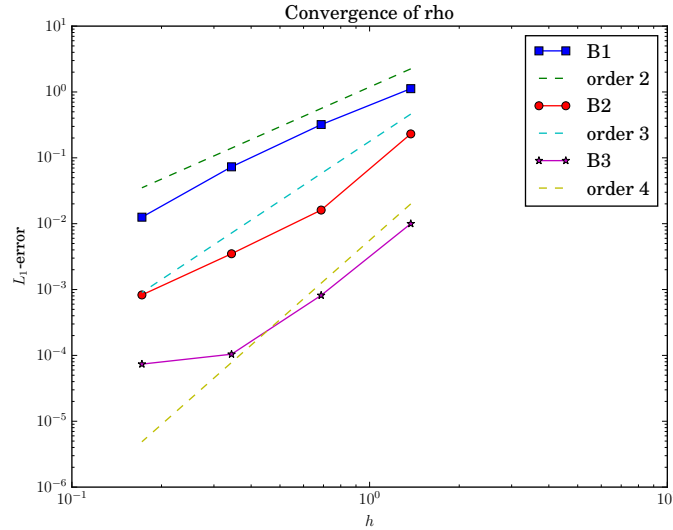


Figure 12: 2D convergence

To start, we want to study the convergence of the method also in 2D. To do so, we test our scheme with a steady vortex test case, so that we can compare the final solution with the initial

one. The domain is a circle of radius 10 and center $(0, 0)$. The initial conditions are

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} \left(1 - \frac{\gamma-1}{\gamma} \frac{1}{2} \left(\frac{5}{2\pi}\right)^2 e^{\frac{1-r^2}{2}}\right)^{\frac{1}{\gamma-1}} \\ \frac{5}{2\pi}(-y)e^{\frac{1-r^2}{2}} \\ \frac{5}{2\pi}(x)e^{\frac{1-r^2}{2}} \\ \rho_0^\gamma \end{pmatrix}.$$

Here $r^2 = x^2 + y^2$ and the boundary conditions are outflow. In our simulations $\gamma = 1.4$ for the EOS (48). Again, we take $\varepsilon = 10^{-9}$, convection coefficient $\lambda = 1.4$ and CFL = 0.1. We stop the simulation at time $T = 1$. We use different refinements of the domain mesh. These are uniform triangular meshes and on the x-axis of figure 12 one can see the maximum diameter of a cell of the mesh. We can see in figure 12 that the convergence is reflecting the theoretical results, even if for \mathbb{B}^3 we need more corrections ($K \approx 7$) to get the order to get closer to the convergence expected. For $\mathbb{B}^1 \theta_1 = 0.1$, for $\mathbb{B}^2 \theta_1 = 0.01$, $\theta_2 = 0$, for $\mathbb{B}^3 \theta_1 = 0.001$, $\theta_2 = 0$.

5.2.2 Euler equation – Sod 2D test case

We tested our method on the analogous of Sod in 2D. This test is again solving Euler equation (47) where $\gamma = 1.4$ in EOS (48). The domain Ω is a circle of radius 1 and center in $(0, 0)$. The initial conditions are:

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} \text{ if } x^2 + y^2 < \frac{1}{4}, \quad \begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 0.125 \\ 0 \\ 0 \\ 0.1 \end{pmatrix} \text{ if } x^2 + y^2 \geq \frac{1}{4}.$$

The parameters used for this test are $\varepsilon = 10^{-9}$, convection coefficient $\lambda = 1.4$, CFL = 0.1, final time $T = 0.25$ and outflow boundary conditions. For $\mathbb{B}^1 \theta_1 = 0.1$, for $\mathbb{B}^2 \theta_1 = 0.1$, $\theta_2 = 0.0001$, for $\mathbb{B}^3 \theta_1 = 0.01$, $\theta_2 = 0.0001$.

We use uniform triangular meshes with different sizes and what is shown in pictures 13 is obtained with $N = 3576$ and $N = 13548$ triangles on the domain.

If we watch pictures 13 and 14, we can see that also in this case the higher the order of polynomial we use, the sharper becomes the solution. In particular, we can say that the solution with \mathbb{B}^2 basis functions for the mesh with $N = 13548$ elements is comparable with the solution for \mathbb{B}^3 with only $N = 3576$ elements. Moreover, we can see that with \mathbb{B}^1 the diffusion is too high and it is smoothening all the discontinuities.

5.2.3 Euler equation – DMR 2D test case

For the last test case, we test our scheme on the DMR (double Mach reflection) problem presented in [15]. The equation we are solving is again the Euler equation (47) with $\gamma = 1.4$ in EOS (48). The domain is a rectangular shape, cut on the bottom right part by an oblique edge. The boundaries of the rectangle are $x = 0$, $x = 2.2$, $y = -0.2$, $y = 3$. The oblique edge is a line passing through points $(0, 0)$ and $(3, 1.7)$. We have wall boundary conditions on the bottom, on the top and on the oblique edge of the mesh, inflow on the left edge and outflow on the right one. The initial conditions are a shock, which divides high density (left-side $x \leq 0$) and low density (right-side $x < 0$). This shock has an initial speed in right direction. As the time passes, the shock crosses the oblique surface and creates more internal shock surfaces. The initial conditions

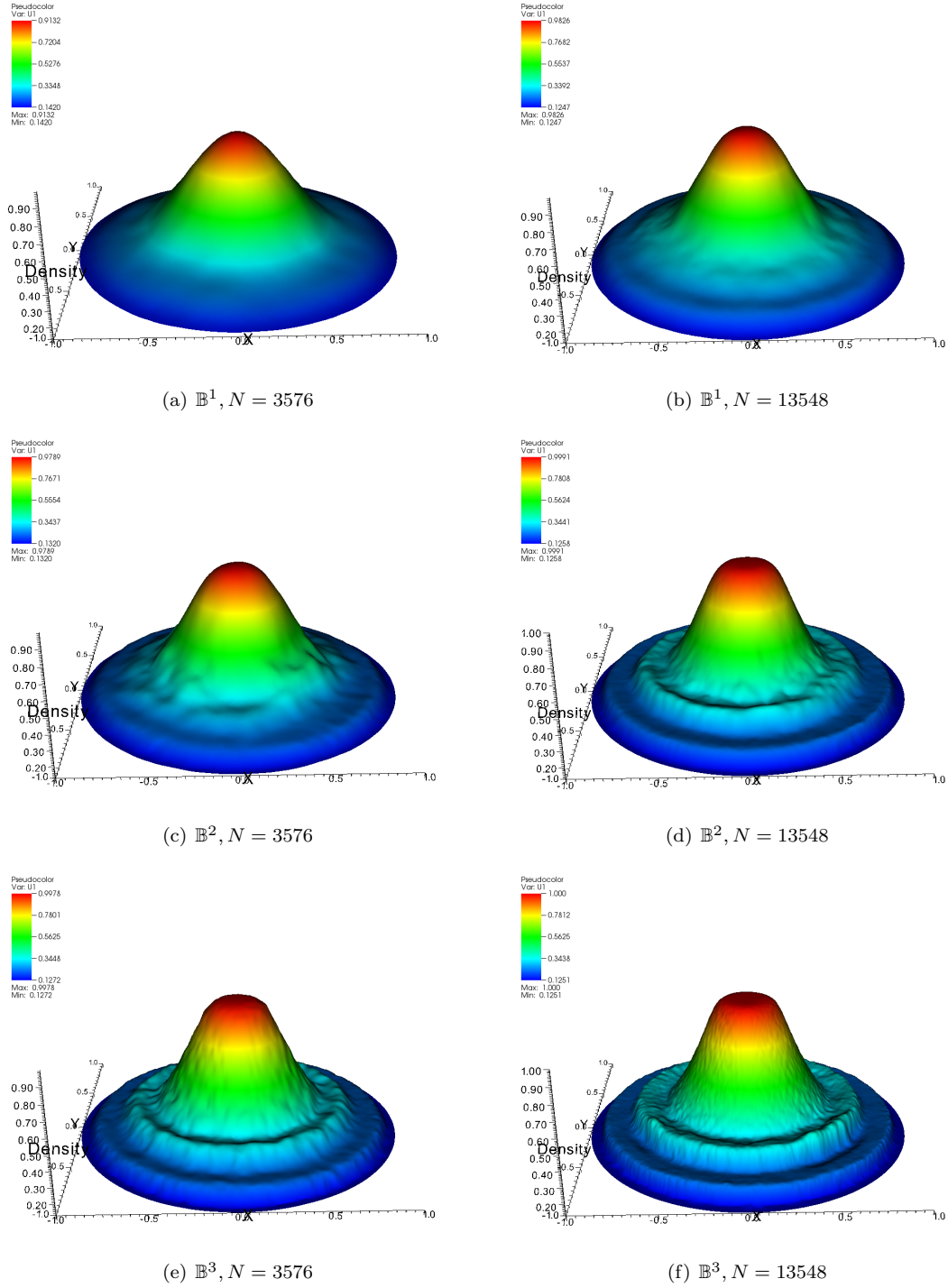
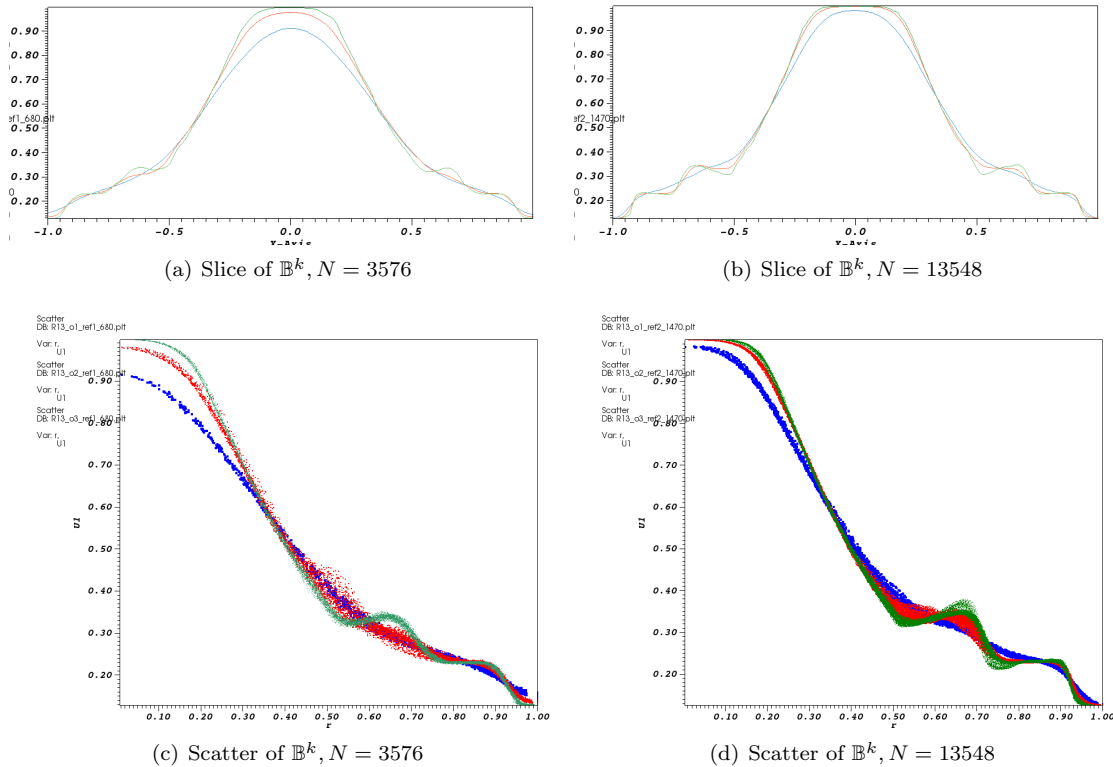


Figure 13: Density of Sod test

are more precisely the following

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 8 \\ 8.25 \\ 0 \\ 116.5 \end{pmatrix} \text{ if } x \leq 0, \quad \begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1.4 \\ 0 \\ 0 \\ 1 \end{pmatrix} \text{ if } x > 0.$$

Figure 14: Density of Sod test (\mathbb{B}^1 blue, \mathbb{B}^2 red and \mathbb{B}^3 green)

The parameters used for this test are $\varepsilon = 10^{-9}$, convection coefficient $\lambda = 15$, CFL = 0.1, final time $T = 0.2$. The mesh we used is composed of $N = 19248$ triangular elements with a maximum diameter of 0.0369. For \mathbb{B}^1 $\theta_1 = 0.1$, for \mathbb{B}^2 $\theta_1 = 0.01$, $\theta_2 = 0.0001$, for \mathbb{B}^3 $\theta_1 = 0.005$, $\theta_2 = 0.0001$.

Again we can see in pictures 15, 16 and 17 that the scheme catches the behaviour of the shock and its reflection against the lower wall. Even now, we can see that the sharpness of the shock is really well captured by the \mathbb{B}^3 scheme, while the others are less precise in defining the shock zone.

6 Conclusions and further investigations

We have presented a residual distribution scheme for hyperbolic system of equations with stiff relaxation source terms for kinetic models. The method proposed takes advantage of the IMEX formulation (implicit for source term and explicit for advection term) to resolve the stiffness of the relaxation source. Nevertheless, we were able to solve computationally explicitly the kinetic model of [8], thanks to an auxiliary equation, which allows us not to recur to nonlinear solver. The high accuracy of the scheme is reached thanks to two ingredients. The first one is the residual distribution framework for spatial discretization [3], which is a finite element based method that is naturally high order because of the choice of different basis functions. The second is the high order time-integration performed in the DeC method, which allows to couple two schemes, the first easy to solve, for us the IMEX scheme, and a second high order scheme, the high order time-integration residual distribution scheme. The result is an iterative method able to reach high order and stability via few iterations. This is the first time, as far as we know, that the residual

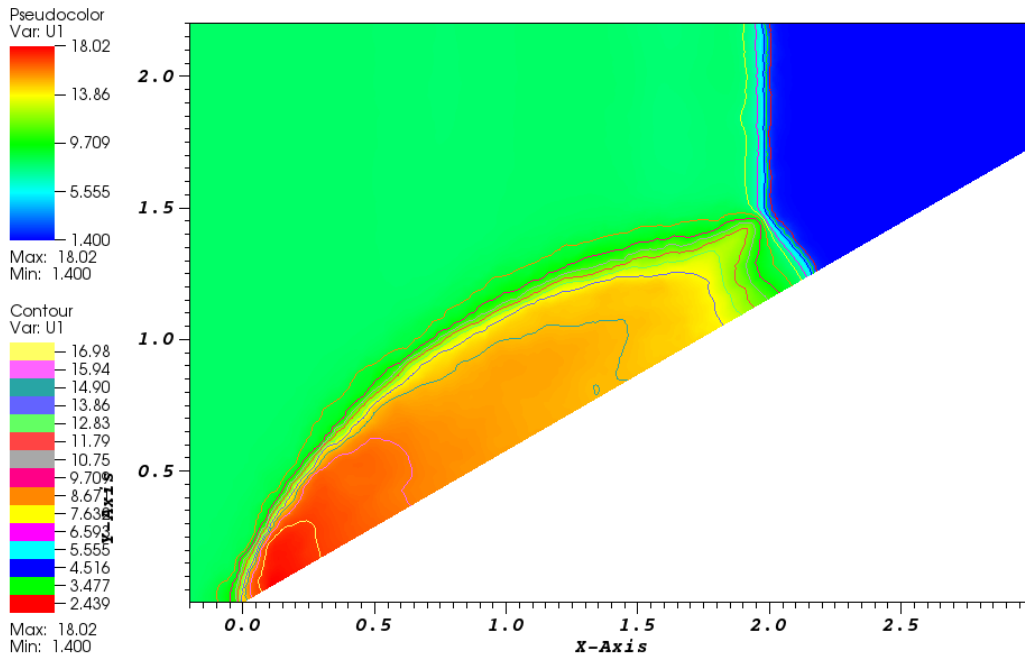


Figure 15: Density of DMR test \mathbb{B}^1

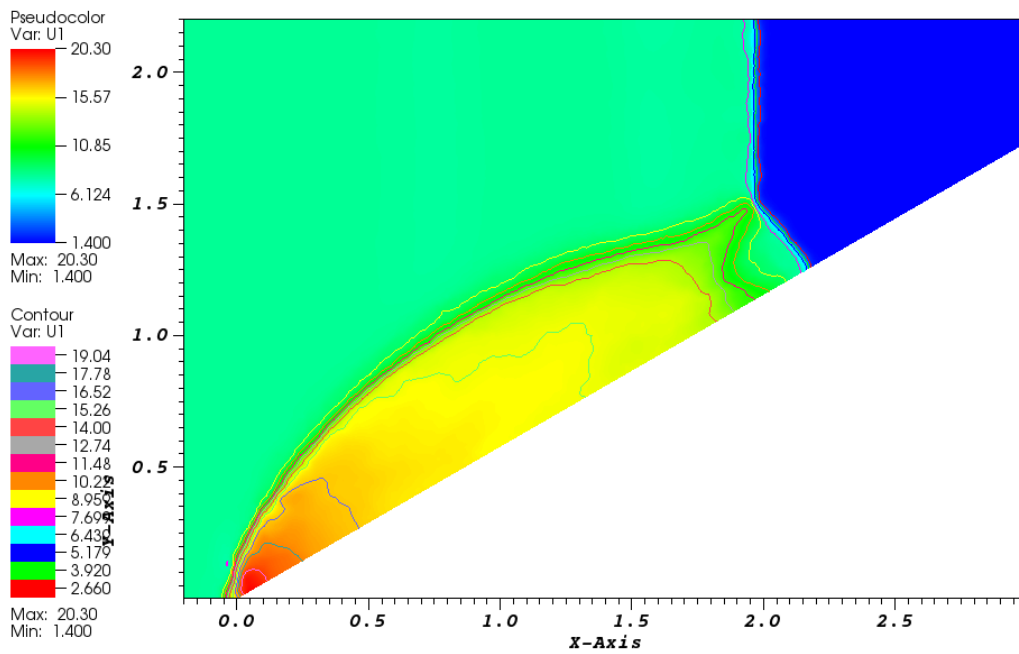
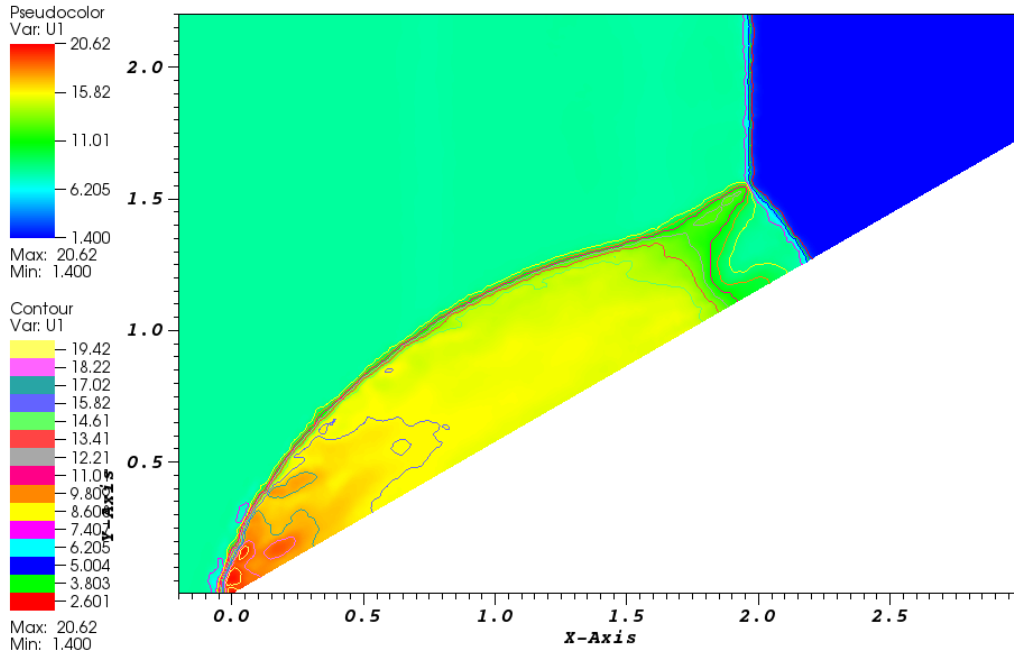


Figure 16: Density of DMR test \mathbb{B}^2

Figure 17: Density of DMR test \mathbb{B}^3

distribution framework is used to solve hyperbolic systems with stiff source terms. Even if in this work we solved only one model, it is easily extensible to different models which present similar properties.

The results obtained both from a theoretical point of view and from the simulation side are satisfactory. Indeed, the theorems proved the asymptotic preserving property for our scheme and the rate of accuracy. In addition, the run simulations are reaching the expected accuracy in 1D and 2D, the correct behaviour of the discontinuities of the solutions is well caught by the scheme and as the order increases we can see big improvements in shapes of solutions.

Further investigations may be in the following directions. There are still some open questions over the complete automation of the scheme. For example, it is still not well known which is the relation between parameters $\theta_1, \theta_2, \Delta t$ and the quality of the solution. There are studies for 1D smooth solutions, where some relations between these quantities are shown, thanks to some von Neumann stability analysis [22]. Nonetheless, these results are not easily extensible to nonlinear flux problems or 2D problems.

Moreover, it is not clear why for \mathbb{B}^3 the scheme needs more corrections than expected to reach the order of convergence, in particular when the mesh is more refined. This is a contradiction of proposition (3.1) as shown in [5].

Finally, we are already working on some extensions of the scheme for multiphase flows equations and we believe that it can be applied also for a large variety of other problems, such as BGK equations, viscoelasticity problems or other kinetic schemes.

Acknowledgments

We acknowledge the support of ITN ModCompShock project funded by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement

No 642768. We acknowledge Paola Bacigaluppi and Svetlana Tokareva for their contributions in coding and discussing the residual distribution formulation.

A Residual Distribution schemes

The key point of the RD schemes is the definition of the splitting of the total residuals into nodal residuals. Through this definition one can actually define the proper scheme to utilise. One can rewrite, for example, the SUPG scheme [16] in this way:

$$\phi_\sigma^K(U_h) = \int_K \varphi_\sigma(\nabla A(U_h) - S(U_h))dx + h_K \int_K (\nabla A(U_h) \cdot \nabla \varphi_\sigma) \tau(\nabla A(U_h) \cdot \nabla U_h). \quad (49)$$

What we use in our code are two types of residuals: one for smooth test cases, one for shock test cases.

A.1 Smooth solutions residuals

When we are dealing with smooth tests and we know a priori that we do not need the extra diffusion to damp oscillations brought by discontinuities, we can use a pure Galerkin discretization with a stabilization of jump of the gradient of the solution [11, 4]. The study of the stability of the scheme in this situation for smooth solutions is shown in [22], through a von Neumann analysis of the scheme.

For an hyperbolic system of equation with source term

$$\partial_t U + \nabla \cdot A(U) - S(U) = 0, \quad (50)$$

the scheme proceeds as follows $\forall \sigma \in \Sigma$

$$\phi_\sigma^{K,1}(U_h) = \int_{\partial K} \varphi_\sigma A(U_h) \cdot \mathbf{n} d\Gamma - \int_K \nabla \varphi_\sigma \cdot A(U_h) d\mathbf{x} - \int_K \varphi_\sigma S(U_h) d\mathbf{x}, \quad (51)$$

and then

$$\phi_\sigma^K = \phi_\sigma^{K,1} + \sum_{k=1}^d \sum_{e \in \text{edge of } K} \theta_k h_e^{2k} \int_e [\nabla^k U_h] \cdot [\nabla^k \varphi_\sigma] d\Gamma. \quad (52)$$

Here d is the degree of the polynomial of the basis functions we use, θ_k are positive coefficients and $[\cdot]$ is the jump across the edge e , namely, if e separates K and K^+ , $[u] = u|_K - u|_{K^+}$. All the derivatives are meant in the direction of the normal to the edge e and h_e is the length of a 1D element of the mesh (the edge e in 2D, the size of a cell $|K|$ in 1D). The schemes just presented are naturally of order $d + 1$ where d is the degree of the polynomial that we are using for the discretization. The parameters θ_k must be chosen carefully if we want the scheme to be stable. The stability analysis of this scheme in [22] suggests some optimal values for these parameters in case of 1D linear fluxes. It is not easy to extend this study to different test cases. In addition, these schemes are not too dissipative and they preserve the order of accuracy. Anyway, they do not guarantee stability in case of shocks and discontinuities.

A.2 Shock solutions residuals

Now, we present the schemes that is used in our simulations in presence of discontinuities or not smooth solutions. More details of these schemes are shown in [6]. The procedure starts defining a local Galerkin Lax–Friedrichs type nodal residual on the steady part of original equation (50):

$$\phi_\sigma^{K,LxF}(U_h) := \int_{\partial K} \varphi_\sigma A(U_h) \cdot \mathbf{n} d\Gamma - \int_K \nabla \varphi_\sigma \cdot A(U_h) dx - \int_K \varphi_\sigma S(U_h) dx + \alpha_K (U_\sigma - \bar{U}_h^K), \quad (53)$$

where \bar{U}_h^K is the average of U_h over the cell K and α_K is defined as

$$\alpha_K := \max_{\sigma \in K} \left(\rho_S \left(\nabla A(\bar{U}_h^K) \cdot \nabla \varphi_\sigma \right) \right), \quad (54)$$

and ρ_S is the function returning the spectral radius of the input matrix. Then, to guarantee monotonicity of the solution near strong discontinuities, we proceed as follows:

$$\begin{aligned} \beta_\sigma^K(U_h) &:= \max \left(\frac{\Phi_\sigma^{K,LxF}}{\Phi^K}, 0 \right) \left(\sum_{j \in K} \max \left(\frac{\Phi_j^{K,LxF}}{\Phi^K}, 0 \right) \right)^{-1}, \\ \phi_\sigma^{*,K} &:= \beta_\sigma^K \phi^K. \end{aligned} \quad (55)$$

These divisions between vectors are meant component-wise in characteristic variables, that implies the computation of the right eigenvectors of the multiplication of the jacobian of the flux and the normal average velocity $\nabla A(U_h) \cdot \mathbf{n}$. Then, we do a blending between this new residual and the Lax-Friedrichs's one. We use a coefficient Θ defined as

$$\Theta := \frac{|\Phi^K|}{\sum_{j \in K} |\Phi_j^{K,LxF}|} \quad (56)$$

and the new residual is

$$\phi_\sigma^{:,K} := (1 - \Theta) \phi_\sigma^{*,K} + \Theta \Phi_\sigma^{K,LxF}. \quad (57)$$

This scheme guarantees the monotonicity principle [3].

After that we add to the scheme the jump stabilization terms

$$\phi_\sigma^K := \phi_\sigma^{:,K} + \sum_{k=1}^d \sum_{e \in \text{edge of } K} \theta_k h_e^{2k} \int_e [\nabla^k U_h] \cdot [\nabla^k \varphi_\sigma] d\Gamma, \quad (58)$$

and this defines the final scheme.

B Deferred Correction properties

B.1 Proof of DeC theorem

Proposition B.1. *Let \mathcal{L}^1 and \mathcal{L}^2 be two operators defined on \mathbb{R}^m , which depend on the parameter Δ , such that*

- \mathcal{L}^1 is coercive for one norm, i.e., $\exists \alpha_1 > 0$ independent of Δ , can be both Δx or Δt since they are linked by CFL conditions, such that for any f, g we have that

$$\alpha_1 \|f - g\| \leq \|\mathcal{L}^1(f) - \mathcal{L}^1(g)\|$$

- $\mathcal{L}^1 - \mathcal{L}^2$ is Lipschitz with constant $\alpha_2 > 0$ uniformly with respect to Δ , i.e., for any U, V

$$\|(\mathcal{L}^1(f) - \mathcal{L}^2(f)) - (\mathcal{L}^1(g) - \mathcal{L}^2(g))\| \leq \alpha_2 \Delta \|f - g\|.$$

We also assume that $\exists! f_\Delta^*$ such that $\mathcal{L}^2(f_\Delta^*) = 0$. Then, if $\eta = \frac{\alpha_2}{\alpha_1} \Delta < 1$, the deferred correction is converging to f^* and after k iterations the error is smaller than η^k .

Proof. Let f^* be the solution of $\mathcal{L}^2(f^*) = 0$. Here, we drop the dependency on f^n in \mathcal{L}^1 , \mathcal{L}^2 , for simplicity. We know that $\mathcal{L}^1(f^*) = \mathcal{L}^1(f^*) - \mathcal{L}^2(f^*)$, so that

$$\mathcal{L}^1(f^{(k+1)}) - \mathcal{L}^1(f^*) = \left(\mathcal{L}^1(f^{(k)}) - \mathcal{L}^1(f^*) \right) - \left(\mathcal{L}^2(f^{(k)}) - \mathcal{L}^2(f^*) \right), \quad (59)$$

then

$$\alpha_1 \|f^{(k+1)} - f^*\| \leq \| \mathcal{L}^1(f^{(k+1)}) - \mathcal{L}^1(f^*) \| = \quad (60)$$

$$= \| \mathcal{L}^1(f^{(k)}) - \mathcal{L}^2(f^{(k)}) - (\mathcal{L}^1(f^*) - \mathcal{L}^2(f^*)) \| \leq \quad (61)$$

$$\leq \alpha_2 \Delta \|f^{(k)} - f^*\|. \quad (62)$$

Hence, we can write

$$\|f^{(k+1)} - f^*\| \leq \left(\frac{\alpha_2}{\alpha_1} \Delta \right) \|f^{(k)} - f^*\| \leq \left(\frac{\alpha_2}{\alpha_1} \Delta \right)^{k+1} \|f^{(0)} - f^*\|. \quad (63)$$

After k iterations we have an error at most of $\eta^k \cdot \|f^{(0)} - f^*\|$. \square

B.2 Lipschitz continuity and coercivity

Let us prove that our \mathcal{L}^1 and \mathcal{L}^2 schemes verify all the hypothesis of proposition (B.1).

Proposition B.2. \mathcal{L}^1 is coercive, i.e., $\exists \alpha_1 > 0$ s.t. $\forall f, g \in V_h$ and $m = 1, \dots, M$

$$\| \mathcal{L}_u^{1,m}(f^0, Pf) - \mathcal{L}_u^{1,m}(f^0, Pg) \| \geq \alpha_1 \|Pf - Pg\|, \quad (64)$$

$$\| \mathcal{L}^{1,m}(f^0, f) - \mathcal{L}^{1,m}(f^0, g) \| \geq \alpha_1 \|f - g\|. \quad (65)$$

Proof. The u part is trivial because

$$\mathcal{L}_{\sigma,u}^{1,m}(f^0, Pf) - \mathcal{L}_{\sigma,u}^{1,m}(f^0, Pg) = Pf_\sigma^m - Pg_\sigma^m. \quad (66)$$

For f part, we have to collect the implicit terms as done in (33c). Then, we can write

$$\mathcal{L}_\sigma^{1,m}(f^0, f) - \mathcal{L}_\sigma^{1,m}(f^0, g) = (f_\sigma^m - g_\sigma^m) - \frac{\Delta t}{\Delta t + \varepsilon} (M(Pf_\sigma^m) - M(Pg_\sigma^m)) = f_\sigma^m - g_\sigma^m. \quad (67)$$

The last step is possible, since the Maxwellians in our scheme are computed from the u equation and they are actually explicit, so they must coincide. If we write the operator explicitly both for u and f , we can see that the coercivity constant $\alpha_1 = 1$, given any norm. \square

Before proving the Lipschitz continuity, we have to introduce some norms. We use the following definition of norm for a function $f \in V_h$, which is consistent with the \mathcal{L}^2 norm,

$$\|f\|^2 = \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| f_\sigma^2. \quad (68)$$

We also define the norm of all the subimesteps as

$$\|\mathbf{f}\| = \| (f^0, \dots, f^M) \| = \sqrt{\sum_{m=1}^M \|f^m\|^2}. \quad (69)$$

Moreover, we will need the definition of the following seminorms

$$|f|_{1,x}^2 := \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| \left(\max_{K|\sigma \in K} \max_{x \in K} \frac{f_\sigma - f(x)}{d(K)} \right)^2, \quad (70)$$

$$|\mathbf{f}|_{1,t}^2 := \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| \left(\max_{m=1,\dots,M} \frac{f^m - f^{m-1}}{\Delta t^m} \right)^2, \quad (71)$$

where $d(K)$ is the diameter of the cell K and it is bounded by $\max_K d(K) = h$. In particular, we note that $|f|_{1,x} \leq |f|_1 = \|\nabla f\|_{L^2}$ for every discretization mesh.

Proposition B.3. *If we assume that*

$$|f|_{1,x} \leq C_1 \|f\|, \quad (72)$$

$$|\mathbf{f}|_{1,t} \leq C_2 \|\mathbf{f}\|, \quad (73)$$

where C_1 and C_2 do not depend on the mesh size h and timestep Δt . And if we require that nodal residuals verify

$$\sum_{\sigma \in D_h} \frac{1}{|\mathcal{C}_\sigma|} \left(\sum_{K|\sigma \in K} \phi_\sigma^K(f) - \phi_\sigma^K(g) \right)^2 \leq C_3 \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| (f_\sigma - g_\sigma)^2 = C_3 \|f - g\|^2. \quad (74)$$

Then, $\mathcal{L}^1 - \mathcal{L}^2$ is Lipschitz continuous, i.e., $\exists \alpha_2 > 0$ s.t. $\forall f, g \in V_h$

$$\| |(\mathcal{L}_u^1(P\mathbf{f}) - \mathcal{L}_u^1(P\mathbf{g})) - (\mathcal{L}_u^2(P\mathbf{f}) - \mathcal{L}_u^2(P\mathbf{g}))| \| \leq \alpha_2 \Delta \|P\mathbf{f} - P\mathbf{g}\|, \quad (75)$$

$$\| |(\mathcal{L}^1(\mathbf{f}) - \mathcal{L}^1(\mathbf{g})) - (\mathcal{L}^2(\mathbf{f}) - \mathcal{L}^2(\mathbf{g}))| \| \leq \alpha_2 \Delta \|\mathbf{f} - \mathbf{g}\|. \quad (76)$$

Remark *The extra hypothesis added are related to the regularity of the solution. Of course, this is not always the case, and, for example, when there are shocks in the solution, (72) does not hold. Anyway, even if we can not prove the convergence for those cases, we see numerically a big improvement in higher order solutions. The inequality (73) is actually given, during the DeC procedure, by the Lipschitz continuity of fluxes and residuals. To keep the proof more general, we add it as a further hypothesis. Equation (74), in our case, is given by the consistency of the nodal residuals, the Lipschitz continuity of the flux F and by the regularity of the solutions f, g as stated in (72).*

Proof. The estimation of (75) is a simplification of the case of (76), so we will skip its proof.

For simplicity, let us define the differences $\delta f := f - g$, $\delta \phi_\sigma^K(f) := \phi_\sigma^K(f) - \phi_\sigma^K(g)$, $\delta M(Pf) := M(Pf) - M(Pg)$, $\delta \mathcal{L} := \mathcal{L}^1 - \mathcal{L}^2$ and $\delta \mathcal{I}(\mathbf{f}) := \mathcal{I}_0(\mathbf{f}) - \mathcal{I}_M(\mathbf{f})$.

Let us split the operators into two parts. The first one is composed of the term related to time derivative and source term \mathcal{L}_{ts} , the second one concerns the advection part \mathcal{L}_{ad} . If we write explicitly the source and time part, we get

$$\begin{aligned} & \delta \mathcal{L}_{ts,\sigma}^m(f) - \delta \mathcal{L}_{ts,\sigma}^m(g) = \\ &= \sum_{K|\sigma \in K} \frac{1}{|\mathcal{C}_\sigma|} \left[\frac{\varepsilon}{\varepsilon + \Delta t^m} \int_K \varphi_\sigma (\delta f_\sigma^m - \delta f^m) - \frac{\Delta t}{\varepsilon} \int_K \varphi_\sigma (\delta M(Pf_\sigma^m) - \delta f_\sigma^m) + \right. \\ & \quad \left. + \frac{\varepsilon}{\varepsilon + \Delta t^m} \frac{1}{\varepsilon} \int_{t^0}^{t^m} \mathcal{I}_M (\delta \phi_{s,\sigma}^K(f^0), \dots, \delta \phi_{s,\sigma}^K(f^M), s) ds \right]. \end{aligned} \quad (77a)$$

Now, let us suppose that the residuals are a consistent discretization of fluxes and source terms, so let us use the Galerkin discretization instead of any other one. Moreover, let us add and subtract the residual in timestep $t^{n,m}$. So, we can write, neglecting $\mathcal{O}(\Delta^2)$,

$$\mathcal{L}_{ts,\sigma}^{1,m}(f) - \mathcal{L}_{ts,\sigma}^{1,m}(g) - \mathcal{L}_{ts,\sigma}^{2,m}(f) + \mathcal{L}_{ts,\sigma}^{2,m}(g) + \mathcal{O}(\Delta^2) = \quad (78a)$$

$$\begin{aligned} &= \frac{1}{|\mathcal{C}_\sigma|} \int_\Omega \varphi_\sigma (\delta f_\sigma^m - \delta f^m) - \frac{1}{|\mathcal{C}_\sigma|} \frac{\Delta t^m}{\varepsilon + \Delta t^m} \int_\Omega \varphi_\sigma (\delta M(Pf_\sigma^m) - \delta M(Pf^m)) + \\ &+ \frac{1}{\varepsilon + \Delta t^m} \int_{t^0}^{t^m} \mathcal{I}_M(\delta \phi_{s,\sigma}^K(f^0) - \delta \phi_{s,\sigma}^K(f^m), \dots, \delta \phi_{s,\sigma}^K(f^M) - \delta \phi_{s,\sigma}^K(f^m), s) ds. \end{aligned} \quad (78b)$$

Now, we sum over the DoFs and we square the previous quantity. We use Lemma A.1 of [4] to pass from coefficients v_σ to pointwise evaluation $v(\sigma)$, with the abuse of notation. It states that $\sum_{\sigma \in K} |v_\sigma - v_{\sigma'}| \leq C_K \sum_{\sigma \in K} |v(\sigma) - v(\sigma')|$ where C_K is the norm of the inverse of the matrix $(\varphi_\sigma(\sigma'))_{\sigma,\sigma'}$ and it depends on K only via the aspect ratio of the element K .

$$\sum_{\sigma \in D_h} |\mathcal{C}_\sigma| \left(\mathcal{L}_{ts,\sigma}^{1,m}(f) - \mathcal{L}_{ts,\sigma}^{1,m}(g) - \mathcal{L}_{ts,\sigma}^{2,m}(f) + \mathcal{L}_{ts,\sigma}^{2,m}(g) \right)^2 \leq \quad (79a)$$

$$\begin{aligned} &\leq C_a h^2 \sum_{\sigma \in D_h} \frac{1}{|\mathcal{C}_\sigma|} \left(\int_\Omega \varphi_\sigma \left(\frac{\delta f_\sigma^m - \delta f^m(x)}{d(K)} \right) \right)^2 + \\ &+ C_b h^2 \frac{\Delta t^m}{(\varepsilon + \Delta t^m)} \sum_{\sigma \in D_h} \frac{1}{|\mathcal{C}_\sigma|} \left(\int_\Omega \varphi_\sigma \frac{\delta M(Pf^m)(\sigma) - \delta M(Pf^m)}{d(K)} \right)^2 + \end{aligned} \quad (79b)$$

$$\begin{aligned} &+ C_c \frac{\Delta t^m}{\varepsilon + \Delta t^m} \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| \max_r (\delta \phi_{s,\sigma}^K(f^r) - \delta \phi_{s,\sigma}^K(f^m))^2 \leq \\ &\leq C_d h^2 (|\delta f^m|_{1,x}^2 + |\delta M(Pf^m)|_{1,x}^2 + \max_r \|\delta f^r - \delta f^m\|) \leq \end{aligned} \quad (79c)$$

$$\leq C_e h^2 (\|\delta f^m\|^2 + \|\delta M(Pf^m)\|^2 + \Delta t^2 \|\delta \mathbf{f}\|_{1,t}) \leq \quad (79d)$$

$$\leq C_f h^2 \|\delta \mathbf{f}\|^2 + \mathcal{O}(h^4) \leq C_4 h^2 \|\mathbf{f} - \mathbf{g}\|^2. \quad (79e)$$

In (79b) we explicitly bring the scale h outside the first two sums, while in the third term we just bound the interpolant polynomial with the maximum of the interpolant values times a constant, in (79c) we use the definition of the seminorm (70), the Lipschitz continuity of residuals (74), the product rule for integrals and the bound $\Delta t^m \leq \Delta t^m + \varepsilon$. In (79d) we use the inequality (72) and the definition of the seminorm (71). In (79e) we use the fact that the Maxwellians M and the projections P are Lipschitz continuous, the inequality (73) and the fact that $\Delta t \sim h$. The constant C_4 does not depend on $h, \Delta t$ nor on ε . It depends on the size of the domain, on the Lipschitz continuity of the Maxwellians, on the regularity of the mesh and on basis functions.

For the advection term a similar computation is carried out, but, in this case the error is a $\mathcal{O}(\Delta t)$. Using the notation of $\phi_\sigma := \sum_{K|\sigma \in K} \phi_\sigma^K$, let us write

$$\|\mathcal{S}_x\|^2 := \sum_{\sigma \in D_h} |\mathcal{C}_\sigma| \left(\delta \mathcal{L}_{ad,\sigma}^{1,m}(f) - \delta \mathcal{L}_{ad,\sigma}^{1,m}(g) \right)^2 = \quad (80a)$$

$$= \sum_{\sigma \in D_h} \frac{1}{|\mathcal{C}_\sigma|} \left(\frac{\varepsilon}{\varepsilon + \Delta t^m} \int_{t^{n,0}}^{t^{n,m}} \delta \mathcal{I}(\delta \phi_{ad,\sigma}(f^0), \dots, \delta \phi_{ad,\sigma}(f^M), s) ds \right)^2 \leq \quad (80b)$$

$$\leq C_l \sum_{\sigma \in D_h} \frac{\Delta t^2}{|\mathcal{C}_\sigma|} \left(\sum_{K|\sigma \in K} \max_{m=1,\dots,M} \frac{|\delta \phi_{ad,\sigma}^K(f^m) - \delta \phi_{ad,\sigma}^K(f^{m-1})|}{\Delta t^m} \right)^2. \quad (80c)$$

In (80c) we use the bound $\varepsilon \leq \varepsilon + \Delta t^m$ and the fact that \mathcal{I}_0 is a zero order approximation of \mathcal{I}_M , so, adding the integration in time, we get the error estimation above.

$$\|\mathcal{S}_x\|^2 \leq C_q \sum_{\sigma \in D_h} \Delta t^2 |\mathcal{C}_\sigma| \left(\max_{m=1, \dots, M} \frac{|\delta f^m - \delta f^{m-1}|}{\Delta t^m} \right)^2 \leq \quad (80d)$$

$$\leq C_p \Delta t^2 \sum_{m=1}^M |f^m - g^m|_{1,t}^2 \leq C_5 \Delta t^2 \|\mathbf{f} - \mathbf{g}\|^2. \quad (80e)$$

In (80d) we use the Lipschitz continuity and consistent hypothesis over the residuals as stated in (74). Finally, in (80e) we use the definition of seminorm (71) and we apply the bound in (73). Again, C_5 does not depend on Δt , h or ε , but only on fluxes, geometry and basis functions.

In conclusion, summing up the inequalities (79e) and (80e), we prove the thesis of the proposition. \square

References

- [1] R. Abgrall. Toward the ultimate conservative scheme: Following the quest. *Journal of Computational Physics*, 167(2):277 – 315, 2001.
- [2] R. Abgrall. Essentially non-oscillatory residual distribution schemes for hyperbolic problems. *J. Comput. Phys.*, 214(2):773–808, 2006.
- [3] R. Abgrall. Residual distribution schemes: Current status and future trends. *Computers & Fluids*, 35(7):641 – 669, 2006. Special Issue Dedicated to Professor Stanley G. Rubin on the Occasion of his 65th Birthday.
- [4] R. Abgrall. High order schemes for hyperbolic problems using globally continuous approximation and avoiding mass matrices. *Journal of Scientific Computing*, 73(2):461–494, Dec 2017.
- [5] R. Abgrall. Some remarks about conservation for residual distribution schemes. *Computational Methods in Applied Mathematics*, 2018. DOI: <https://doi.org/10.1515/cmam-2017-0056>.
- [6] R. Abgrall, P. Bacigaluppi, and S. Tokareva. High-order residual distribution scheme for the time-dependent euler equations of fluid dynamics. *Computers & Mathematics with Applications*, 2018.
- [7] R. Abgrall, A. Larat, and M. Ricchiuto. Construction of very high order residual distribution schemes for steady inviscid flow problems on hybrid unstructured meshes. *J. Comput. Phys.*, 230(11):4103–4136, 2011.
- [8] D. Aregba-Driollet and R. Natalini. *Discrete Kinetic Schemes for Systems of Conservation Laws*, pages 1–10. Birkhäuser Basel, Basel, 1999.
- [9] D. Aregba-Driollet and R. Natalini. Discrete kinetic schemes for multidimensional systems of conservation laws. *SIAM J. Numer. Anal.*, 37(6):1973–2004, 2000.
- [10] S. Boscarino, J. Qiu, and G. Russo. Implicit-explicit integral deferred correction methods for stiff problems. 40, 01 2017.
- [11] E. Burman and P. Hansbo. Edge stabilization for galerkin approximations of convection–diffusion–reaction problems. *Computer Methods in Applied Mechanics and Engineering*, 193(15):1437 – 1453, 2004. Recent Advances in Stabilized and Multiscale Finite Element Methods.
- [12] P. Colella and P. R. Woodward. The Piecewise Parabolic Method (PPM) for Gas-Dynamical Simulations. *Journal of Computational Physics*, 54:174–201, September 1984.
- [13] H. Deconinck and M. Ricchiuto. *Residual Distribution Schemes: Foundations and Analysis*. John Wiley & Sons, Ltd, 2004.
- [14] A. Dutt, L. Greengard, and V. Rokhlin. Spectral Deferred Correction Methods for Ordinary Differential Equations. *BIT Numerical Mathematics*, 40(2):241–266, 2000.
- [15] H. Glaz, P. Colella, I. I. Glass, and L. R. Deschambault. A numerical study of oblique shock-wave reflections with experimental comparisons. 398:117–140, 03 1985.
- [16] T. J. R. Hughes, L. P. Franca, and G. M. Hulbert. A new finite element formulation for computational fluid dynamics: Viii. the galerkin/least-squares method for advective-diffusive equations. *Computer Methods in Applied Mechanics and Engineering*, 73(2):173 – 189, 1989.

- [17] S. Jin and P. Xin. The relaxation schemes for systems of conservation laws in arbitrary space dimensions. *Comm. Pure Appl. Math.*, 48:235–276, 1995.
- [18] M. L. Minion. Semi-implicit spectral deferred correction methods for ordinary differential equations. *Commun. Math. Sci.*, 1(3):471–500, 09 2003.
- [19] L. Pareschi and G. Russo. Implicit–explicit runge–kutta schemes and applications to hyperbolic systems with relaxation. *Journal of Scientific Computing*, 25(1):129–155, Oct 2005.
- [20] M. Ricchiuto and R. Abgrall. Explicit runge-kutta residual distribution schemes for time dependent problems: Second order case. *J. Comput. Phys.*, 229(16):5653–5691, August 2010.
- [21] C. W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock-capturing schemes. *Journal of Computational Physics*, 77(2):439 – 471, 1988.
- [22] D. Torlo and R. Abgrall. Von Neumann analysis for residual distribution galerkin scheme with jump stabilization terms. working paper or preprint, December 2018.